



## UvA-DARE (Digital Academic Repository)

### Adaptive Spectral Galerkin Methods with Dynamic Marking

Canuto, C.; Nochetto, R.H.; Stevenson, R.P.; Verani, M.

**DOI**

[10.1137/15M104579X](https://doi.org/10.1137/15M104579X)

**Publication date**

2016

**Document Version**

Final published version

**Published in**

SIAM journal on numerical analysis

[Link to publication](#)

**Citation for published version (APA):**

Canuto, C., Nochetto, R. H., Stevenson, R. P., & Verani, M. (2016). Adaptive Spectral Galerkin Methods with Dynamic Marking. *SIAM journal on numerical analysis*, *54*(6), 3193–3213. <https://doi.org/10.1137/15M104579X>

**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

**Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

## ADAPTIVE SPECTRAL GALERKIN METHODS WITH DYNAMIC MARKING\*

CLAUDIO CANUTO<sup>†</sup>, RICARDO H. NOCHETTO<sup>‡</sup>, ROB STEVENSON<sup>§</sup>, AND  
MARCO VERANI<sup>¶</sup>

**Abstract.** The convergence and optimality theory of adaptive Galerkin methods is almost exclusively based on the Dörfler marking. This entails a fixed parameter and leads to a contraction constant bounded below away from zero. For spectral Galerkin methods this is a severe limitation which affects performance. We present a dynamic marking strategy that allows for a superlinear relation between consecutive discretization errors, and show exponential convergence with linear computational complexity whenever the solution belongs to a Gevrey approximation class.

**Key words.** spectral methods, adaptivity, convergence, optimal cardinality

**AMS subject classifications.** 65M70, 65N35

**DOI.** 10.1137/15M104579X

**1. Introduction.** The modern analysis of adaptive discretizations of partial differential equations aims at establishing rigorous results of *convergence* and *optimality*. The former results concern the convergence of the approximate solutions produced by the successive iterations of the adaptive algorithm towards the exact solution  $u$ , with an estimate of the error decay rate measured in an appropriate norm. On the other hand, optimality results compare the cardinality of the active set of basis functions used to expand the discrete solution to the minimal cardinality needed to approximate the exact solution with similar accuracy; this endeavor borrows ideas from nonlinear approximation theory. Confining ourselves in the sequel to second-order elliptic boundary value problems, such kind of analysis has been carried out first for wavelet discretizations [11, 14], then for  $h$ -type finite elements [15, 1, 18, 10, 9, 12], [15] dealing just with convergence, and more recently for spectral-type methods [4, 5, 7]; we refer to the surveys [3, 8, 16, 19]. In contrast, the state of the art for  $hp$ -type finite elements [17] is still in evolution; see [13, 2] and the more recent paper [6] which includes optimality estimates.

For all these cases, convergence is proven to be *linear*, i.e., a certain expression controlling the error (a norm, or a combination of norm and estimator) contracts with some fixed parameter  $\rho < 1$  from one iteration to the next one, e.g.,  $\|u - u_{k+1}\| \leq \rho \|u - u_k\|$ . This is typically achieved if the adaptation strategy is based on some

---

\*Received by the editors October 28, 2015; accepted for publication (in revised form) August 26, 2016; published electronically November 3, 2016.

<http://www.siam.org/journals/sinum/54-6/M104579.html>

**Funding:** The first and the fourth authors are partially supported by the Italian research grant *Prin* 2012 2012HBLYE4 “Metodologie innovative nella modellistica differenziale numerica” and by INdAM-GNCS. The second author is partially supported by NSF grants DMS-1109325 and DMS-1411808.

<sup>†</sup>Dipartimento di Scienze Matematiche, Politecnico di Torino, Corso Duca degli Abruzzi 24, I-10129 Torino, Italy (claudio.canuto@polito.it).

<sup>‡</sup>Department of Mathematics and Institute for Physical Science and Technology, University of Maryland, College Park, MD 20742 (rhn@math.umd.edu).

<sup>§</sup>Korteweg-de Vries Institute for Mathematics, University of Amsterdam, P.O. Box 94248, 1090 GE Amsterdam, The Netherlands (r.p.stevenson@uva.nl).

<sup>¶</sup>MOX-Dipartimento di Matematica, Politecnico di Milano, P.zza Leonardo Da Vinci 32, I-20133 Milano, Italy (marco.verani@polimi.it).

form of *Dörfler marking* (or *bulk chasing*) with fixed parameter  $\theta < 1$ : assuming that  $\sum_{i \in \mathcal{J}} \eta_i^2$  is some additive error estimator at iteration  $k$ , one identifies a minimal subset  $\mathcal{J}' \subset \mathcal{J}$  such that

$$\sum_{i \in \mathcal{J}'} \eta_i^2 \geq \theta^2 \sum_{i \in \mathcal{J}} \eta_i^2$$

and utilizes  $\mathcal{J}'$  for the construction of the new discretization at iteration  $k + 1$ . For wavelet or  $h$ -type finite element discretizations, optimality is guaranteed by performing cautious successive adaptations, i.e., by choosing a moderate value of  $\theta$ , say  $0 < \theta \leq \theta_{\max} < 1$  [18]. This avoids the need for cleaning up the discrete solution from time to time, by subjecting it to a *coarsening* stage.

On the other hand, the resulting contraction factor  $\rho = \rho(\theta)$  turns out to be bounded from below by a positive constant, say  $0 < \rho_{\min} \leq \rho < 1$  (related to the “condition number” of the exact problem), regardless of the choice of  $\theta$ . This entails a limitation on the speed of convergence for infinite-order methods [4, 5], but is not restrictive for fixed-order methods [18, 10].

It has been shown in [4] that such an obstruction can be avoided if a specific property of the differential operator holds, namely, the so-called *quasi-sparsity* of the inverse of the associated stiffness matrix. Upon exploiting this information, a more aggressive marking strategy can be adopted, which judiciously enlarges the set  $\mathcal{J}'$  coming out of Dörfler’s stage. The resulting contraction factor  $\rho$  can be now made arbitrarily close to 0 by choosing  $\theta$  arbitrarily close to 1.

When a method of spectral type is used, one expects a fast (possibly, exponentially fast) decay of the discretization error for smooth solutions. In such a situation, a slow convergence of the iterations of the adaptive algorithm would endanger the overall performance of the method; from this perspective, it is useful to be able to make the contraction factor as close to 0 as desired. Yet, linear convergence of the adaptive iterations is not enough to guarantee the optimality of the method. Let us explain why this occurs, and why a *superlinear* convergence is preferable, using the following idealized setting.

As is customary in nonlinear approximation, we consider the *best  $N$ -term approximation error*  $E_N(u)$  of the exact solution  $u$ , in a suitable norm, using combinations of at most  $N$  functions taken from a chosen basis. We prescribe a decay law of  $E_N(u)$  as  $N$  increases, which classically for fixed-order approximations is *algebraic* and reads

$$(1.1) \quad \sup_N N^s E_N(u) < \infty$$

for some positive  $s$ . However, for infinite-order methods such as spectral approximations an *exponential* law is relevant that reads

$$(1.2) \quad \sup_N e^{\eta N^\alpha} E_N(u) < \infty$$

for some  $\eta > 0$  and  $\alpha \in (0, 1]$ , where  $\alpha < 1$  accommodates the inclusion of  $C^\infty$ -functions that are not analytic. This defines corresponding algebraic and exponential *sparsity classes* for the exact solution  $u$ . These classes are related to Besov and Gevrey regularity of  $u$  respectively.

We now assume the ideal situation that at each iteration of our adaptive algorithm<sup>1</sup>

$$(1.3) \quad \|u - u_k\| \approx N_k^{-s} \quad \text{or} \quad \|u - u_k\| \approx e^{-\eta N_k^\alpha},$$

<sup>1</sup>Throughout the paper, we write  $A_k \lesssim B_k$  to indicate that  $A_k$  can be bounded by a multiple of  $B_k$ , independently of the iteration counter  $k$  and other parameters which  $A_k$  and  $B_k$  may depend on;  $A_k \approx B_k$  means  $A_k \lesssim B_k$  and  $B_k \lesssim A_k$ .

where  $N_k$  is the cardinality of the discrete solution  $u_k$ , i.e., the dimension of the approximation space activated at iteration  $k$ . We assume, in addition, that the error decays linearly from one iteration to the next, i.e., it satisfies precisely

$$(1.4) \quad \|u - u_{k+1}\| = \rho \|u - u_k\|.$$

If  $u$  belongs to a sparsity class of algebraic type, then one easily gets  $N_k \approx \rho^{-k/s}$ , i.e., cardinalities grow exponentially fast and

$$\Delta N_k := N_{k+1} - N_k \approx N_k \approx \|u - u_k\|^{-1/s},$$

i.e., the increment of cardinality between consecutive iterations is proportional to the current cardinality as well as to the error raised to the power  $-1/s$ . The important message stemming from this ideal setting is that for a *practical* adaptive algorithm one should be able to derive the estimates  $\|u - u_{k+1}\| \leq \rho \|u - u_k\|$  and  $\Delta N_k \lesssim \|u - u_k\|^{-1/s}$ , because they yield

$$N_n = \sum_{k=0}^{n-1} \Delta N_k \lesssim \sum_{k=0}^{n-1} \|u - u_k\|^{-1/s} \leq \|u - u_n\|^{-1/s} \sum_{k=0}^{n-1} \rho^{(n-k)/s} \lesssim \|u - u_n\|^{-1/s}.$$

This geometric-series argument is precisely the strategy used in [18, 10] and gives an estimate similar to (1.1). The performance of a practical adaptive algorithm is thus quasi-optimal.

If  $u$  belongs to a sparsity class of exponential type, instead, the situation changes radically. In fact, assuming (1.3) and (1.4), one has  $e^{-\eta N_k^\alpha} \approx \rho^k$ , and so

$$\lim_{k \rightarrow \infty} k^{-1/\alpha} N_k = (|\log \rho|/\eta)^{1/\alpha},$$

i.e., the cardinality  $N_k$  grows polynomially. For a practical adaptive algorithm, proving such a growth is very hard if not impossible. This obstruction has motivated the insertion of a coarsening stage in the adaptive algorithm presented in [4]. Coarsening removes the negligible components of the discrete solution possibly activated by the marking strategy and guarantees that the final cardinality is nearly optimal [11, 4], but it does not account for the workload to create  $u_k$ .

One of the key points of the present contribution is the observation that if the convergence of the adaptive algorithm is superlinear, then one is back to the simpler case of exponential growth of cardinalities which is amenable to a sharper performance analysis. To see this, let us assume a superlinear relation between consecutive errors:

$$(1.5) \quad \|u - u_{k+1}\| = \|u - u_k\|^q$$

for some  $q > 1$ . If additionally  $u_k$  satisfies (1.3), then one infers that  $e^{-\eta N_{k+1}^\alpha} \approx e^{-\eta q N_k^\alpha}$ , whence

$$\lim_{k \rightarrow \infty} \frac{\Delta N_k}{N_k} = q^{1/\alpha} - 1, \quad \lim_{k \rightarrow \infty} \frac{|\log \|u - u_k\||^{1/\alpha}}{N_k} = \eta^{1/\alpha},$$

the latter being just a consequence of (1.3). This suggests that the geometric-series argument may be invoked again in the optimality analysis of the adaptive algorithm.

This ideal setting does not apply directly to our *practical* adaptive algorithm. We will be able to prove estimates that are consistent with the preceding derivation to some extent, namely,

$$\|u - u_{k+1}\| \leq \|u - u_k\|^q, \quad \Delta N_k \leq Q |\log \|u - u_k\||^{1/\bar{\alpha}}$$

with constants  $Q > 0$  and  $\bar{\alpha} \in (0, \alpha]$ . Invoking  $\|u - u_n\| \leq \|u - u_k\|^{q^{n-k}}$ , we then realize that

$$N_n = \sum_{k=0}^{n-1} \Delta N_k \leq Q \sum_{k=0}^{n-1} |\log \|u - u_k\||^{1/\bar{\alpha}} \leq \frac{Qq^{1/\bar{\alpha}}}{q^{1/\bar{\alpha}} - 1} |\log \|u - u_n\||^{1/\bar{\alpha}}.$$

Setting  $\bar{\eta} := \left(\frac{Qq^{1/\bar{\alpha}}}{q^{1/\bar{\alpha}} - 1}\right)^{-\bar{\alpha}}$ , we deduce the estimate

$$\sup_n e^{\bar{\eta} N_n^{\bar{\alpha}}} \|u - u_n\| \leq 1,$$

which is similar to (1.2), albeit with different class parameters. The most important parameter is  $\bar{\alpha}$ . Its possible degradation relative to  $\alpha$  is mainly caused by the fact that the residual, the only computable quantity accessible to our practical algorithm, belongs to a sparsity class with a main parameter generally smaller than that of the solution  $u$ . This perhaps unexpected property is typical of the exponential class and has been elucidated in [4].

In order for the marking strategy to guarantee superlinear convergence, one needs to adopt a dynamic choice of Dörfler's parameter  $\theta$ , which pushes its value towards 1 as the iterations proceed. We accomplish this requirement by equating the quantity  $1 - \theta_k^2$  to some function of the dual norm of the residual  $r_k$ , which is monotonically increasing and vanishing at the origin. This defines our *dynamic marking strategy*. The order of the root at the origin dictates the exponent  $q$  in the superlinear convergence estimate of our adaptive algorithm.

The paper is organized as follows. In section 2 we introduce the model elliptic problem and its spectral Galerkin approximation based on either multidimensional Fourier or (modified) Legendre expansions. In particular, we highlight properties of the resulting stiffness matrix that will be fundamental in what follows. We present the adaptive algorithm in section 3, first for the static marking ( $\theta$  fixed) and later for the dynamic marking ( $\theta$  tending towards 1); superlinear convergence is proven. With the optimality analysis in mind, we next recall in section 4 the definition and crucial properties of a family of sparsity classes of exponential type, related to Gevrey regularity of the solution, and we investigate how the sparsity class of the Galerkin residual deteriorates relative to that of the exact solution. Finally, in section 5 we relate the cardinality of the adaptive discrete solutions, as well as the workload needed to compute them, to the expected accuracy of the approximation. Our analysis confirms that the proposed dynamic marking strategy avoids any form of coarsening, while providing exponential convergence with linear computational complexity, assuming optimal linear solvers.

**2. Model elliptic problem and Galerkin methods.** Let  $d \geq 1$  and consider the following elliptic PDE in a  $d$ -dimensional rectangular domain  $\Omega$  with periodic or homogeneous Dirichlet boundary conditions:

$$(2.1) \quad Lu = -\nabla \cdot (\nu \nabla u) + \sigma u = f \quad \text{in } \Omega,$$

where  $\nu$  and  $\sigma$  are sufficiently smooth real coefficients satisfying  $0 < \nu_* \leq \nu(x) \leq \nu^* < \infty$  and  $0 < \sigma_* \leq \sigma(x) \leq \sigma^* < \infty$  in  $\Omega$ ; let us set

$$\alpha_* = \min(\nu_*, \sigma_*) \quad \text{and} \quad \alpha^* = \max(\nu^*, \sigma^*) .$$

Depending on the boundary conditions, let  $V$  be equal to  $H_0^1(\Omega)$  or  $H_p^1(\Omega)$ , being the space of periodic functions with square integrable weak gradient, and denote by  $V^*$  its dual space. We formulate (2.1) variationally as

$$(2.2) \quad u \in V \quad : \quad a(u, v) = \langle f, v \rangle \quad \forall v \in V ,$$

where  $a(u, v) = \int_{\Omega} \nu \nabla u \cdot \nabla \bar{v} + \int_{\Omega} \sigma u \bar{v}$  (bar indicating as usual complex conjugate). We denote by  $\|v\| = \sqrt{a(v, v)}$  the energy norm of any  $v \in V$ , which satisfies

$$(2.3) \quad \sqrt{\alpha_*} \|v\|_V \leq \|v\| \leq \sqrt{\alpha^*} \|v\|_V .$$

**2.1. Riesz basis.** We start with an abstract formulation which encompasses the two examples of interest: trigonometric functions and Legendre polynomials. Let  $\phi = \{\phi_k : k \in \mathcal{K}\}$  be a Riesz basis of  $V$ . Thus, we assume the following relation between a function  $v = \sum_{k \in \mathcal{K}} \hat{v}_k \phi_k \in V$  and its coefficients:

$$(2.4) \quad \|v\|_V^2 \simeq \sum_{k \in \mathcal{K}} |\hat{v}_k|^2 d_k =: \|v\|_{\phi}^2$$

for suitable weights  $d_k > 0$ . Correspondingly, any element  $f \in V^*$  can be expanded along the *dual basis*  $\phi^* = \{\phi_k^*\}$  as  $f = \sum_{k \in \mathcal{K}} \hat{f}_k \phi_k^*$ , with  $\hat{f}_k = \langle f, \phi_k \rangle$ , yielding the dual norm representation

$$(2.5) \quad \|f\|_{V^*}^2 \simeq \sum_{k \in \mathcal{K}} |\hat{f}_k|^2 d_k^{-1} =: \|v\|_{\phi^*}^2 .$$

For future reference, we introduce the vectors  $\mathbf{v} = (\hat{v}_k d_k^{1/2})_{k \in \mathcal{K}}$  and  $\mathbf{f} = (\hat{f}_k d_k^{-1/2})_{k \in \mathcal{K}}$  as well as the constants  $\beta_* \leq \beta^*$  of the norm equivalence in (2.4)

$$(2.6) \quad \beta_* \|v\|_V \leq \|v\|_{\phi} = \|\mathbf{v}\|_{\ell^2} \leq \beta^* \|v\|_V \quad \forall v \in V .$$

This implies

$$(2.7) \quad \frac{1}{\beta^*} \|f\|_{V^*} \leq \|f\|_{\phi^*} = \|\mathbf{f}\|_{\ell^2} \leq \frac{1}{\beta_*} \|f\|_{V^*} \quad \forall f \in V^* .$$

The two key examples to keep in mind are the trigonometric basis and tensor products of the Babuška–Shen (BS) basis. We discuss them briefly below.

**Trigonometric basis.** Let  $\Omega = (0, 2\pi)^d$  and the trigonometric basis be  $\phi_k(x) = \frac{1}{(2\pi)^{d/2}} e^{ik \cdot x}$  for any  $k \in \mathcal{K} = \mathbb{Z}^d$  and  $x \in \Omega$ . Any function  $v \in L^2(\Omega)$  can be expanded in terms of  $\{\phi_k\}_{k \in \mathbb{Z}^d}$  as follows:

$$(2.8) \quad v = \sum_k \hat{v}_k \phi_k , \quad \hat{v}_k = \langle v, \phi_k \rangle , \quad \|v\|_{L^2(\Omega)}^2 = \sum_k |\hat{v}_k|^2 .$$

The space  $V := H_p^1(\Omega)$  can now be easily characterized as the subspace of those  $v \in L^2(\Omega)$  for which

$$\|v\|_V^2 = \|v\|_{H_p^1(\Omega)}^2 = \sum_k |\hat{V}_k|^2 < \infty \quad (\text{where } \hat{V}_k := \hat{v}_k d_k^{1/2} \text{ with } d_k := 1 + |k|^2).$$

This induces an *isomorphism* between  $H_p^1(\Omega)$  and  $\ell^2(\mathbb{Z}^d)$ : for each  $v \in H_p^1(\Omega)$  let  $\mathbf{v} = (\hat{V}_k)_{k \in \mathcal{K}} \in \ell^2(\mathbb{Z}^d)$  and note that  $\|v\|_{H_p^1(\Omega)} = \|\mathbf{v}\|_{\ell^2}$ . Likewise, the dual space  $H_p^{-1}(\Omega) = (H_p^1(\Omega))'$  is characterized as the space of those functionals  $f$  for which

$$\|f\|_{V^*}^2 = \|f\|_{H_p^{-1}(\Omega)}^2 = \sum_k |\hat{F}_k|^2 \quad \text{with} \quad \hat{F}_k := \hat{f}_k d_k^{-1/2}.$$

We also have an isomorphism between  $H_p^{-1}(\Omega)$  and  $\ell^2(\mathbb{Z}^d)$  upon setting  $\mathbf{f} = (\hat{F}_k)_{k \in \mathcal{K}}$  for  $f \in H_p^{-1}(\Omega)$  and realizing that  $\|f\|_{H_p^{-1}(\Omega)} = \|\mathbf{f}\|_{\ell^2}$ .

**Babuška–Shen (BS) basis.** Let us start with the one-dimensional case  $d = 1$ . Set  $I = (-1, 1)$ ,  $V := H_0^1(I)$ , and let  $L_k(x)$ ,  $k \geq 0$ , stand for the  $k$ th Legendre orthogonal polynomial in  $I$ , which satisfies  $\deg L_k = k$ ,  $L_k(1) = 1$ , and

$$(2.9) \quad \int_I L_k(x)L_m(x) dx = \frac{2}{2k+1} \delta_{km}, \quad m \geq 0.$$

The natural modal basis in  $H_0^1(I)$  is the *BS basis*, whose elements are defined as

$$(2.10) \quad \eta_k(x) = \sqrt{\frac{2k-1}{2}} \int_x^1 L_{k-1}(s) ds = \frac{1}{\sqrt{4k-2}} (L_{k-2}(x) - L_k(x)), \quad k \geq 2.$$

The basis elements satisfy  $\deg \eta_k = k$  and

$$(2.11) \quad (\eta_k, \eta_m)_{H_0^1(I)} = \int_I \eta'_k(x)\eta'_m(x) dx = \delta_{km}, \quad k, m \geq 2,$$

i.e., they form an orthonormal basis for the  $H_0^1(I)$  inner product. Equivalently, the (semi-infinite) stiffness matrix  $S_\eta$  of the BS basis with respect to this inner product is the identity matrix.

We now consider, for simplicity, the two-dimensional case  $d = 2$  since the case  $d > 2$  is similar. Let  $\Omega = (-1, 1)^2$ ,  $V = H_0^1(\Omega)$ , and consider the *tensorized BS basis*, whose elements are defined as

$$(2.12) \quad \eta_k(x) = \eta_{k_1}(x_1)\eta_{k_2}(x_2), \quad k_1, k_2 \geq 2,$$

where we set  $k = (k_1, k_2)$  and  $x = (x_1, x_2)$ ; indices vary in the set

$$\mathcal{K} = \{k \in \mathbb{N}^2 : k_i \geq 2 \text{ for } i = 1, 2\},$$

which is ordered “a la Cantor” by increasing total degree  $k_{\text{tot}} = k_1 + k_2$  and, for the same total degree, by increasing  $k_1$ . The tensorized BS basis is no longer orthogonal, since

$$(2.13) \quad (\eta_k, \eta_m)_{H_0^1(\Omega)} = (\eta_{k_1}, \eta_{m_1})_{H_0^1(I)} (\eta_{k_2}, \eta_{m_2})_{L^2(I)} + (\eta_{k_1}, \eta_{m_1})_{L^2(I)} (\eta_{k_2}, \eta_{m_2})_{H_0^1(I)},$$

whence  $(\eta_k, \eta_m)_{H_0^1(\Omega)} \neq 0$  if and only if  $k_1 = m_1$  and  $k_2 - m_2 \in \{-2, 0, 2\}$ , or  $k_2 = m_2$  and  $k_1 - m_1 \in \{-2, 0, 2\}$ . Obviously, we cannot have a Parseval representation of the  $H_0^1(\Omega)$ -norm of  $v = \sum_{k \in \mathcal{K}} \hat{v}_k \eta_k$  in terms of the coefficients  $\hat{v}_k$ . With the aim of getting (2.4), we follow [7] and we first perform the orthonormalization of the BS basis via a Gram–Schmidt procedure. This allows us to build a sequence of functions

$$(2.14) \quad \Phi_k = \sum_{m \leq k} g_{mk} \eta_m,$$

such that  $g_{kk} \neq 0$  and

$$(\Phi_k, \Phi_m)_{H_0^1(\Omega)} = \delta_{km} \quad \forall k, m \in \mathcal{K}.$$

We will refer to the collection  $\Phi := \{\Phi_k : k \in \mathcal{K}\}$  as the *orthonormal Babuška–Shen* (OBS) basis, for which the associated stiffness matrix  $S_\Phi$  with respect to the  $H_0^1(\Omega)$ -inner product is the identity matrix. Equivalently, if  $G = (g_{mk})$  is the upper triangular matrix which collects the coefficients generated by the Gram–Schmidt algorithm above, one has

$$(2.15) \quad G^T S_\eta G = S_\Phi = I,$$

that is the validity of (2.6) with  $d_k = 1$ . However, unlike  $S_\eta$ , which is very sparse, the upper triangular matrix  $G$  is full; in view of this, we next apply a thresholding procedure to wipe out a significant portion of the nonzero entries sitting in the left-most columns of  $G$ . This leads to a modified basis whose computational efficiency is quantitatively improved, without significantly deteriorating the properties of the OBS basis. To be more precise, we use the following notation:  $G_t$  indicates the matrix obtained from  $G$  by setting to zero a certain finite set of off-diagonal entries, so that in particular  $\text{diag}(G_t) = \text{diag}(G)$ ; correspondingly,  $E := G_t - G$  is the matrix measuring the truncation quality, for which  $\text{diag}(E) = 0$ . Finally, we introduce the matrix

$$(2.16) \quad S_\phi = G_t^T S_\eta G_t$$

which we interpret as the stiffness matrix associated with the modified BS basis defined in analogy to (2.14) as

$$(2.17) \quad \phi_k = \sum_{m \in \mathcal{M}_t(k)} g_{mk} \eta_m,$$

where  $\mathcal{M}_t(k) = \{m \leq k : E_{mk} = 0\}$ . This forms a new basis in  $H_0^1(\Omega)$  (because  $k \in \mathcal{M}_t(k)$  and  $g_{kk} \neq 0$ ). We will term it a *nearly orthonormal Babuška–Shen* (NOBS) basis. Note that only the basis functions  $\phi_k$  having total degree not exceeding a certain value, say  $p$ , may be affected by the compression, while all the others coincide with the corresponding orthonormal basis functions  $\Phi_k$  defined in (2.14).

If  $D_\phi = \text{diag } S_\phi$ , then for any value of  $p$  there are strategies to build  $G_t$  (depending on  $p$ ) such that the eigenvalues  $\lambda$  of

$$(2.18) \quad S_\phi x = \lambda D_\phi x$$

are close to one and bounded from above and away from 0, independently of  $p$  [7]. This guarantees the validity for the NOBS basis of (2.4) with  $d_k$  equal to the diagonal elements of the matrix  $D_\phi$  and (2.6) for suitable choice of constants  $\beta_*$ ,  $\beta^*$  depending on the eigenvalues of (2.18) (see [7] for more details).

**2.2. Infinite dimensional algebraic problem.** Let us identify the solution  $u = \sum_k \hat{u}_k \phi_k$  of problem (2.1) with the vector  $\mathbf{u} = (\hat{U}_k)_{k \in \mathcal{K}} = (\hat{u}_k d_k^{1/2})_{k \in \mathcal{K}}$  of its normalized coefficients with respect to the basis  $\{\phi_k\}_{k \in \mathcal{K}}$ . Similarly, let us identify the right-hand side  $f$  with the vector  $\mathbf{f} = (\hat{F}_\ell)_{\ell \in \mathcal{K}} = (\hat{f}_\ell d_\ell^{-1/2})_{\ell \in \mathcal{K}}$  of its normalized dual coefficients. Finally, let us introduce the bi-infinite, symmetric, and positive-definite stiffness matrix

$$(2.19) \quad \mathbf{A} = (a_{\ell,k})_{\ell,k \in \mathcal{K}} \quad \text{with} \quad a_{\ell,k} = \frac{a(\phi_k, \phi_\ell)}{\sqrt{d_k} \sqrt{d_\ell}}.$$



Then, problem (2.1) can be equivalently written as

$$(2.20) \quad \mathbf{A}\mathbf{u} = \mathbf{f} ,$$

where, thanks to (2.3) and the norm equivalences (2.6)–(2.7),  $\mathbf{A}$  defines a bounded invertible operator in  $\ell^2(\mathcal{K})$ .

**Decay properties of  $\mathbf{A}$  and  $\mathbf{A}^{-1}$ .** The decay of the entries of  $\mathbf{A}$  away from the diagonal depends on the regularity of the coefficients  $\nu$  and  $\sigma$  of  $L$ . If  $\nu$  and  $\sigma$  are real analytic in a neighborhood of  $\Omega$ , then  $a_{k,m}$  decays exponentially away from the diagonal [4, 5, 7]: there exist parameters  $c_L, \eta_L > 0$  such that

$$(2.21) \quad |a_{k,m}| \leq c_L \exp(-\eta_L |k - m|) \quad \forall k, m \in \mathcal{K};$$

we then say that  $\mathbf{A}$  belongs to the exponential class  $\mathcal{D}_e(\eta_L)$ , in particular  $\mathbf{A}$  is quasi-sparse. This justifies the *symmetric truncation*  $\mathbf{A}_J$  of  $\mathbf{A}$  with *parameter*  $J$ , defined as  $(\mathbf{A}_J)_{\ell,k} = a_{\ell,k}$  if  $|\ell - k| \leq J$  and  $(\mathbf{A}_J)_{\ell,k} = 0$  otherwise, which satisfies [4, 5, 7]

$$(2.22) \quad \|\mathbf{A} - \mathbf{A}_J\|_{\ell^2 \rightarrow \ell^2} \leq C_{\mathbf{A}}(J+1)^{d-1} e^{-\eta_L J}$$

for some  $C_{\mathbf{A}} > 0$  depending only on  $c_L$ .

Most notably, the inverse matrix  $\mathbf{A}^{-1}$  is also quasi-sparse [4, 5, 7]. Precisely,  $\mathbf{A}^{-1} \in \mathcal{D}_e(\bar{\eta}_L)$  for some  $\bar{\eta}_L \in (0, \eta_L]$  and  $\bar{c}_L$  only dependent on  $c_L$  and  $\eta_L$ . Thus, there exists an explicit constant  $C_{\mathbf{A}^{-1}}$  (depending only on  $c_L$  and  $\eta_L$ ) such that the symmetric truncation  $(\mathbf{A}^{-1})_J$  of  $\mathbf{A}^{-1}$  satisfies

$$(2.23) \quad \|\mathbf{A}^{-1} - (\mathbf{A}^{-1})_J\|_{\ell^2 \rightarrow \ell^2} \leq C_{\mathbf{A}^{-1}}(J+1)^{d-1} e^{-\bar{\eta}_L J} \leq C_{\mathbf{A}^{-1}} e^{-\bar{\eta}_L J}$$

for a suitable exponent  $\tilde{\eta}_L < \bar{\eta}_L$ .

**Galerkin method.** Given any finite index set  $\Lambda \subset \mathcal{K}$ , we define the subspace  $V_{\Lambda} = \text{span}\{\phi_k \mid k \in \Lambda\}$  of  $V$ ; we set  $|\Lambda| = \text{card } \Lambda$ , so that  $\dim V_{\Lambda} = |\Lambda|$ . If  $v \in V$  admits the expansion  $v = \sum_{k \in \mathcal{K}} \hat{v}_k \phi_k$ , then we define its projection  $P_{\Lambda} v$  upon  $V_{\Lambda}$  by setting  $P_{\Lambda} v := \sum_{k \in \Lambda} \hat{v}_k \phi_k$ . Similarly, we define the subspace  $V_{\Lambda}^* = \text{span}\{\phi_k^* \mid k \in \Lambda\}$  of  $V^*$ . If  $f$  admits an expansion  $f = \sum_{k \in \mathcal{K}} \hat{f}_k \phi_k^*$ , then we define its projection  $P_{\Lambda}^* f$  onto  $V_{\Lambda}^*$  upon setting  $P_{\Lambda}^* f := \sum_{k \in \Lambda} \hat{f}_k \phi_k^*$ .

Given any finite  $\Lambda \subset \mathcal{K}$ , the Galerkin approximation of (2.1) is defined as

$$(2.24) \quad u_{\Lambda} \in V_{\Lambda} \quad : \quad a(u_{\Lambda}, v_{\Lambda}) = \langle f, v_{\Lambda} \rangle \quad \forall v_{\Lambda} \in V_{\Lambda} .$$

Let  $\mathbf{u}_{\Lambda}$  be the vector collecting the coefficients of  $u_{\Lambda}$  indexed in  $\Lambda$ ; let  $\mathbf{f}_{\Lambda}$  be the analogous restriction for the vector of the coefficients of  $f$ . Finally, denote by  $\mathbf{R}_{\Lambda}$  the matrix that restricts a vector indexed in  $\mathcal{K}$  to the portion indexed in  $\Lambda$ , so that  $\mathbf{R}_{\Lambda}^H$  is the corresponding extension matrix. If

$$(2.25) \quad \mathbf{A}_{\Lambda} := \mathbf{R}_{\Lambda} \mathbf{A} \mathbf{R}_{\Lambda}^H ,$$

then problem (2.24) can be equivalently written as

$$(2.26) \quad \mathbf{A}_{\Lambda} \mathbf{u}_{\Lambda} = \mathbf{f}_{\Lambda} .$$

For any  $w \in V_{\Lambda}$ , we define the residual  $r(w) \in V^*$  as

$$r(w) = f - Lw = \sum_{k \in \mathcal{K}} \hat{r}_k(w) \phi_k^*$$

where  $\hat{r}_k(w) = \langle f - Lw, \phi_k \rangle = \langle f, \phi_k \rangle - a(w, \phi_k)$ . The definition (2.24) of  $u_{\Lambda}$  is

equivalent to the condition  $P_\Lambda^* r(u_\Lambda) = 0$ , i.e.,  $\hat{r}_k(u_\Lambda) = 0$  for every  $k \in \Lambda$ . By the continuity and coercivity of the bilinear form, one has

$$(2.27) \quad \frac{1}{\alpha^*} \|r(u_\Lambda)\|_{V^*} \leq \|u - u_\Lambda\|_V \leq \frac{1}{\alpha_*} \|r(u_\Lambda)\|_{V^*} ,$$

which in view of (2.3) and (2.7) can be rephrased as

$$(2.28) \quad \frac{\beta_*}{\sqrt{\alpha^*}} \|r(u_\Lambda)\|_{\phi^*} \leq \|u - u_\Lambda\| \leq \frac{\beta^*}{\sqrt{\alpha_*}} \|r(u_\Lambda)\|_{\phi^*} .$$

Therefore, if  $(\hat{r}_k(u_\Lambda))_{k \in \mathcal{K}}$  are the coefficients of  $r(u_\Lambda)$  with respect to the dual basis  $\phi^*$ , the quantity

$$\|r(u_\Lambda)\|_{\phi^*} = \left( \sum_{k \notin \Lambda} |\hat{R}_k(u_\Lambda)|^2 \right)^{1/2} \quad \text{with} \quad \hat{R}_k(u_\Lambda) = \hat{r}_k(u_\Lambda) d_k^{-1/2}$$

is an *error estimator* from above and from below. However, this quantity is not computable because it involves infinitely many terms. A *feasible version* of the present algorithm can be efficiently realized by combining the ideas of [11] (see also [19]) and the results contained in [4, 5, 7].

**Equivalent formulation of the Galerkin problem.** For future reference, we now rewrite the Galerkin problem (2.24) in an equivalent (infinite-dimensional) manner. Let

$$\mathbf{P}_\Lambda : \ell^2(\mathcal{K}) \rightarrow \ell^2(\mathcal{K})$$

be the projector defined as

$$(\mathbf{P}_\Lambda \mathbf{v})_\lambda := \begin{cases} v_\lambda & \text{if } \lambda \in \Lambda , \\ 0 & \text{if } \lambda \notin \Lambda . \end{cases}$$

Note that  $\mathbf{P}_\Lambda$  can be represented as a diagonal bi-infinite matrix whose diagonal elements are 1 for indexes belonging to  $\Lambda$ , and zero otherwise. We set  $\mathbf{Q}_\Lambda := \mathbf{I} - \mathbf{P}_\Lambda$  and introduce the bi-infinite matrix  $\hat{\mathbf{A}}_\Lambda := \mathbf{P}_\Lambda \mathbf{A} \mathbf{P}_\Lambda + \mathbf{Q}_\Lambda$  which is equal to  $\mathbf{A}_\Lambda$  for indexes in  $\Lambda$  and to the identity matrix, otherwise. The definitions of the projectors  $\mathbf{P}_\Lambda$  and  $\mathbf{Q}_\Lambda$  yield the following property:

$$(2.29) \quad \text{If } \mathbf{A} \text{ is invertible with } \mathbf{A} \in \mathcal{D}_e(\eta_L), \text{ then the same holds for } \hat{\mathbf{A}}_\Lambda .$$

Furthermore, the constants  $C_{\hat{\mathbf{A}}_\Lambda}$  and  $C_{(\hat{\mathbf{A}}_\Lambda)^{-1}}$  which appear in the inequalities (2.22) and (2.23) for  $\hat{\mathbf{A}}_\Lambda$  can be bounded uniformly in  $\Lambda$ , since in turn they can be bounded in terms of  $\eta_L$  and  $c_L$ , respectively.

Now, let us consider the following extended Galerkin problem: find  $\hat{\mathbf{u}} \in \ell^2$  such that

$$(2.30) \quad \hat{\mathbf{A}}_\Lambda \hat{\mathbf{u}} = \mathbf{P}_\Lambda \mathbf{f} .$$

Let  $\mathbf{u}_\Lambda$  be the Galerkin solution to (2.26); then, it is easy to check that  $\hat{\mathbf{u}} = \mathbf{R}_\Lambda^H \mathbf{u}_\Lambda$ .

**3. Adaptive spectral Galerkin method.** In this section we present our adaptive spectral Galerkin method, named **DYN-GAL**, that is based on a new notion of marking strategy, namely, a *dynamic marking*. In section 3.1 we recall the enriched Dörfler marking strategy, introduced in [4], which represents an enhancement of the classic Dörfler marking strategy. In section 3.2 we introduce the dynamic marking strategy, present **DYN-GAL**, and prove its quadratic convergence.

**3.1. Static Dörfler marking.** Fix any  $\theta \in (0, 1)$  and set  $\Lambda_0 = \emptyset, u_{\Lambda_0} = 0$ . For  $n = 0, 1, \dots$ , assume that  $\Lambda_n$  and  $u_n := u_{\Lambda_n} \in V_{\Lambda_n}$  and  $r_n := r(u_n) = Lu_n - f$  are already computed and choose  $\Lambda_{n+1} := \Lambda_n \cup \partial\Lambda_n$ , where the set  $\partial\Lambda_n$  is built by a two-step procedure that we call **E-DÖRFLER** for *enriched Dörfler*:

**Function**  $\partial\Lambda_n = \mathbf{E-DÖRFLER}(\Lambda_n, \theta)$   
 $\widetilde{\partial\Lambda}_n = \mathbf{DÖRFLER}(r_n, \theta)$   
 $\partial\Lambda_n = \mathbf{ENRICH}(\widetilde{\partial\Lambda}_n, J)$

The first step is the usual *Dörfler's marking* with parameter  $\theta$ :

$$(3.1) \quad \|P_{\widetilde{\partial\Lambda}_n}^* r_n\|_{\phi^*} = \|P_{\Lambda_{n+1}}^* r_n\|_{\phi^*} \geq \theta \|r_n\|_{\phi^*} \quad \text{or} \quad \sum_{k \in \widetilde{\partial\Lambda}_n} |\hat{R}_k(u_n)|^2 \geq \theta^2 \sum_{k \in \mathcal{K}} |\hat{R}_k(u_n)|^2$$

with  $\widetilde{\Lambda}_{n+1} = \Lambda_n \cup \widetilde{\partial\Lambda}_n$ . This also reads

$$(3.2) \quad \|r_n - P_{\Lambda_{n+1}}^* r_n\|_{\phi^*} \leq \sqrt{1 - \theta^2} \|r_n\|_{\phi^*}$$

and can be implemented by rearranging the coefficients  $\hat{R}_k(u_n)$  in decreasing order of modulus and picking the largest ones (*greedy approach*). However, this is only an idealized algorithm because the number of coefficients  $\hat{R}_k(u_n)$  is infinite. This marking is known to yield a contraction property between  $u_n$  and the Galerkin solution  $\tilde{u}_{n+1} \in V_{\widetilde{\Lambda}_{n+1}}$  of the form

$$\|u - \tilde{u}_{n+1}\| \leq \rho(\theta) \|u - u_n\|$$

with  $\rho(\theta) = \sqrt{1 - (\frac{\beta_*}{\beta^*})^2 \frac{\alpha_*}{\alpha^*} \theta^2}$  [4, 5]. When  $\alpha_* < \alpha^*$  we see that in contrast to (3.2),  $\rho(\theta)$  is bounded below away from 0 by  $\sqrt{1 - (\frac{\beta_*}{\beta^*})^2 \frac{\alpha_*}{\alpha^*}}$ .

The second step of **E-DÖRFLER** is meant to remedy this situation and hinges on the a priori structure of  $\mathbf{A}^{-1}$  already alluded to in section 2.2. The goal is to augment the set  $\widetilde{\partial\Lambda}_n$  to  $\partial\Lambda_n$  judiciously. This is contained in the following proposition (see [4]) whose proof is reported here for completeness.

**PROPOSITION 3.1** (enrichment). *Let  $\widetilde{\partial\Lambda}_n = \mathbf{DÖRFLER}(r_n, \theta)$ , and let  $J = J(\theta) > 0$  satisfy*

$$(3.3) \quad C_{\mathbf{A}^{-1}} e^{-\tilde{\eta}_L J} \leq \frac{\beta_* \beta^*}{\sqrt{\alpha_* \alpha^*}} \sqrt{1 - \theta^2},$$

where  $C_{\mathbf{A}^{-1}}$  and  $\tilde{\eta}_L$  are defined in (2.23). Let  $\partial\Lambda_n = \mathbf{ENRICH}(\widetilde{\partial\Lambda}_n, J)$  be built as follows:

$$\partial\Lambda_n := \{k \in \mathcal{K} : \text{there exists } \ell \in \widetilde{\partial\Lambda}_n \text{ such that } |k - \ell| \leq J\}.$$

Then for  $\Lambda_{n+1} = \Lambda_n \cup \partial\Lambda_n$ , the Galerkin solution  $u_{n+1} \in V_{\Lambda_{n+1}}$  satisfies

$$(3.4) \quad \|u - u_{n+1}\| \leq \bar{\rho}(\theta) \|u - u_n\|$$

with

$$(3.5) \quad \bar{\rho}(\theta) = 2 \frac{\beta^* \sqrt{\alpha^*}}{\beta_* \sqrt{\alpha_*}} \sqrt{1 - \theta^2}.$$

*Proof.* Let  $g_n := P_{\partial\tilde{\Lambda}_n}^* r_n = P_{\tilde{\Lambda}_{n+1}}^* r_n$  which, according to (3.2), satisfies

$$\|r_n - g_n\|_{\phi^*} \leq \sqrt{1 - \theta^2} \|r_n\|_{\phi^*} .$$

Let  $w_n \in V$  be the solution of  $Lw_n = g_n$ , which, in general, will have infinitely many components, and let us split it as

$$w_n = P_{\Lambda_{n+1}} w_n + P_{\Lambda_{n+1}^c} w_n =: y_n + z_n \in V_{\Lambda_{n+1}} \oplus V_{\Lambda_{n+1}^c} .$$

The minimality property in the energy norm of the Galerkin solution  $u_{n+1}$  over the set  $\Lambda_{n+1}$  yet to be defined, in conjunction with (2.3) and (2.28), implies

$$\begin{aligned} \|u - u_{n+1}\| &\leq \|u - (u_n + y_n)\| \leq \|u - u_n - w_n + z_n\| \\ &\leq \frac{1}{\sqrt{\alpha_*}} \|L(u - u_n - w_n)\|_{V^*} + \sqrt{\alpha^*} \|z_n\|_V \\ &= \frac{\beta^*}{\sqrt{\alpha_*}} \|r_n - g_n\|_{\phi^*} + \sqrt{\alpha^*} \|z_n\|_V , \end{aligned}$$

whence

$$\|u - u_{n+1}\| \leq \frac{\beta^*}{\sqrt{\alpha_*}} \sqrt{1 - \theta^2} \|r_n\|_{\phi^*} + \sqrt{\alpha^*} \|z_n\|_V .$$

Let  $\mathbf{z}_n$  and  $\mathbf{r}_n$  be the vectors of normalized coefficients of  $z_n$  and  $r_n$ , respectively. Since  $\mathbf{z}_n = \mathbf{P}_{\Lambda_{n+1}^c} \mathbf{A}^{-1} \mathbf{P}_{\partial\tilde{\Lambda}_n} \mathbf{r}_n$  and  $\|z_n\|_V \leq \frac{1}{\beta_*} \|z_n\|_{\phi} = \frac{1}{\beta_*} \|\mathbf{z}_n\|_{\ell^2}$ , we now construct  $\Lambda_{n+1}^c$  to control  $\|\mathbf{z}_n\|_{\ell^2}$ . If

$$k \in \Lambda_{n+1}^c \quad \text{and} \quad \ell \in \partial\tilde{\Lambda}_n \quad \Rightarrow \quad |k - \ell| > J ,$$

then we have  $\mathbf{P}_{\Lambda_{n+1}^c} (\mathbf{A}^{-1})_J \mathbf{P}_{\partial\tilde{\Lambda}_n} \mathbf{r}_n = 0$  which yields

$$\begin{aligned} \sqrt{\alpha^*} \|z_n\|_V &\leq \frac{\sqrt{\alpha^*}}{\beta_*} \|\mathbf{z}_n\|_{\ell^2} = \frac{\sqrt{\alpha^*}}{\beta_*} \|\mathbf{P}_{\Lambda_{n+1}^c} (\mathbf{A}^{-1} - (\mathbf{A}^{-1})_J) \mathbf{P}_{\partial\tilde{\Lambda}_n} \mathbf{r}_n\|_{\ell^2} \\ &\leq \frac{\sqrt{\alpha^*}}{\beta_*} \|\mathbf{A}^{-1} - (\mathbf{A}^{-1})_J\|_{\ell^2 \rightarrow \ell^2} \|\mathbf{r}_n\|_{\ell^2} \\ (3.6) \quad &\leq \frac{\sqrt{\alpha^*}}{\beta_*} C_{\mathbf{A}^{-1}} e^{-\tilde{\eta}_L J} \|r_n\|_{\phi^*} , \end{aligned}$$

where we have used (2.23). We now choose  $J = J(\theta) > 0$  to satisfy (3.3), and we obtain

$$(3.7) \quad \|u - u_{n+1}\| \leq 2 \frac{\beta^*}{\sqrt{\alpha_*}} \sqrt{1 - \theta^2} \|r_n\|_{\phi^*} \leq 2 \frac{\beta^* \sqrt{\alpha^*}}{\beta_* \sqrt{\alpha_*}} \sqrt{1 - \theta^2} \|u - u_n\| ,$$

as asserted. □

We observe that, as desired, the new error reduction rate

$$(3.8) \quad \bar{\rho}(\theta) = 2 \frac{\beta^* \sqrt{\alpha^*}}{\beta_* \sqrt{\alpha_*}} \sqrt{1 - \theta^2}$$

can be made arbitrarily small by choosing  $\theta$  suitably close to 1. This observation was already made in [4, 5], but we improve it in section 3.2 upon choosing  $\theta$  dynamically.

*Remark 3.2* (cardinality of  $\partial\Lambda_n$  for trigonometric basis). Since we add a ball of radius  $J$  around each point of  $\widetilde{\partial\Lambda}_n$  we get a crude estimate

$$(3.9) \quad |\partial\Lambda_n| \leq |B_d(0, J) \cap \mathbb{Z}^d| |\widetilde{\partial\Lambda}_n| \approx \omega_d J^d |\widetilde{\partial\Lambda}_n|,$$

where  $\omega_d$  is the measure of the  $d$ -dimensional Euclidean unit ball  $B(0, 1)$  centered at the origin.

**3.2. Dynamic Dörfler marking and adaptive spectral algorithm.** In this section we improve on the above marking strategy upon making the choice of  $\theta$  dynamic. At each iteration  $n$  let us select the Dörfler parameter  $\theta_n$  such that

$$(3.10) \quad \sqrt{1 - \theta_n^2} = C_0 \frac{\|r_n\|_{\phi^*}}{\|r_0\|_{\phi^*}}$$

for a proper choice of the positive constant  $C_0$  that will be made precise later. This implies

$$(3.11) \quad J(\theta_n) = -\frac{1}{\tilde{\eta}_L} \log \frac{\|r_n\|_{\phi^*}}{\|r_0\|_{\phi^*}} + K_1$$

according to (3.3), where  $K_1 := -\frac{1}{\tilde{\eta}_L} \log\left(\frac{\beta_*\beta^*}{\sqrt{\alpha_*\alpha^*} C_{\mathbf{A}^{-1}}}\right) + \delta_n$  and  $\delta_n \in [0, 1)$ .

We thus have the following adaptive spectral Galerkin method with dynamic choice (3.10) of the marking parameter  $\theta_n = (1 - C_0^2 \|r_n\|_{\phi^*}^2 / \|r_0\|_{\phi^*}^2)^{1/2}$ :

**DYN-GAL**( $\varepsilon$ )

set  $r_0 := f$ ,  $\Lambda_0 := \emptyset$ ,  $n = -1$

do

$n \leftarrow n + 1$

$\partial\Lambda_n := \mathbf{E-DÖRFLER}(\Lambda_n, (1 - C_0^2 \|r_n\|_{\phi^*}^2 / \|r_0\|_{\phi^*}^2)^{1/2})$

$\Lambda_{n+1} := \Lambda_n \cup \partial\Lambda_n$

$u_{n+1} := \mathbf{GAL}(\Lambda_{n+1})$

$r_{n+1} := \mathbf{RES}(u_{n+1})$

while  $\|r_{n+1}\|_{\phi^*} > \varepsilon \|r_0\|_{\phi^*}$

where **GAL** computes the Galerkin solution and **RES** the residual. The following result shows the quadratic convergence of **DYN-GAL**.

**THEOREM 3.3** (quadratic convergence). *Let the constant  $C_0$  of (3.10) satisfy  $C_0 \leq \frac{1}{4} \frac{\alpha_*}{\alpha^*} \left(\frac{\beta_*}{\beta^*}\right)^2$  and  $C_1 := \frac{\sqrt{\alpha^*}}{2\beta_* \|f\|_{\phi^*}}$ . Then the residual  $r_n$  of **DYN-GAL** satisfies*

$$(3.12) \quad \frac{\|r_{n+1}\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \leq \left(\frac{\|r_n\|_{\phi^*}}{2\|r_0\|_{\phi^*}}\right)^2 \quad \forall n \geq 0,$$

and the algorithm terminates in finite steps for any tolerance  $\varepsilon$ . In addition, two consecutive solutions of **DYN-GAL** satisfy

$$(3.13) \quad \|u - u_{n+1}\| \leq C_1 \|u - u_n\|^2 \quad \forall n \geq 0.$$

*Proof.* Invoke (3.7) and (3.10) to figure out that

$$(3.14) \quad \begin{aligned} \frac{\|r_{n+1}\|_{\phi^*}}{\|r_0\|_{\phi^*}} &\leq \frac{\sqrt{\alpha^*} \|u - u_{n+1}\|}{\beta_* \|r_0\|_{\phi^*}} \leq 2 \frac{\alpha^*}{\alpha_*} \left(\frac{\beta^*}{\beta_*}\right)^2 \sqrt{1 - \theta_n^2} \frac{\|r_n\|_{\phi^*}}{\|r_0\|_{\phi^*}} \\ &\leq 2C_0 \frac{\alpha^*}{\alpha_*} \left(\frac{\beta^*}{\beta_*}\right)^2 \left(\frac{\|r_n\|_{\phi^*}}{\|r_0\|_{\phi^*}}\right)^2 \leq \frac{1}{2} \left(\frac{\|r_n\|_{\phi^*}}{\|r_0\|_{\phi^*}}\right)^2 \end{aligned}$$

which implies (3.12). We thus realize that **DYN-GAL** converges quadratically and terminates in finite steps for any tolerance  $\varepsilon$ . Finally, combining (2.28) with (3.14), we readily obtain (3.13) upon using  $r_0 = f$ .  $\square$

*Remark 3.4* (superlinear rate). If the dynamic marking parameter  $\theta_n$  is chosen so that  $\sqrt{1 - \theta_n^2} = C_0 (\frac{\|r_n\|_{\phi^*}}{\|r_0\|_{\phi^*}})^\sigma$  for some  $\sigma > 0$ , then we arrive at the rate  $\|u - u_{n+1}\| \leq C_1 \|u - u_n\|^{1+\sigma}$ .

It seems to us that the quadratic rate (3.13) is the first one in adaptivity theory. The relation (3.12) reads equivalently

$$(3.15) \quad \frac{\|r_{n+1}\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \leq \left( \frac{\|r_{n+1-k}\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right)^{2^k}, \quad 0 \leq k \leq n + 1,$$

and implies that  $\|r_n\|_{\phi^*}/\|r_0\|_{\phi^*}$  is within machine precision in about  $n = 6$  iterations. This fast decay is consistent with spectral methods. Upon termination, we obtain the relative error

$$\|u - u_n\| \leq \frac{\beta^* \sqrt{\alpha^*}}{\beta_* \sqrt{\alpha_*}} \|u\| \varepsilon,$$

because  $\|f\|_{\phi^*} \leq \frac{\sqrt{\alpha^*}}{\beta_*} \|u\|$ .

The algorithm **DYN-GAL** entails exact computation of the residual  $r_n$ , which in general has infinitely many terms. We do not dwell here with inexact or feasible versions of **DYN-GAL** but refer to the ideas of [11] (see also [19]) and the discussion in [4, 5, 7] for a possible approach to extend our present *idealized* setting.

**4. Nonlinear approximation and Gevrey sparsity classes.** Given any  $v \in V$  we define its *best N-term approximation error* as

$$E_N(v) = \inf_{\Lambda \subset \mathcal{K}, |\Lambda|=N} \|v - P_\Lambda v\|_\phi.$$

We are interested in classifying functions  $v$  according to the decay law of  $E_N(v)$  as  $N \rightarrow \infty$ , i.e., according to the “sparsity” of their expansions in terms of the basis  $\{\phi_k\}_{k \in \mathcal{K}}$ . Of special interest to us is the following exponential Gevrey class.

**DEFINITION 4.1** (exponential class of functions). *For  $\eta > 0$  and  $0 < t \leq d$ , we denote by  $\mathcal{A}_G^{\eta,t}$  the subset of  $V$  defined as*

$$\mathcal{A}_G^{\eta,t} := \left\{ v \in V : \|v\|_{\mathcal{A}_G^{\eta,t}} := \sup_{N \geq 0} \left( E_N(v) \exp \left( \eta \omega_d^{-t/d} N^{t/d} \right) \right) < +\infty \right\},$$

where  $\omega_d$  is the measure of the  $d$ -dimensional Euclidean unit ball  $B_d(0, 1)$  centered at the origin.

**DEFINITION 4.2** (exponential class of sequences). *Let  $\ell_G^{\eta,t}(\mathcal{K})$  be the subset of sequences  $\mathbf{v} \in \ell^2(\mathcal{K})$  so that*

$$\|\mathbf{v}\|_{\ell_G^{\eta,t}(\mathcal{K})} := \sup_{n \geq 1} \left( n^{(1-t/d)/2} \exp \left( \eta \omega_d^{-t/d} n^{t/d} \right) |v_n^*| \right) < +\infty,$$

where  $\mathbf{v}^* = (v_n^*)_{n=1}^\infty$  is the nonincreasing rearrangement of  $\mathbf{v}$ .

The relationship between  $\mathcal{A}_G^{\eta,t}$  and  $\ell_G^{\eta,t}(\mathcal{K})$  is stated in the following [4, Proposition 4.2].

PROPOSITION 4.3 (equivalence of exponential classes). *Given a function  $v \in V$  and the sequence  $\mathbf{v} = (\hat{v}_k \sqrt{d_k})_{k \in \mathcal{K}}$  of its coefficients, one has  $v \in \mathcal{A}_G^{\eta,t}$  if and only if  $\mathbf{v} \in \ell_G^{\eta,t}(\mathcal{K})$  with*

$$\|v\|_{\mathcal{A}_G^{\eta,t}} \lesssim \|\mathbf{v}\|_{\ell_G^{\eta,t}(\mathcal{K})} \lesssim \|v\|_{\mathcal{A}_G^{\eta,t}}.$$

For functions  $v$  in  $\mathcal{A}_G^{\eta,t}$  one can estimate the minimal cardinality of a set  $\Lambda$  such that  $\|v - P_\Lambda v\|_\phi \leq \varepsilon$  as follows: since  $\|v - P_{\tilde{\Lambda}} v\|_\phi > \varepsilon$  for any set  $\tilde{\Lambda}$  with cardinality  $|\tilde{\Lambda}| = |\Lambda| - 1$ , we deduce

$$(4.1) \quad |\Lambda| \leq \omega_d \left( \frac{1}{\eta} \log \frac{\|v\|_{\mathcal{A}_G^{\eta,t}}}{\varepsilon} \right)^{d/t} + 1.$$

For the analysis of the optimality of our algorithm it is important to investigate the sparsity class of the image  $Lv$  for the operator  $L$  defined in (2.1), when the function  $v$  belongs to the sparsity class  $\mathcal{A}_G^{\eta,t}$ . Sparsity classes of exponential type for functionals  $f \in V^*$  can be defined analogously as above, using now the best  $N$ -term approximation error in  $V^*$ ,

$$E_N^*(f) = \inf_{\Lambda \subset \mathcal{K}, |\Lambda|=N} \|f - P_\Lambda^* f\|_{\phi^*}.$$

The following result is based on [4, Proposition 5.2].

PROPOSITION 4.4 (continuity of  $L$  in  $\mathcal{A}_G^{\eta,t}$ ). *Let  $L$  be such that the associated stiffness matrix  $\mathbf{A}$  satisfies the decay condition (2.21). Given  $\eta > 0$  and  $t \in (0, d]$ , there exist  $\bar{\eta} > 0$ ,  $\bar{t} \in (0, t]$ , and a constant  $C_L \geq 1$  such that*

$$(4.2) \quad \|Lv\|_{\mathcal{A}_G^{\bar{\eta},\bar{t}}} \leq C_L \|v\|_{\mathcal{A}_G^{\eta,t}} \quad \forall v \in \mathcal{A}_G^{\eta,t}.$$

*Proof.* Let  $\mathbf{A}$  be the stiffness matrix associated with the operator  $L$ . In [4] it is proven that if  $\mathbf{A}$  is banded with  $2p + 1$  nonzero diagonals, then the result holds with  $\bar{\eta} = \frac{\eta}{(2p+1)^{t/d}}$  and  $\bar{t} = t$ ; on the other hand, if  $\mathbf{A} \in \mathcal{D}_e(\eta_L)$  is dense, but the coefficients  $\eta_L$  and  $\eta$  satisfy the inequality  $\eta < \eta_L \omega_d^{t/d}$ , then the result holds with  $\bar{\eta} = \zeta(t)\eta$  and  $\bar{t} = \frac{t}{1+t}$ , where  $\zeta(t) = (\frac{1+t}{2^d \omega_d^{1+t}})^{\frac{t}{d(1+t)}}$ . Finally, if  $\eta \geq \eta_L \omega_d^{t/d}$ , we introduce an arbitrary  $\hat{\eta} > 0$  satisfying  $\hat{\eta} < \eta_L \omega_d^{t/d}$ ; then the result holds with  $\bar{\eta} = \zeta(t)\hat{\eta}$  and  $\bar{t} = \frac{t}{1+t}$ , since  $\|v\|_{\mathcal{A}_G^{\hat{\eta},t}} \leq \|v\|_{\mathcal{A}_G^{\eta,t}}$ .  $\square$

Taking into account that  $\zeta(t) \leq 1$  for  $1 \leq d \leq 10$  (see again [4]), this result indicates that the residual is expected to belong to a less favorable sparsity class than the one of the solution. Counterexamples in [4] show that (4.2) cannot be improved.

Finally, we discuss the sparsity class of the residual  $r = r(u_\Lambda)$  for any Galerkin solution  $u_\Lambda$ .

PROPOSITION 4.5 (sparsity class of the residual). *Let  $\mathbf{A} \in \mathcal{D}_e(\eta_L)$  and  $\mathbf{A}^{-1} \in \mathcal{D}_e(\bar{\eta}_L)$ , for constants  $\eta_L > 0$  and  $\bar{\eta}_L \in (0, \eta_L]$  so that (2.22) and (2.23) hold. If  $u \in \mathcal{A}_G^{\eta,t}$  for some  $\eta > 0$  and  $t \in (0, d]$ , then there exist suitable positive constants  $\bar{\eta} \leq \eta$  and  $\bar{t} \leq t$  such that  $r(u_\Lambda) \in \mathcal{A}_G^{\bar{\eta},\bar{t}}$  for any index set  $\Lambda$  and*

$$\|r(u_\Lambda)\|_{\mathcal{A}_G^{\bar{\eta},\bar{t}}} \lesssim \|u\|_{\mathcal{A}_G^{\eta,t}}.$$

*Proof.* Proposition 4.4 yields the existence of  $\bar{\eta} > 0$  and  $\bar{t} \in (0, t]$  such that

$$(4.3) \quad \|f\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}} = \|Lu\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}} \lesssim \|u\|_{\mathcal{A}_G^{\eta, t}}.$$

In order to bound  $\|r(u_\Lambda)\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}}$  in terms of  $\|f\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}}$ , let us write

$$\mathbf{r}_\Lambda = \mathbf{A}(\mathbf{u} - \mathbf{u}_\Lambda) = \mathbf{f} - \mathbf{A}\mathbf{u}_\Lambda,$$

then use  $\mathbf{u}_\Lambda = (\widehat{\mathbf{A}}_\Lambda)^{-1}(\mathbf{P}_\Lambda \mathbf{f})$  from (2.30) to get

$$\mathbf{r}_\Lambda = \mathbf{f} - \mathbf{A}(\widehat{\mathbf{A}}_\Lambda)^{-1}(\mathbf{P}_\Lambda \mathbf{f}).$$

Now, assuming just for simplicity that the indices in  $\Lambda$  come first (this can be realized by a permutation), we have

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_\Lambda & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{pmatrix} \quad \text{and} \quad \widehat{\mathbf{A}}_\Lambda = \begin{pmatrix} \mathbf{A}_\Lambda & \mathbf{O} \\ \mathbf{O}^T & \mathbf{I} \end{pmatrix},$$

whence  $(\widehat{\mathbf{A}}_\Lambda)^{-1} = \begin{pmatrix} (\mathbf{A}_\Lambda)^{-1} & \mathbf{O} \\ \mathbf{O}^T & \mathbf{I} \end{pmatrix}.$

Setting  $\mathbf{f} = (\mathbf{f}_\Lambda \ \mathbf{f}_{\Lambda^c})^T$ , so that  $\mathbf{P}_\Lambda \mathbf{f} = (\mathbf{f}_\Lambda \ \mathbf{0})^T$ , we have

$$\begin{aligned} \mathbf{A}(\widehat{\mathbf{A}}_\Lambda)^{-1}(\mathbf{P}_\Lambda \mathbf{f}) &= \begin{pmatrix} \mathbf{A}_\Lambda & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{pmatrix} \begin{pmatrix} (\mathbf{A}_\Lambda)^{-1} & \mathbf{O} \\ \mathbf{O}^T & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{f}_\Lambda \\ \mathbf{0} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{A}_\Lambda & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{pmatrix} \begin{pmatrix} (\mathbf{A}_\Lambda)^{-1} \mathbf{f}_\Lambda \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{f}_\Lambda \\ \mathbf{B}^T (\mathbf{A}_\Lambda)^{-1} \mathbf{f}_\Lambda \end{pmatrix}. \end{aligned}$$

Then,

$$\mathbf{r}_\Lambda = \begin{pmatrix} \mathbf{0} \\ \mathbf{f}_{\Lambda^c} - \mathbf{B}^T (\mathbf{A}_\Lambda)^{-1} \mathbf{f}_\Lambda \end{pmatrix} = \begin{pmatrix} \mathbf{O} & \mathbf{O} \\ -\mathbf{B}^T (\mathbf{A}_\Lambda)^{-1} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{f}_\Lambda \\ \mathbf{f}_{\Lambda^c} \end{pmatrix} =: \mathbf{R} \mathbf{f}.$$

Now, since  $\mathbf{A} \in \mathcal{D}_e(\eta_L)$  and  $(\mathbf{A}_\Lambda)^{-1} \in \mathcal{D}_e(\bar{\eta}_L)$ , it is easily seen that  $\mathbf{R} \in \mathcal{D}_e(\tilde{\eta}_L)$  with  $\tilde{\eta}_L = \bar{\eta}_L$  if  $\bar{\eta}_L < \eta_L$ , or  $\tilde{\eta}_L < \eta_L$  arbitrary if  $\bar{\eta}_L = \eta_L$ .

Finally, we apply Proposition 4.4 to the operator  $R$  defined by the matrix  $\mathbf{R}$ , obtaining the existence of constants  $\tilde{\eta} > 0$  and  $\tilde{t} \in (0, \bar{t}]$  such that

$$\|r(u_\Lambda)\|_{\mathcal{A}_G^{\tilde{\eta}, \tilde{t}}} \lesssim \|f\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}},$$

whence the result. □

**5. Optimality properties of DYN-GAL.** In this section we derive an exponential rate of convergence for  $\|u - u_n\|$  in terms of the number of degrees of freedom  $\Lambda_n$  activated by **DYN-GAL** and assess the computational work necessary to achieve this rate. This is made precise in the following theorem.

**THEOREM 5.1** (exponential convergence rate). *Let  $u \in \mathcal{A}_G^{\eta, t}$ , where the Gevrey class  $\mathcal{A}_G^{\eta, t}$  is introduced in Definition 4.1. Upon termination of **DYN-GAL**, the iterate  $u_{n+1} \in V_{\Lambda_{n+1}}$  and set of active coefficients  $\Lambda_{n+1}$  satisfy  $\|u - u_{n+1}\| \leq \frac{\beta^*}{\sqrt{\alpha^*}} \|f\|_{\phi^* \varepsilon}$*



and

$$(5.1) \quad |\Lambda_{n+1}| \leq \omega_d \left( \frac{1}{\eta_*} \log \frac{C_* \frac{\|u\|_{\mathcal{A}_G^{\eta_*,t}}}{\|f\|_{\phi^*}}}{\varepsilon} \right)^{d/t_*}$$

with parameters  $C_* > 0$ ,  $\eta_* < \eta$ , and  $t_* < t$ . Moreover, if the number of arithmetic operations needed to solve a linear system scales linearly with its dimension, then the workload  $W_\varepsilon$  of **DYN-GAL** upon completion satisfies

$$(5.2) \quad W_\varepsilon \leq \omega_d \left( \frac{1}{\eta^*} \log \frac{C^* \frac{\|u\|_{\mathcal{A}_G^{\eta^*,t}}}{\|f\|_{\phi^*}}}{\varepsilon^4 |\log |\log \varepsilon||^{-1}} \right)^{d/t_*},$$

where  $\eta^* < \eta_*$  and  $C^* > C_*$  but  $t_*$  remains the same as in (5.1).

*Proof.* We proceed in several steps.

1. *Expression of  $J(\theta_k)$ :* Our first task is to simplify the expression (3.11) for  $J(\theta_k)$ , namely, to absorb the term  $K_1$ : there is  $C_2 > 1$  such that

$$(5.3) \quad J(\theta_k) \leq \frac{C_2}{\tilde{\eta}_L} \left| \log \frac{\|r_k\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right|.$$

In fact, if  $C_2$  is given by  $C_2 = 1 + \frac{\tilde{\eta}_L \max(0, K_1)}{\log 2}$ , then

$$K_1 \leq \frac{C_2 - 1}{\tilde{\eta}_L} \log 2 \leq \frac{C_2 - 1}{\tilde{\eta}_L} \left| \log \frac{\|r_k\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right|$$

because  $\|r_k\|_{\phi^*} \leq \|r_0\|_{\phi^*}$  for all  $k \geq 0$  according to (3.14). This in turn implies (5.3)

$$J(\theta_k) \leq \frac{1}{\tilde{\eta}_L} \left| \log \frac{\|r_k\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right| + \frac{C_2 - 1}{\tilde{\eta}_L} \left| \log \frac{\|r_k\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right| = \frac{C_2}{\tilde{\eta}_L} \left| \log \frac{\|r_k\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right|.$$

2. *Active set  $\Lambda_k$ :* We now examine the output  $\Lambda_k$  of **E-DÖRFLER**. Employing the minimality of Dörfler marking and (3.2), we deduce

$$(5.4) \quad E_{|\partial\Lambda_k|}^*(r_k) = \|r_k - P_{\partial\Lambda_k}^* r_k\|_{\phi^*} \leq \sqrt{1 - \theta_k^2} \|r_k\|_{\phi^*},$$

which clearly implies  $E_{|\partial\Lambda_k|-1}^* > \sqrt{1 - \theta_k^2} \|r_k\|_{\phi^*}$ . The latter inequality, together with the Definition 4.2 of  $\|r_k\|_{\mathcal{A}_G^{\bar{\eta}, \bar{\varepsilon}}}$ , yields

$$\|r_k\|_{\mathcal{A}_G^{\bar{\eta}, \bar{\varepsilon}}} > \sqrt{1 - \theta_k^2} \|r_k\|_{\phi^*} \exp\left(\bar{\eta} \omega_d^{-\bar{\varepsilon}/d} (|\partial\Lambda_k| - 1)^{\bar{\varepsilon}/d}\right).$$

We note that this, along with  $|\partial\Lambda_k| \geq 1$ , ensures  $\frac{\|r_k\|_{\mathcal{A}_G^{\bar{\eta}, \bar{\varepsilon}}}}{\sqrt{1 - \theta_k^2} \|r_k\|_{\phi^*}} > 1$  whence

$$|\partial\Lambda_k| \leq \omega_d \left( \frac{1}{\bar{\eta}} \log \frac{\|r_k\|_{\mathcal{A}_G^{\bar{\eta}, \bar{\varepsilon}}}}{\sqrt{1 - \theta_k^2} \|r_k\|_{\phi^*}} \right)^{d/\bar{\varepsilon}} + 1.$$

We now recall the membership of the residual  $r_k := r(u_k)$  to the Gevrey class  $\mathcal{A}_G^{\bar{\eta}, \bar{t}}$  for all  $k \geq 0$ , established in Proposition 4.5: there exists  $C_3 > 0$  independent of  $k$  and  $u$  such that

$$\|r_k\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}} \leq C_3 \|u\|_{\mathcal{A}_G^{\eta, t}}.$$

Combining this with the dynamic marking (3.10) implies

$$\frac{\|r_k\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}}}{\sqrt{1 - \theta_k^2} \|r_k\|_{\phi^*}} = \frac{\|f\|_{\phi^*} \|r_k\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}}}{C_0 \|r_k\|_{\phi^*}^2} \leq \frac{C_4 \|f\|_{\phi^*} \|u\|_{\mathcal{A}_G^{\eta, t}}}{\|r_k\|_{\phi^*}^2}$$

with  $C_4 = C_3/C_0$ , whence

$$|\widetilde{\partial\Lambda}_k| \leq \omega_d \left( \frac{1}{\bar{\eta}} \log \frac{C_4 \|f\|_{\phi^*} \|u\|_{\mathcal{A}_G^{\eta, t}}}{\|r_k\|_{\phi^*}^2} \right)^{d/\bar{t}} + 1.$$

We let  $C_5$  satisfy  $1 = \omega_d (\frac{1}{\bar{\eta}} \log C_5)^{d/\bar{t}}$ , and use that  $\bar{t} < t \leq d$ , to obtain the simpler expression

$$(5.5) \quad |\widetilde{\partial\Lambda}_k| \leq \omega_d \left( \frac{1}{\bar{\eta}} \log \frac{C_4 C_5 \|f\|_{\phi^*} \|u\|_{\mathcal{A}_G^{\eta, t}}}{\|r_k\|_{\phi^*}^2} \right)^{d/\bar{t}}.$$

On the other hand, in view of (5.3), the enrichment step (3.9) of **E-DÖRFLER** yields

$$(5.6) \quad |\partial\Lambda_k| \leq C_6 \omega_d J(\theta_k)^d |\widetilde{\partial\Lambda}_k| \leq \frac{C_6 C_2^d \omega_d^2}{\bar{\eta}_L^d \bar{\eta}^{d/\bar{t}}} \left| \log \frac{\|r_k\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right|^d \left( \log \frac{C_4 C_5 \|f\|_{\phi^*} \|u\|_{\mathcal{A}_G^{\eta, t}}}{\|r_k\|_{\phi^*}^2} \right)^{d/\bar{t}}.$$

We exploit the quadratic convergence (3.12) to write for  $k \leq n$

$$(5.7) \quad \left| \log \frac{\|r_k\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right|^d \leq 2^{d(k-n)} \left| \log \frac{\|r_n\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right|^d = 2^{d(k-n)} \left( \log \frac{2\|f\|_{\phi^*}}{\|r_n\|_{\phi^*}} \right)^d.$$

Using the bound  $\|f\|_{\phi^*} \leq \|f\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}} = \|r_0\|_{\mathcal{A}_G^{\bar{\eta}, \bar{t}}} \leq C_3 \|u\|_{\mathcal{A}_G^{\eta, t}}$  in the two previous inequalities yields

$$\log \frac{C_4 C_5 \|f\|_{\phi^*} \|u\|_{\mathcal{A}_G^{\eta, t}}}{\|r_k\|_{\phi^*}^2} \leq 2 \log \frac{C_7 \|u\|_{\mathcal{A}_G^{\eta, t}}}{\|r_k\|_{\phi^*}} \quad \text{and} \quad \log \frac{2\|f\|_{\phi^*}}{\|r_n\|_{\phi^*}} \leq \log \frac{2C_3 \|u\|_{\mathcal{A}_G^{\eta, t}}}{\|r_n\|_{\phi^*}}$$

with  $C_7 = (C_3 C_4 C_5)^{1/2}$ . Introducing the constants

$$C_8 = \max(2C_3, C_7), \quad \frac{1}{\hat{\eta}} = \frac{2^{d/\bar{t}} C_6 C_2^d \omega_d}{\bar{\eta}_L^d \bar{\eta}^{d/\bar{t}}}, \quad t_* = \frac{\bar{t}}{1 + \bar{t}},$$

the derived upper bound for  $|\partial\Lambda_k|$  can be simplified as follows:

$$|\partial\Lambda_k| \leq 2^{d(k-n)} \frac{\omega_d}{\hat{\eta}} \left( \log \frac{C_8 \|u\|_{\mathcal{A}_G^{\eta, t}}}{\|r_n\|_{\phi^*}} \right)^{d/t_*}.$$

Recalling now that  $|\Lambda_0| = 0$  and for  $n \geq 0$

$$|\Lambda_{n+1}| = \sum_{k=0}^n |\partial\Lambda_k|,$$

we have

$$|\Lambda_{n+1}| \leq \frac{\omega_d}{\hat{\eta}} \left( \sum_{k=0}^n 2^{d(k-n)} \right) \left( \log \frac{C_8 \|u\|_{\mathcal{A}_G^{\eta,t}}}{\|r_n\|_{\phi^*}} \right)^{d/t_*}.$$

This can be written equivalently as

$$(5.8) \quad |\Lambda_{n+1}| \leq \omega_d \left( \frac{1}{\eta_*} \log \frac{C_8 \|u\|_{\mathcal{A}_G^{\eta,t}}}{\|r_n\|_{\phi^*}} \right)^{d/t_*}$$

with  $\eta_* = \left( \frac{\hat{\eta}}{\sum_{k=1}^{\infty} 2^{-dk}} \right)^{t_*/d}$ . At last, we make use of  $\|r_n\|_{\phi^*} > \varepsilon \|f\|_{\phi^*}$  to get the desired estimate (5.1) with  $C_* = C_8$ .

3. *Computational work:* Let us finally discuss the total computational work  $\mathcal{W}_\varepsilon$  of **DYN-GAL**. We start with some useful notation. We set  $\delta_n := \frac{\|r_n\|_{\phi^*}}{\|r_0\|_{\phi^*}}$  for  $n \geq 0$  and  $\varepsilon_{\ell+1} = \varepsilon_\ell^2$  for  $\ell \geq 1$  with  $\varepsilon_1 = \frac{1}{2}$ . We note that there exists an integer  $L > 0$  such that  $\varepsilon_{L+1} < \varepsilon \leq \varepsilon_L$  with  $\varepsilon$  being the tolerance of **DYN-GAL**. In addition, we observe that for every iteration  $n > 0$  of **DYN-GAL** there exists  $\ell > 0$  such that  $\delta_n \in I_\ell := (\varepsilon_{\ell+1}, \varepsilon_\ell]$  and that for each interval  $I_\ell$  there exists at most one  $\delta_n \in I_\ell$  because  $\delta_{n+1} \leq \frac{1}{2} \delta_n^2$  according to (3.12). Finally, to each interval  $I_\ell$  we associate the following computational work  $W_\ell$  to find and store  $u_{n+1}$ :

$$W_\ell = \begin{cases} C_\# |\Lambda_{n+1}| & \text{if there exists } n \text{ such that } \delta_n \in I_\ell, \\ 0 & \text{otherwise.} \end{cases}$$

This assumes that the number of arithmetic operations needed to solve the linear system for  $u_n$  scales linearly with its dimension and  $C_\#$  is an absolute constant that may depend on the specific solver. The total computational work of **DYN-GAL** is bounded by

$$\mathcal{W}_\varepsilon = \sum_{\ell=1}^L W_\ell.$$

We now get a bound for  $W_\varepsilon$ . In view of (5.1) we have

$$W_\ell \leq C_\# \omega_d \left( \frac{1}{\eta_*} \log \frac{C_* \frac{\|u\|_{\mathcal{A}_G^{\eta,t}}}{\|f\|_{\phi^*}}}{\varepsilon_{\ell+1}} \right)^{d/t_*} = C_\# \omega_d \left( \frac{1}{\eta_*} \log \frac{C_* \frac{\|u\|_{\mathcal{A}_G^{\eta,t}}}{\|f\|_{\phi^*}}}{\varepsilon_1^{2^\ell}} \right)^{d/t_*}.$$

Therefore, upon adding over  $\ell$  and using that  $d/t_* \geq 1$ , we obtain

$$\begin{aligned} \mathcal{W}_\varepsilon &\leq \frac{C_\# \omega_d}{\eta_*^{d/t_*}} \left( \sum_{\ell=1}^L \log C_* \frac{\|u\|_{\mathcal{A}_G^{\eta,t}}}{\|f\|_{\phi^*}} + \sum_{\ell=1}^L \log \varepsilon_1^{-2^\ell} \right)^{d/t_*} \\ &\leq \frac{C_\# \omega_d}{\eta_*^{d/t_*}} \left( \log \frac{LC_* \frac{\|u\|_{\mathcal{A}_G^{\eta,t}}}{\|f\|_{\phi^*}}}{\varepsilon_1^{2^{L+1}}} \right)^{d/t_*}. \end{aligned}$$

Since  $\varepsilon \leq \varepsilon_L = \varepsilon_1^{2^{L-1}}$ , we deduce  $\varepsilon^4 \leq \varepsilon_1^{2^{L+1}}$  and  $L \leq \frac{\log \frac{|\log \varepsilon|}{\log 2}}{\log 2} + 1 \leq C_9 \log |\log \varepsilon|$ . Inserting this bound in the preceding expression yields

$$(5.9) \quad \mathcal{W}_\varepsilon \leq \omega_d \left( \frac{1}{\eta^*} \log \frac{C^* \frac{\|u\|_{\mathcal{A}_2^{\eta,t}}}{\|f\|_{\phi^*}}}{\varepsilon^4 |\log |\log \varepsilon||^{-1}} \right)^{d/t_*}$$

with

$$\eta^* = \frac{\eta_*}{C_\#^{t_*/d}}, \quad C^* = C_9 C_*$$

which is the asserted estimate (5.2). The proof is thus complete.  $\square$

*Remark 5.2.* Note that the bound on the workload, given in (5.2), is at most an absolute multiple of the bound, given in (5.1), on the number of active coefficients.

*Remark 5.3* (superlinear convergence). If  $\sqrt{1 - \theta_n^2} = C_0 \left( \frac{\|r_n\|_{\phi^*}}{\|r_0\|_{\phi^*}} \right)^\sigma$  with  $\sigma > 0$ , then (5.1) still holds with the same parameters  $\eta_*, t_*$ .

*Remark 5.4* (algebraic class). Let us consider the case when  $u$  belongs to the algebraic class

$$\mathcal{A}_B^s := \left\{ v \in V : \|v\|_{\mathcal{A}_B^s} := \sup_{N \geq 0} E_N(v) (N+1)^{s/d} < +\infty \right\},$$

which is related to Besov regularity. We can distinguish two cases:

1.  $\mathbf{A}$  belongs to an exponential class but the residuals belong to an algebraic class;
2.  $\mathbf{A}$  belongs to an algebraic class  $\mathcal{D}_a(\eta_L)$ , i.e., there exists a constant  $c_L > 0$  such that its elements satisfy  $|a_{\ell,k}| \leq c_L (1 + |\ell - k|)^{-\eta_L}$ , and the residual belongs to an algebraic class.

We now study the optimality properties of **DYN-GAL** for these two cases. Let us first observe that whenever the residuals belong to the algebraic class  $\mathcal{A}_B^s$ , the bound (5.5) becomes

$$(5.10) \quad |\widetilde{\partial\Lambda}_k| \lesssim \|u - u_k\|_V^{-2d/s}.$$

This results from the dynamic marking (3.10) together with (5.4) in the algebraic case.

Let us start with Case 1. Using (5.10) and (5.3), which is still valid here as we assume that  $\mathbf{A}$  belongs to an exponential class, the bound (5.6) is replaced by

$$(5.11) \quad |\partial\Lambda_k| \lesssim \|u - u_k\|_V^{-2d/s} \left| \log \frac{\|r_k\|_{\phi^*}}{2\|r_0\|_{\phi^*}} \right|^d \lesssim \|u - u_k\|_V^{-2d/s} 2^{d(k-n)} \left( \log \frac{2\|f\|_{\phi^*}}{\|r_n\|_{\phi^*}} \right)^d,$$

where in the last inequality we have employed (5.7). Hence, we have

$$\begin{aligned} |\Lambda_{n+1}| &= \sum_{k=0}^n |\partial\Lambda_k| \lesssim \left( \log \frac{2\|f\|_{\phi^*}}{\|r_n\|_{\phi^*}} \right)^d \sum_{k=0}^n \|u - u_n\|_V^{-2\frac{d}{s}2^{k-n}} 2^{d(k-n)} \\ &\gtrsim \left( \log \frac{2\|f\|_{\phi^*}}{\|r_n\|_{\phi^*}} \right)^d \|u - u_n\|_V^{-2\frac{d}{s}} \sum_{k=0}^n 2^{d(k-n)} \\ &\gtrsim \left( \log \frac{2\|f\|_{\phi^*}}{\|r_n\|_{\phi^*}} \right)^d \|u - u_n\|_V^{-2\frac{d}{s}} \lesssim \left( \log \frac{2\|f\|_{\phi^*}}{\|u - u_n\|_V} \right)^d \|u - u_{n+1}\|_V^{-\frac{d}{s}}, \end{aligned}$$

where in the last inequality we have employed the quadratic convergence of **DYN-GAL**. The above result implies that **DYN-GAL** is optimal for Case 1 (up to a logarithmic factor).

Let us now consider Case 2. Since  $\mathbf{A}$  belongs to an algebraic class, (5.3) is replaced by

$$J(\theta_k) \approx \|u - u_k\|_V^{-1/s}.$$

Thus, the bound (5.6) is replaced by

$$|\partial\Lambda_k| \lesssim \|u - u_k\|_V^{-3d/s},$$

which implies

$$|\Lambda_{n+1}| = \sum_{k=0}^n |\partial\Lambda_k| \lesssim \|u - u_n\|_V^{-3\frac{d}{s}} \lesssim \|u - u_{n+1}\|_V^{-\frac{3}{2}\frac{d}{s}},$$

where in the last inequality we have again employed the quadratic convergence of **DYN-GAL**. This result is *not* optimal for Case 2, due to the factor  $2/3$  multiplying  $s$  in the exponent. We recall that the algorithm **FA-ADFOUR** of [4] is similar to **DYN-GAL** but with static marking parameter  $\theta$ . The theory of **FA-ADFOUR** requires neither a restriction on  $\theta$  nor coarsening and is proven to be optimal in the algebraic case; see [4, Theorem 7.2].

#### REFERENCES

- [1] P. BINEV, W. DAHMEN, AND R. DEVORE, *Adaptive finite element methods with convergence rates*, Numer. Math., 97 (2004), pp. 219–268.
- [2] M. BÜRG AND W. DÖRFLER, *Convergence of an adaptive hp finite element strategy in higher space-dimensions*, Appl. Numer. Math., 61 (2011), pp. 1132–1146.
- [3] C. CANUTO, R.H. NOCHETTO, R.P. STEVENSON, AND M. VERANI, *High-order adaptive Galerkin methods*, in Spectral and High Order Methods for Partial Differential Equations ICOSAHOM 2014, Lect. Notes Comput. Sci. Eng. 106, M. Berzins, J.S. Hesthaven, and R.M. Kirby, eds., Springer, Cham, Switzerland, 2014, pp. 51–72.
- [4] C. CANUTO, R.H. NOCHETTO, AND M. VERANI, *Adaptive Fourier-Galerkin methods*, Math. Comp., 83 (2014), pp. 1645–1687.
- [5] C. CANUTO, R.H. NOCHETTO, AND M. VERANI, *Contraction and optimality properties of adaptive Legendre-Galerkin methods: The one-dimensional case*, Comput. Math. Appl., 67 (2014), pp. 752–770.
- [6] C. CANUTO, R.H. NOCHETTO, R. STEVENSON, AND M. VERANI, *Convergence and optimality of hp-AFEM*, Numer. Math., doi:10.1007/S00211-016-0826-X (2016).
- [7] C. CANUTO, V. SIMONCINI, AND M. VERANI, *Contraction and optimality properties of an adaptive Legendre-Galerkin method: The multi-dimensional case*, J. Sci. Comput., 63 (2015), pp. 769–798, doi:10.1007/s10915-014-9912-3.

- [8] C. CANUTO AND M. VERANI, *On the numerical analysis of adaptive spectral/hp methods for elliptic problems*, in Analysis and Numerics of Partial Differential Equations, Springer INdAM Ser. 4, Springer, Milan, 2013, pp. 165–192.
- [9] C. CARSTENSEN, M. FEISCHL, M. PAGE, AND D. PRAETORIUS, *Axioms of adaptivity*, Comput. Math. Appl., 67 (2014), pp. 1195–1253.
- [10] J.M. CASCON, C. KREUZER, R.H. NOCHETTO, AND K.G. SIEBERT, *Quasi-optimal convergence rate for an adaptive finite element method*, SIAM J. Numer. Anal., 46 (2008), pp. 2524–2550.
- [11] A. COHEN, W. DAHMEN, AND R. DEVORE, *Adaptive wavelet methods for elliptic operator equations – convergence rates*, Math. Comp, 70 (1998), pp. 27–75.
- [12] L. DIENING, CH. KREUZER, AND R. STEVENSON, *Instance optimality of the adaptive maximum strategy*, Found. Comput. Math., 16 (2016), pp. 33–68.
- [13] W. DÖRFLER AND V. HEUVELINE, *Convergence of an adaptive hp finite element strategy in one space dimension*, Appl. Numer. Math., 57 (2007), pp. 1108–1124.
- [14] T. GANTUMUR, H. HARBRECHT, AND R. STEVENSON, *An optimal adaptive wavelet method without coarsening of the iterands*, Math. Comp., 76 (2007), pp. 615–629.
- [15] P. MORIN, R.H. NOCHETTO, AND K.G. SIEBERT, *Data oscillation and convergence of adaptive FEM*, SIAM J. Numer. Anal., 38 (2000), pp. 466–488.
- [16] R.H. NOCHETTO, K.G. SIEBERT, AND A. VEESER, *Theory of adaptive finite element methods: An introduction*, in Multiscale, Nonlinear and Adaptive Approximation, Springer, Berlin, 2009, pp. 409–542.
- [17] CH. SCHWAB, *p- and hp-Finite Element Methods: Theory and Applications in Solid and Fluid Mechanics*, Numerical Mathematics and Scientific Computation, Clarendon, New York, 1998.
- [18] R. STEVENSON, *Optimality of a standard adaptive finite element method*, Found. Comput. Math., 7 (2007), pp. 245–269.
- [19] R. STEVENSON, *Adaptive wavelet methods for solving operator equations: An overview*, in Multiscale, Nonlinear and Adaptive Approximation, Springer, Berlin, 2009, pp. 543–597.