# Nonlinear Galerkin Method Using Chebyshev and Legendre Polynomials I. The One-Dimensional Case

Jie Shen; Roger Temam

# NONLINEAR GALERKIN METHOD USING CHEBYSHEV AND LEGENDRE POLYNOMIALS I. THE ONE-DIMENSIONAL CASE*

JIE SHEN† AND ROGER TEMAM‡

**Abstract.** A new strategy, stemming from the nonlinear Galerkin method [M. Marion and R. Temam, *SIAM J. Numer. Anal.*, 26 (1989), pp. 1139–1157], for solving linear elliptic and nonlinear dissipative evolution equations by using Chebyshev and Legendre polynomials is presented. The essential idea is to decompose the solution into a low mode part and a high mode part and to treat them separately. The robustness of the method over the classical Galerkin method is substantiated by rigorous error estimates and preliminary numerical experiments.

**Key words.** spectral-Galerkin method, spectral-tau method, Chebyshev polynomial, Legendre polynomial, nonlinear Galerkin method

**AMS subject classifications.** 35A40, 65J15, 65M15, 65M70

**1. Introduction.** The exchange of energy between the low and high mode components of a flow is an important aspect of nonlinear phenomena that needs to be understood. An attempt to address this question from the computational point of view appeared with the nonlinear Galerkin method, which was introduced in relation with dynamical systems theory and the concept of inertial manifolds (see, e.g., [4], [5], [7], [8], [14], [20], [22]. The nonlinear Galerkin method, which has proven to be computationally efficient (cf. [6], [12], [11], is based upon a differentiated treatment of the low and high mode components of a flow.

More generally, even for linear evolution equations and in the absence of turbulence, certain forms of discretization produce, as we shall see, a coupling between the small and large scale components of a flow. In this case too there is need, for computational efficiency, to study the coupling between the small and large scale components. Indeed it is computationally inefficient to allocate as much computing resources to compute the small scale component of the flow carrying little energy as we do with the large scale component of the flow which carries most of the energy.

Our aim in this article is to study the coupling of the low and high mode components of a flow when spectral discretizations, using Chebyshev or Legendre polynomials, are implemented. This question is addressed for both linear elliptic equations and nonlinear dissipative evolution equations. We restrict ourselves in this article to the one-dimensional (1-D) case. In a subsequent article, we shall address the multidimensional cases. We consider the special approximation methods based on Chebyshev and Legendre polynomials, namely, the Galerkin approximation and the tau approximation: *our objective is to derive simplified versions of the classical algorithms that*

---

*produce better conditioned systems and a reduction in computing time without affecting the discretization error of the scheme under consideration.*

The paper is organized as follows. In the next section, we introduce several discrete spaces and related projectors that will be used in this work and we recall some basic properties of approximations by Legendre and Chebyshev polynomials. In §3, we describe in detail how we decompose the 1-D Poisson equation and Helmholtz equation. We have also performed some preliminary numerical experiments that confirm the theoretical results in this section. Finally, in §4, we propose and analyze two schemes of nonlinear Galerkin type, followed by a discussion and a comparison with the classical Galerkin scheme.

**2. Preliminaries.** Let $\Phi_n(x)$ be either the $n$th order Chebyshev $T_n(x)$ or the Legendre $L_n(x)$ polynomial; we set

$S_m = \text{span}\{\Phi_0(x), \Phi_1(x), \ldots, \Phi_m(x)\}$,

$Q_{dm} = \text{span}\{\Phi_{m-1}(x), \Phi_m(x), \ldots, \Phi_{dm}(x)\}$, where $d$ is some integer of our choice that may depend on $m$ (the relation between $d$ and $m$ will be clarified later). We set also

$\tilde{Q}_{dm} = \text{span}\{\Phi_{m-1}(x), \Phi_m(x), \ldots, \Phi_{dm-2}(x)\}$,

$V_m = \{v \in S_m : v(\pm 1) = 0\}$,

$W_{dm} = \{w \in Q_{dm} : w(\pm 1) = 0\}$.

We recall that $\Phi_n(x)$ is a polynomial of degree $n$ and therefore

(2.1) $$S_m = \text{span}\{1, x, \ldots, x^m\}.$$

It is also an easy matter to show that

(2.2) $$S_{m-2} \oplus Q_{dm} = S_{dm}, \quad S_{m-2} \oplus \tilde{Q}_{dm} = S_{dm-2}, \quad V_m \oplus W_{dm} = V_{dm}.$$

For the latter relation, we observe that $V_m \cap W_{dm} = \varnothing$. Indeed this space is spanned by $\Phi_{m-1}$ and $\Phi_m$. Taking into account the homogeneous boundary condition at $x = \pm 1$, we conclude that any element in this space is 0.

Let $I = (-1, 1)$, we denote by

$L_\omega^2(I)$: the weighted $L^2$ Hilbert space with the scalar product

$$(u, v)_\omega = \int_I u(x) v(x) \omega(x) dx \quad \forall u, v \in L_\omega^2(I),$$

and the norm $\|u\|_\omega = (u, u)_\omega^{1/2}$, where $\omega(x) \equiv 1$ in the Legendre case and $\omega = (1 - x^2)^{-1/2}$ in the Chebyshev case. In some cases, we will drop the subscript $\omega$ when $\omega \equiv 1$.

$H_\omega^n(I) = \{v \in L_\omega^2(I) : d^k v/dx^k \in L_\omega^2(I), \; k = 0, 1, \ldots, n\}$ the weighted Sobolev spaces with the norm $\|v\|_{n,\omega} = (\sum_{k=0}^n \int_I |d^k v/dx^k|^2 \omega dx)^{1/2}$. For any $s \geq 0$, $H_\omega^s(I)$ can then be defined by interpolation. We set in particular

$$H_{0,\omega}^1(I) = \{v \in H_\omega^1(I) : v(\pm 1) = 0\},$$

and for $u \in L_\omega^2(I)$, we define

$$\|u\|_{-1,\omega} = \sup_{v \in H_{0,\omega}^1(I)} \frac{(u, v)_\omega}{\|v\|_{1,\omega}}.$$

We recall that the $\{\Phi_n(x)\}$ satisfy the following orthogonality relations

(2.3) $$(\Phi_i(x), \Phi_j(x))_\omega = c_{i,\omega} \delta_{ij} \quad \forall i, j \geq 0,$$

for some constants $\{c_{i,\omega}\}$. From the above orthogonality relation and (2.1), we derive that

$$(2.4) \quad (y,z)_\omega = 0 \quad \forall y \in S_{m-2}, \; z \in Q_{dm}; \quad \text{and} \quad (y_{xx},z)_\omega = 0 \quad \forall y \in S_m, z \in Q_{dm}.$$

These two simple relations are essential for the decompositions that we will address in the next section.

Let us denote

$$a_\omega(u,v) = \int_I u_x(v\omega)_x dx.$$

We recall that $a_\omega(\cdot,\cdot)^{1/2}$ is a norm in $H^1_{0,\omega}(I)$ equivalent to $\|\cdot\|_{1,\omega}$ and furthermore $a_\omega(\cdot,\cdot)$ is coercive and continuous on $H^1_{0,\omega}(I) \times H^1_{0,\omega}(I)$. Namely, there exist $\alpha, \beta > 0$, such that

$$(2.5) \qquad a_\omega(u,v) \le \alpha \|u\|_{1,\omega} \|v\|_{1,\omega} \quad \forall u,v \in H^1_{0,\omega}(I),$$

$$(2.6) \qquad a_\omega(u,u) \ge \beta \|u\|^2_{1,\omega} \quad \forall u \in H^1_{0,\omega}(I).$$

The above two inequalities are trivially satisfied when $\omega \equiv 1$. For the Chebyshev weight, we refer to [2]. Hence $a_\omega(\cdot,\cdot)^{1/2}$ can be viewed as an equivalent norm on $H^1_{0,\omega}(I)$. To simplify our presentation, we set $\|u\|_{1,\omega} = a_\omega(u,u)^{1/2}$ and we will use $C$ to denote a generic constant that does not depend on $d$, $m$, and on any function.

We now introduce several projectors that will be used in the sequel:

$P_m$: the orthogonal projector in $L^2_\omega(I)$ onto $S_m$, namely,

$$(u - P_m u, v)_\omega = 0 \quad \forall v \in S_m, \quad u \in L^2_\omega(I);$$

$\Pi_m$: the projector in $H^1_{0,\omega}(I)$ onto $V_m$ defined by

$$a_\omega(u - \Pi_m u, v) = 0 \quad \forall v \in V_m, \quad u \in H^1_{0,\omega}(I);$$

$\Pi^2_m$: the projector in $H^2_\omega(I) \cap H^1_{0,\omega}(I)$ onto $V_m$ defined by

$$(2.7) \qquad -((u - \Pi^2_m u)_{xx}, v)_\omega = 0 \quad \forall v \in S_{m-2}, \quad u \in H^2_\omega(I) \cap H^1_{0,\omega}(I).$$

We recall that for $s \ge 0$ (see [3]),

$$(2.8) \qquad \|u - P_m u\|_\omega \le C m^{-s} \|u\|_{s,\omega} \quad \forall u \in H^s_\omega(I),$$

and for $s \ge 1$ (see, for instance, [16, Thm. 4.2]),

$$(2.9) \qquad \inf_{v \in V_m} \|u - v\|_{\nu,\omega} \le C m^{\nu-s} \|u\|_{s,\omega} \quad \forall 0 \le \nu \le s, \quad u \in H^s_\omega(I) \cap H^1_{0,\omega}(I).$$

As for the two other projectors, we can prove Lemma 2.1.

LEMMA 2.1. *For $s \ge 1$ and $u \in H^s_\omega(I) \cap H^1_{0,\omega}(I)$,*

$$(2.10) \qquad \|u - \Pi_m u\|_{\nu,\omega} \le C m^{\nu-s} \|u\|_{s,\omega} \quad \forall 0 \le \nu \le 1.$$

*For $s \ge 2$ and $u \in H^s_\omega(I) \cap H^1_{0,\omega}(I)$,*

$$(2.11) \qquad \|u - \Pi^2_m u\|_{\nu,\omega} \le C m^{\nu-s} \|u\|_{s,\omega} \quad \forall 0 \le \nu \le 2.$$

*Proof.* The proof of (2.10) is standard and can be found, for instance, in [13]. Hence we will only prove (2.11).

For $\nu = 2$, using (2.7) and since $v_{xx} \in S_{m-2}$ for all $v \in V_m$, we find

$$
\begin{aligned}
\|u - \Pi_m^2 u\|_{2,\omega}^2 &\leq C\|(u - \Pi_m^2 u)_{xx}\|_\omega^2 = ((u - \Pi_m^2 u)_{xx}, (u - \Pi_m^2 u)_{xx})_\omega \\
&= \inf_{v \in V_m} ((u - \Pi_m^2 u)_{xx}, (u - v)_{xx})_\omega \\
&\leq C\|u - \Pi_m^2 u\|_{2,\omega} \inf_{v \in V_m} \|u - v\|_{2,\omega}.
\end{aligned}
$$

Hence, by (2.9), we obtain

$$
(2.12) \qquad \|u - \Pi_m^2 u\|_{2,\omega} \leq C \inf_{v \in V_m} \|u - v\|_{2,\omega} \leq Cm^{2-s}\|u\|_{s,\omega}.
$$

For $\nu = 0$, we use the following modified duality argument. Let $v$ be the solution of

$$
(2.13) \qquad -(v\omega)_{xx} = (u - \Pi_m^2 u)\omega, \qquad v(\pm 1) = 0.
$$

In fact (2.13) is equivalent to

$$
(2.14) \qquad a_\omega(\phi, v) = (\phi_x, (v\omega)_x) = (u - \Pi_m^2 u, \phi)_\omega \quad \forall \phi \in H_{0,\omega}^1(I).
$$

By (2.5) and (2.6), we find that (2.14) has a unique solution $v \in H_{0,\omega}^1(I)$ and, furthermore,

$$
(2.15) \qquad \|v\|_{2,\omega} \leq C\|u - \Pi_m^2 u\|_\omega.
$$

Hence, by (2.13), (2.12), (2.10), and (2.15), we derive

$$
\begin{aligned}
\|u - \Pi_m^2 u\|_\omega^2 &= (u - \Pi_m^2 u, u - \Pi_m^2 u)_\omega = -(u - \Pi_m^2 u, (v\omega)_{xx}) \\
&= -((u - \Pi_m^2 u)_{xx}, v)_\omega = \inf_{\phi \in S_{m-2}} (-(u - \Pi_m^2 u)_{xx}, v - \phi)_\omega \\
&\leq C\|(u - \Pi_m^2 u)_{xx}\|_\omega \inf_{\phi \in S_{m-2}} \|v - \phi\|_\omega \\
&\leq Cm^{2-s}\|u\|_{s,\omega}(m-2)^{-2}\|v\|_{2,\omega} \leq Cm^{-s}\|u\|_{s,\omega}\|u - \Pi_m^2 u\|_\omega.
\end{aligned}
$$

The result for any $\nu \in (0, 2)$ can be derived by interpolation.     □

Finally, let us define a new operator $B_m$ from $H_{0,\omega}^1(I)$ into $V_m$ by

$$
(2.16) \qquad B_m u \in V_m, \quad (u - B_m u, v)_\omega = 0 \quad \forall v \in S_{m-2},
$$

i.e., $u$ and $B_m u$ share the first $m - 2$ modes, and the last two modes of $B_m u$ are determined by imposing the boundary conditions.

LEMMA 2.2. *For $s \geq 1$ and $u \in H_\omega^s(I) \cap H_{0,\omega}^1(I)$,*

$$
(2.17) \qquad \|u - B_m u\|_{\nu,\omega} \leq Cm^{\nu-s}\|u\|_{s,\omega} \quad \forall -1 \leq \nu \leq 1.
$$

*Proof.* By definition, we have

$$
(u - B_m u, v)_\omega = 0 \quad \forall v \in S_{m-2}.
$$

Since $v_{xx} \in S_{m-2}$ for any $v \in V_m$, we get

$$
(u - B_m u, v_{xx})_\omega = 0 \quad \forall v \in V_m.
$$

Hence, by (2.5) and (2.6), we find

$$
\begin{aligned}
\|u - B_m u\|_{1,\omega}^2 &= a_\omega(u - B_m u, u - B_m u) = -(u - B_m u, (u - B_m u)_{xx})_\omega \\
&= \inf_{v \in V_m} (-(u - B_m u), (u - v)_{xx})_\omega = \inf_{v \in V_m} a_\omega(u - v, u - B_m u) \\
&\leq \alpha \inf_{v \in V_m} \|u - v\|_{1,\omega}\|u - B_m u\|_{1,\omega}.
\end{aligned}
$$

Therefore, by (2.9),

$$(2.18) \qquad \|u - B_m u\|_{1,\omega} \leq \alpha \inf_{v \in V_m} \|u - v\|_{1,\omega} \leq Cm^{1-s}\|u\|_{s,\omega}.$$

For $\nu = -1$, we use the standard duality argument. For any $v \in H^1_{0,\omega}(I)$, let $w$ be the solution of

$$(2.19) \qquad -w_{xx} = v, \qquad w(\pm 1) = 0.$$

Then, by using (2.5), (2.9), (2.18), and the regularity of $w$, we find

$$\|u - B_m u\|_{-1,\omega} = \sup_{v \in H^1_{0,\omega}(I)} \frac{(u - B_m u, v)_\omega}{\|v\|_{1,\omega}} = \sup_{v \in H^1_{0,\omega}(I)} \frac{(u - B_m u, -w_{xx})_\omega}{\|v\|_{1,\omega}}$$

$$= \sup_{v \in H^1_{0,\omega}(I)} \inf_{\phi \in V_m} \frac{(u - B_m u, (\phi - w)_{xx})_\omega}{\|v\|_{1,\omega}}$$

$$= \sup_{v \in H^1_{0,\omega}(I)} \inf_{\phi \in V_m} \frac{a_\omega(w - \phi, u - B_m u)}{\|v\|_{1,\omega}}$$

$$\leq \alpha \|u - B_m u\|_{1,\omega} \sup_{v \in H^1_{0,\omega}(I)} \inf_{\phi \in V_m} \frac{\|w - \phi\|_{1,\omega}}{\|v\|_{1,\omega}}$$

$$\leq Cm^{1-s}\|u\|_{s,\omega} m^{-2} \sup_{v \in H^1_{0,\omega}(I)} \frac{\|w\|_{3,\omega}}{\|v\|_{1,\omega}} \leq Cm^{-1-s}\|u\|_{s,\omega}.$$

We then conclude by using the interpolation for $\nu \in (-1, 1)$. $\qquad \square$

*Remark* 2.1. It is interesting to note that in the Legendre case, $B_m \equiv \Pi_m$. Indeed by definition of $\Pi_m$,

$$0 = a(u - \Pi_m u, v) = -(u - \Pi_m u, v_{xx}) \quad \forall v \in V_m,$$

and it is obvious that the map $\frac{d^2}{dx^2} : V_m \to S_{m-2}$ is surjective and therefore

$$(u - \Pi_m u, v) = 0 \quad \forall v \in S_{m-2}.$$

Hence $\Pi_m u$ coincides with $B_m u$.

For an introduction and a more detailed presentation of the spectral methods using Chebyshev or Legendre polynomials, we refer to [10] and [1].

## 3. Low and high modes decomposition for solution of the 1-D Helmholtz equation. In this section, we shall approximate the 1-D Helmholtz equation

$$(3.1) \qquad \lambda u - u_{xx} = f \text{ in } I; \qquad u(\pm 1) = 0$$

with a new strategy. In fact, we shall decompose the approximate solution of (3.1) into two parts: one with only "low modes" and one with only "high modes" (we refer to the low (resp., high) mode as the low (resp., high) order Chebyshev or Legendre polynomials). We then obtain the approximate solution by adding the solutions of two decoupled subsystems.

### 3.1. Decomposition of the 1-D Poisson equation. The variational formulation of the Poisson equation ((3.1) with $\lambda = 0$) with homogeneous boundary condition in the 1-D case is the following.

Find $u \in H^1_{0,\omega}(I)$ such that

$$(3.2) \qquad a_\omega(u, v) = (f, v) \quad \forall v \in H^1_{0,\omega}(I).$$

The Galerkin approximation of (3.2) in $V_{dm}$ is the following.
    Find $u_{dm} \in V_{dm}$ such that

(3.3) $$a_\omega(u_{dm}, v) = (f, v)_\omega \quad \forall v \in V_{dm}.$$

**Legendre–Galerkin case.** We decompose (3.3) as follows.
    Find $y_m \in V_m$ and $z_{dm} \in W_{dm}$ such that

(3.4) $$a(y_m, v) = (f, v) \quad \forall v \in V_m,$$

(3.5) $$a(z_{dm}, w) = (f, w) \quad \forall w \in W_{dm}.$$

By integration by parts and thanks to (2.4), we find $a(y_m, w) = -(d^2 y_m/dx^2, w) = 0 \ \forall w \in W_{dm}$. Similarly we have $a(z_{dm}, v) = -(d^2 z_{dm}/dx^2, v) = -(z_{dm}, d^2 v/dx^2) = 0 \ \forall v \in V_m$. Therefore, we derive that $y_m + z_{dm} = u_{dm}$, which is the solution of the Legendre–Galerkin approximation (3.3) in $V_{dm}$.
    **Chebyshev–Galerkin case.** We decompose (3.3) as follows.
    Find $y_m \in V_m$ and $z_{dm} \in W_{dm}$ such that

(3.6) $$a_\omega(y_m + z_{dm}, v) = (f, v)_\omega \quad \forall v \in V_m,$$

(3.7) $$a_\omega(z_{dm}, w) = (f, w)_\omega \quad \forall w \in W_{dm}.$$

Once again since $a_\omega(y_m, w) = -((d^2 y_m/dx^2), w)_\omega = 0 \quad \forall w \in W_{dm}$, we conclude also that $y_m + z_{dm} = u_{dm}$.
    Therefore, from the error estimates for (3.3) (see, for instance, [10], [2]), we have the following results for both Legendre (3.4)–(3.5) and Chebyshev (3.6)–(3.7) approximations.
    LEMMA 3.1.

$$\|y_m + z_{dm} - u\|_{\nu,\omega} \le C(dm)^{\nu-s}\|f\|_{s-2,\omega} \quad \forall 0 \le \nu \le 1 < s.$$

We note that the system (3.4)–(3.5) or (3.6)–(3.7) can be solved efficiently by using the new technique for spectral-Galerkin methods introduced in [17].
    **Chebyshev-tau and Legendre-tau cases.** To avoid redundancy, we will treat Legendre and Chebyshev cases simultaneously. The tau approximation of the 1-D Poisson equation in $V_{dm}$ is the following (see, for instance, [10]).
    Find $u_{dm} \in V_{dm}$ such that

(3.8) $$-\left(\frac{d^2 u_{dm}}{dx^2}, v\right)_\omega = (f, v)_\omega \quad \forall v \in S_{dm-2}.$$

We can decompose (3.8) as follows.
    Find $y_m \in V_m$ and $z_{dm} \in W_{dm}$ such that

(3.9) $$-\left(\frac{d^2 y_m}{dx^2} + \frac{d^2 z_{dm}}{dx^2}, v\right)_\omega = (f, v)_\omega \quad \forall v \in S_{m-2},$$

(3.10) $$-\left(\frac{d^2 z_{dm}}{dx^2}, w\right)_\omega = (f, w)_\omega \quad \forall w \in \tilde{Q}_{dm}.$$

Since $((d^2 y_m/dx^2), w)_\omega = 0$ for all $y_m \in V_m$, $w \in \tilde{Q}_{dm}$, and $S_{m-2} \oplus \tilde{Q}_{dm} = S_{dm-2}$, it is an easy matter to verify that $y_m + z_{dm} = u_{dm}$, which is the solution of the tau approximation (3.8) in $V_{dm}$. Therefore, from the error estimate of the tau approximation (cf., for instance, [18]), we have

LEMMA 3.2.

$$\|y_m + z_{dm} - u\|_{\nu,\omega} \le C(dm)^{\nu-s}\|f\|_{s-2,\omega} \quad \forall 0 \le \nu \le 2 < s.$$

It is well known that the system (3.8) (hence (3.9) once $z_{dm}$ is known from (3.10)) can be transformed into an essentially diagonal system (cf. [10]). Similarly, we can also transform (3.10) into an essentially diagonal system. In fact, in the Chebyshev case, let $f = \sum_{n=0}^\infty f_n T_n(x)$ and $z_{dm} = \sum_{n=m-1}^{dm} a_n T_n(x)$, then (3.10) is equivalent to the following essentially diagonal system:

$$(3.11) \quad \begin{cases} a_n = \dfrac{f_{n-2}}{(4n(n-1))} - \dfrac{e_{n+2}f_n}{2(n^2-1)} + \dfrac{e_{n+4}f_{n+2}}{4n(n+1)}, \quad m+1 \le n \le dm, \\ \displaystyle\sum_{n=m-1}^{dm} (\pm 1)^n a_n = 0, \end{cases}$$

where $e_n = 1$ for $n \le dm$ and $e_n = 0$ for $n > dm$.

Knowing $z_{dm}$, we let $f + (d^2 z_{dm}/dx^2) = \sum_{n=0}^\infty \tilde{f}_n T_n(x)$ and $y_m = \sum_{n=0}^m b_n T_n(x)$, then (3.9) is equivalent to:

$$(3.12) \quad \begin{cases} b_n = \dfrac{c_{n-2}\tilde{f}_{n-2}}{(4n(n-1))} - \dfrac{\tilde{e}_{n+2}\tilde{f}_n}{2(n^2-1)} + \dfrac{\tilde{e}_{n+4}\tilde{f}_{n+2}}{4n(n+1)}, \quad 2 \le n \le m, \\ \displaystyle\sum_{n=0}^N (\pm 1)^n b_n = 0, \end{cases}$$

where $c_0 = 2$ and $c_n = 1$ for $n \ge 1$; $\tilde{e}_n = 1$ for $n \le m$; and $\tilde{e}_n = 0$ for $n > m$.

Hence we can determine $y_m$ and $z_{dm}$ by solving two essentially diagonal systems.

### 3.2. Decomposition of the 1-D Helmholtz equation.

*Case* a (Galerkin approximation). As before, we will treat Chebyshev and Legendre approximations together. The variational formulation of (3.1) (with $\lambda = 1$ for simplicity) is the following.

Find $u \in H^1_{0,\omega}(I)$ such that

$$(3.13) \qquad (u,v)_\omega + a_\omega(u,v) = (f,v)_\omega \quad \forall v \in H^1_{0,\omega}(I).$$

The Galerkin approximation of (3.13) in $V_{dm}$ is

$$(3.14) \qquad (u_{dm},v)_\omega + a_\omega(u_{dm},v) = (f,v)_\omega \quad \forall v \in V_{dm}.$$

We recall that the error estimate in Lemma 3.1 holds as well in this case. Unlike the previous case, we will not seek an exact decomposition of (3.14). Instead we are searching for an approximate decomposition maintaining the accuracy of (3.14). We propose to approximate (3.13) as follows:

Find $y_m \in V_m$ and $z_{dm} \in W_{dm}$ such that

$$(3.15) \qquad (y_m,v)_\omega + a_\omega(y_m + z_{dm},v) = (f,v)_\omega, \qquad v \in V_m,$$

$$(3.16) \qquad (z_{dm},w)_\omega + a_\omega(z_{dm},w) = (f,w)_\omega, \qquad w \in W_{dm}.$$

We note that in the Legendre case, since $\omega \equiv 1$, the term $a(z_{dm}, v)$ in (3.15) is equal to zero and hence can be deleted from the formulation. Once again the system (3.15)–(3.16) can be solved efficiently by using the new technique for spectral-Galerkin methods introduced in [17].

Since $(y, w)_\omega \neq 0$ and $(z, v)_\omega \neq 0$, $y_m + z_{dm}$ is no longer the solution of the Galerkin approximation $u_{dm}$ in (3.14). However, we have the following result.

THEOREM 3.1.

$$(3.17) \qquad \|u - y_m - z_{dm}\|_{1,\omega} \leq C((dm)^{1-s} + m^{-1-s})\|f\|_{s-2,\omega}.$$

*Proof.* Let us prove first the following inverse triangular inequality:

$$(3.18) \qquad \exists \delta > 0, \quad \text{s.t.} \ \|y + z\|_{1,\omega}^2 \geq \delta(\|y\|_{1,\omega}^2 + \|z\|_{1,\omega}^2) \quad \forall y \in V_m, \ z \in W_{dm}.$$

In fact by the definition of $B_m$, we have $B_m y = y$ and $B_m z = 0$ for any $y \in V_m$, $z \in W_{dm}$. Hence

$$y = B_m(y + z), \qquad z = (I - B_m)(y + z).$$

It then follows from Lemma 2.2 that there exists $\gamma > 0$, such that

$$\|y\|_{1,\omega} \leq \gamma\|y + z\|_{1,\omega}, \qquad \|z\|_{1,\omega} \leq \gamma\|y + z\|_{1,\omega}.$$

(3.18) follows from the above inequalities with $\delta = \gamma^2/2$.

Using the operator $\Pi_{dm}$, we derive from (3.13) that

$$(3.19) \qquad (u, v + w)_\omega + a_\omega(\Pi_{dm}u, v + w) = (f, v + w)_\omega \quad \forall v \in V_m, \quad w \in W_{dm}.$$

Let us denote $\xi = B_m u - y_m$, $\eta = (B_{dm} - B_m)u - z_{dm}$, and $e = \xi + \eta = B_{dm}u - y_m - z_{dm}$. First we need to prove $\eta \in W_{dm}$. By definition, we have

$$(u - B_m u, v)_\omega = 0 \quad \forall v \in S_{m-2}; \quad (u - B_{dm}u, v)_\omega = 0 \quad \forall v \in S_{dm-2}.$$

Therefore $(B_{dm}u - B_m u, v) = 0 \ \forall v \in S_{m-2}$. Hence $B_{dm}u - B_m u \in W_{dm}$ and $\eta \in W_{dm}$.

Now subtracting (3.15) and (3.16) from (3.19), we obtain

$$(3.20) \qquad \begin{aligned} (u - y_m, v)_\omega + (u - z_{dm}, w)_\omega &+ a_\omega(\Pi_{dm}u - B_{dm}u, v + w) \\ &+ a_\omega(e, v + w) = 0 \quad \forall v \in V_m, \quad w \in W_{dm}. \end{aligned}$$

Set $v = \xi$ and $w = \eta$ in (3.20), since

$$(u - y_m, v)_\omega = (u - B_m u, v)_\omega + (\xi, v)_\omega,$$

$$(u - z_{dm}, w)_\omega = (u - (B_{dm} - B_m)u, w)_\omega + (\eta, w)_\omega.$$

Thanks to (3.18), we find

$$(3.21) \qquad \begin{aligned} \|\xi\|_\omega^2 + \|\eta\|_\omega^2 &+ \frac{\delta}{2}(\|\xi\|_{1,\omega}^2 + \|\eta\|_{1,\omega}^2) + \frac{1}{2}\|e\|_{1,\omega}^2 \\ &\leq \|\xi\|_\omega^2 + \|\eta\|_\omega^2 + \|e\|_{1,\omega}^2 = -(u - B_m u, \xi)_\omega - (u, \eta)_\omega \\ &\quad + ((B_{dm} - B_m)u, \eta)_\omega - a_\omega(\Pi_{dm}u - B_{dm}u, e). \end{aligned}$$

Using the Schwarz inequality, we can bound the right-hand side of the above equation as follows:

$$-(u - B_m u, \xi)_\omega \leq \|u - B_m u\|_{-1,\omega}\|\xi\|_{1,\omega} \leq \frac{1}{2}(\delta^{-1}\|u - B_m u\|_{-1,\omega}^2 + \delta\|\xi\|_{1,\omega}^2);$$

$$-(u, \eta)_\omega = ((I - P_{m-2})u, \eta)_\omega$$

$$\leq \|u - P_{m-2}u\|_{-1,\omega}\|\eta\|_{1,\omega} \leq 2\delta^{-1}\|u - P_{m-2}u\|_{-1,\omega}^2 + \frac{\delta}{4}\|\eta\|_{1,\omega}^2.$$

Similarly,

$$
\begin{aligned}
(B_{dm}u - B_m u, \eta)_\omega &= (B_{dm}u - u, \eta)_\omega + (u - B_m u, \eta)_\omega \\
&\leq \frac{\delta}{4}\|\eta\|_{1,\omega}^2 + 2\delta^{-1}(\|u - B_m u\|_{-1,\omega}^2 + \|u - B_{dm}u\|_{-1,\omega}^2);
\end{aligned}
$$

$$
\begin{aligned}
-a_\omega(\Pi_{dm}u - B_{dm}u, e) &\leq \frac{1}{4}\|e\|_{1,\omega}^2 + \|\Pi_{dm}u - B_{dm}u\|_{1,\omega}^2 \\
&\leq \frac{1}{4}\|e\|_{1,\omega}^2 + (\|\Pi_{dm}u - u\|_{1,\omega}^2 + \|u - B_{dm}u\|_{1,\omega}^2).
\end{aligned}
$$

To bound $\|u - P_{m-2}u\|_{-1,\omega}$, we need to prove

$$(3.22) \qquad \|u - P_m u\|_{-1,\omega} \leq C m^{-1-s}\|u\|_{s,\omega}.$$

Indeed, by (2.8),

$$
\begin{aligned}
\|u - P_m u\|_{-1,\omega} &= \sup_{v \in H_{0,\omega}^1(I)} \frac{(u - P_m u, v)_\omega}{\|v\|_{1,\omega}} = \sup_{v \in H_{0,\omega}^1(I)} \frac{(u - P_m u, v - P_m v)_\omega}{\|v\|_{1,\omega}} \\
&\leq C\|u - P_m u\|_\omega \sup_{v \in H_{0,\omega}^1(I)} \frac{\|v - P_m v\|_\omega}{\|v\|_{1,\omega}} \leq C m^{-1-s}\|u\|_{s,\omega}.
\end{aligned}
$$

Combining these inequalities into (3.21), by virtue of (2.8), Lemmas 2.1 and 2.2, and (3.22), we find

$$
\begin{aligned}
\|\xi\|_\omega^2 + \|\eta\|_\omega^2 + \|e\|_{1,\omega}^2 &\leq C(\|u - B_m u\|_{-1,\omega}^2 + \|u - B_{dm}u\|_\omega^2 + \|u - P_{m-2}u\|_{-1,\omega}^2) \\
&\quad + C(\|u - \Pi_{dm}u\|_{1,\omega}^2 + \|u - B_{dm}u\|_{1,\omega}^2) \\
&\leq C(m^{2(-1-s)} + (dm)^{2(1-s)})\|u\|_{s,\omega}^2.
\end{aligned}
$$

Therefore, using the triangular inequality and Lemma 2.2, we conclude that

$$
\begin{aligned}
\|u - y_m - z_{dm}\|_{1,\omega} &\leq \|e\|_{1,\omega} + \|u - B_{dm}u\|_{1,\omega} \\
&\leq C((dm)^{1-s} + m^{-1-s})\|u\|_{s,\omega} \leq C((dm)^{1-s} + m^{-1-s})\|f\|_{s-2,\omega}. \quad \square
\end{aligned}
$$

*Remark* 3.1. We note that the first term of the error estimate (3.17) is inherited from the classical Galerkin approximation (3.14), and the second term of the error estimate (3.17) is introduced by dropping $(y_m, w)_\omega$ and $(z_{dm}, v)_\omega$ from (3.15)–(3.16). However, as long as $d^{s-1} \lesssim m^2$, (3.15)–(3.16) is as accurate as (3.14). See also some related comments in Remark 4.5.

*Case* b (Tau approximations). The original tau approximation in $V_{dm}$ for (3.1) with $\lambda = 1$ is

$$(3.23) \qquad (u_{dm}, v)_\omega - \left(\frac{d^2 u_{dm}}{dx^2}, v\right)_\omega = (f, v)_\omega \quad \forall v \in S_{dm-2}.$$

It is easy to show that the optimal error estimate in Lemma 3.2 holds as well in this case. We propose to decompose approximately (3.23) by the following decoupled system.

Find $y_m \in V_m$ and $z_{dm} \in W_{dm}$ such that

$$(3.24) \qquad (y_m, v)_\omega - \left(\frac{d^2 y_m}{dx^2} + \frac{d^2 z_{dm}}{dx^2}, v\right)_\omega = (f, v)_\omega \quad \forall v \in S_{m-2},$$

$$(3.25) \qquad (z_{dm}, w)_\omega - \left(\frac{d^2 z_{dm}}{dx^2}, w\right)_\omega = (f, w)_\omega \quad \forall w \in \tilde{Q}_m.$$

We note that similar to (3.9)–(3.10), the systems for (3.24)–(3.25) can be transformed into two essentially tridiagonal form and can be efficiently solved in practice. As in the Galerkin case, $y_m + z_{dm} \neq u_{dm}$. However, we can prove the following result in Theorem 3.2.

THEOREM 3.2. *For $m$ sufficiently large, we have*

$$\|u - y_m - z_{dm}\|_{2,\omega} \leq C((dm)^{2-s} + m^{-s})\|f\|_{s-2,\omega}.$$

*Proof.* Using the operator $\Pi^2_{dm}$, we derive from (3.1) with $\lambda = 1$ that

$$(u, v + w)_\omega - \left(\frac{d^2 \Pi^2_{dm} u}{dx^2}, v + w\right)_\omega = (f, v + w)_\omega \quad \forall v \in S_{m-2}, \quad w \in \tilde{Q}_{dm}.$$

Let us denote $\xi = \Pi^2_m u - y_m$, $\eta = (\Pi^2_{dm} - \Pi^2_m)u - z_{dm}$, and $e = \xi + \eta = \Pi^2_{dm} u - y_m - z_{dm}$. Subtracting (3.24) and (3.25) from the above equation, we obtain

$$(3.26) \quad (u - y_m, v)_\omega + (u - z_{dm}, w)_\omega - \left(\frac{d^2 e}{dx^2}, v + w\right)_\omega = 0 \quad \forall v \in S_{m-2}, \ w \in \tilde{Q}_{dm}.$$

Now we set $v = -P_{m-2}e_{xx}$ and $w = -(I - P_{m-2})\eta_{xx}$ in (3.26) and treat each terms as follows.

Since $\xi_{xx} \in S_{m-2}$, we derive that $w = -(I - P_{m-2})\eta_{xx} = -(I - P_{m-2})e_{xx}$ and $v + w = -e_{xx}$. Therefore, for the last term in (3.26), we have

$$\begin{aligned}
-(e_{xx}, -e_{xx})_\omega &= \frac{1}{2}(e_{xx}, e_{xx})_\omega + \frac{1}{2}(v + w, v + w)_\omega \\
&= \frac{1}{2}(\|e_{xx}\|^2_\omega + \|v\|^2_\omega + \|w\|^2_\omega) \\
&= \frac{1}{2}(\|e_{xx}\|^2_\omega + \|P_{m-2}e_{xx}\|^2_\omega + \|(I - P_{m-2})\eta_{xx}\|^2_\omega).
\end{aligned}$$

For the first and second terms in (3.26), we have

$$(3.27) \qquad (u - y_m, -P_{m-2}e_{xx})_\omega = -(u - \Pi^2_m u, P_{m-2}e_{xx})_\omega - (\xi, P_{m-2}e_{xx})_\omega,$$

$$(3.28) \qquad (u - z_{dm}, w)_\omega = (u, w)_\omega - ((\Pi^2_{dm} - \Pi^2_m)u, w)_\omega + (\eta, w)_\omega.$$

To treat the five terms on the right-hand sides of (3.27)–(3.28), we will use frequently integration by part, (2.8), and the Schwarz inequality. Indeed,

$$(\eta, w)_\omega = (\eta, -(I - P_{m-2})\eta_{xx})_\omega = (\eta, -\eta_{xx})_\omega = \|\eta\|^2_{1,\omega};$$

$$\begin{aligned}
(\xi, -P_{m-2}e_{xx})_\omega &= -(\xi, \xi_{xx})_\omega - (\xi, P_{m-2}\eta_{xx})_\omega = \|\xi\|^2_{1,\omega} + ((I - P_{m-2})\xi, \eta_{xx})_\omega \\
&\geq \|\xi\|^2_{1,\omega} - \|(I - P_{m-2})\xi\|_\omega \|(I - P_{m-2})\eta_{xx}\|_\omega \\
&\geq \|\xi\|^2_{1,\omega} - \frac{1}{8}\|(I - P_{m-2})\eta_{xx}\|^2_\omega - 2\|(I - P_{m-2})\xi\|_\omega \\
&\geq \|\xi\|^2_{1,\omega} - \frac{1}{8}\|(I - P_{m-2})\eta_{xx}\|^2_\omega - C(m - 2)^{-1}\|\xi\|_{1,\omega}.
\end{aligned}$$

Therefore, for $m$ sufficiently large, we have

$$(\xi, -P_{m-2}e_{xx})_\omega \geq C\|\eta\|^2_{1,\omega} - \frac{1}{8}\|(I - P_{m-2})\eta_{xx}\|^2_\omega.$$

Similarly,

$$(u, w)_\omega = (u, -(I - P_{m-2})\eta_{xx})_\omega \geq - \|(I - P_{m-2})u\|_\omega \|(I - P_{m-2})\eta_{xx}\|_\omega$$
$$\geq - \frac{1}{8}\|(I - P_{m-2})\eta_{xx}\|_\omega^2 - 2\|(I - P_{m-2})u\|_\omega^2.$$

The last two terms can be easily bounded:

$$|(u - \Pi_m^2 u, P_{m-2}e_{xx})_\omega| \leq 2\|u - \Pi_m^2 u\|_\omega^2 + \frac{1}{8}\|P_{m-2}e_{xx}\|^2,$$

$$|((\Pi_{dm}^2 - \Pi_m^2)u, w)_\omega| = |((\Pi_{dm}^2 - \Pi_m^2)u, (I - P_{m-2})\eta_{xx})_\omega|$$
$$\leq 2\|(\Pi_{dm}^2 - \Pi_m^2)u\|_\omega^2 + \frac{1}{8}\|(I - P_{m-2})\eta_{xx}\|_\omega^2.$$

Combining these inequalities into (3.26), using (2.8), Poincaré's inequality, and Lemma 2.1, we find

$$\|e_{xx}\|_\omega^2 + \|\xi\|_{1,\omega}^2 + \|\eta\|_{1,\omega}^2 + \|(I - P_{m-2})\eta_{xx}\|_\omega^2$$
$$\leq C \left\{ \|u - \Pi_m^2 u\|_\omega^2 + \|u - \Pi_{dm}^2 u\|_\omega^2 + \|(I - P_{m-2})u\|_\omega^2 \right\}$$
$$\leq C \left\{ m^{-2s} + (dm)^{-2s} + (m - 2)^{-2s} \right\} \|u\|_{s,\omega}^2$$
$$\leq Cm^{-2s}\|u\|_{s,\omega}^2.$$

Finally, by Lemma 2.1, we conclude that

$$\|u - y_m - z_{dm}\|_{2,\omega} \leq \|e\|_{2,\omega} + \|u - \Pi_{dm}^2 u\|_{2,\omega} \leq \|e_{xx}\|_\omega + \|u - \Pi_{dm}^2 u\|_{2,\omega}$$
$$\leq C(m^{-s} + (dm)^{2-s})\|u\|_{s,\omega} \leq C(m^{-s} + (dm)^{2-s})\|f\|_{s-2,\omega}. \quad \square$$

*Remark 3.2.*

(a) From Theorem 3.2, we realize that as long as $d^{s-2} \lesssim m^2$, (3.24)–(3.25) is as accurate as the tau approximation in $V_{dm}$ (see also Remark 3.1).

(b) For each of the cases considered above, we have successfully decomposed the system in $V_{dm}$ into two *decoupled* subsystems. One of the nice features of the decompositions is that the two subsystems have substantially smaller condition numbers than the original system. Hence, in general, the larger $d$ is, the better conditioned the system becomes. However, the smoother the solution is, the larger $s$ is and the smaller can we take $d$.

(c) Similarly, we can also apply our decomposition techniques to the biharmonic equation or more generally to an equation of the form $D_x^{2k}u + Lu = f$ with $k = 1$ or 2 and with a lower order (compared with $D_x^{2k}$) constant coefficient linear operator $L$. In fact, instead of using two coupled modes between $V_m$ and $W_{dm}$ for the Laplacian operator, we should use four coupled modes between $V_m$ and $W_{dm}$ for the biharmonic operator.

(d) See some related comments in Remark 4.5.

**3.3. Numerical results.** We have implemented the algorithm (3.15)–(3.16) using the Legendre polynomials for solving the following Helmholtz equation

$$(3.29) \qquad u - u_{xx} = f(x) \text{ in } I, \quad u(-1) = 0, \quad u(+1) = 1,$$

where $f(x)$ is given by

$$f(x) = \begin{cases} 0, & x \in [-1, 0], \\ x^\gamma - \gamma(\gamma - 1)x^{\gamma-2}, & x \in (0, +1]. \end{cases}$$

The exact solution of this problem is

$$u(x) = \begin{cases} 0, & x \in [-1,0], \\ x^\gamma, & x \in (0,+1]. \end{cases}$$

We note that in solving (3.29), we actually use the so-called pseudo-spectral Galerkin method, which is a common practice in the implementation of the spectral method. More precisely, $f$ is replaced by $I_{dm}f$, the polynomial of degree less than or equal to $dm$, which interpolates $f$ at the Gauss–Lobatto points. Hence the pseudo-spectral treatment introduces an extra error term: $\|f - I_{dm}f\|_\omega \leq C(dm)^{-\sigma}\|f\|_{\sigma,\omega}$ (see, for instance, [1]). In summary let $u_{dm}$ (resp., $y_m + z_{dm}$) be the approximate solution given by the pseudo-spectral Galerkin (resp., nonlinear Galerkin) scheme, then (cf. Theorem 3.1)

$$\|u - u_{dm}\|_{1,\omega} \leq C\left((dm)^{1-s}\|f\|_{s-2,\omega} + (dm)^{-\sigma}\|f\|_{\sigma,\omega}\right),$$

$$\|u - y_m - z_{dm}\|_{1,\omega} \leq C\left(((dm)^{1-s} + m^{-1-s})\|f\|_{s-2,\omega} + (dm)^{-\sigma}\|f\|_{\sigma,\omega}\right).$$

Since, for any $\gamma \geq 2$, we have

$$f(x) \in H^{\gamma-3/2-\varepsilon}(I), \quad u(x) \in H^{\gamma+1/2-\varepsilon}(I) \quad \forall \varepsilon > 0.$$

Hence, we should expect a convergence rate of order $\gamma - \frac{3}{2}$ for the error in $H_0^1(I)$ norm.

In Table 1, we have listed the errors in $H_0^1(I)$ by using the classical Galerkin scheme (denoted by GAL) (3.14) in $V_{dm}$ and the "nonlinear" Galerkin scheme (denoted by NLG) (3.15)–(3.16) in $V_m \times W_{dm}$ with various choices of $d$ for $\gamma = 2$. For the Galerkin scheme (3.14), the errors in $H_0^1(I)$ norm are defined as

$$Err(dm) := \left\{\frac{1}{dm+1}\sum_{i=0}^{dm}\frac{|u_x(x_i) - (u_{dm})_x(x_i)|^2}{\sup_{0\leq i\leq dm}|u_x(x_i)|^2}\right\}^{\frac{1}{2}}.$$

For the scheme (3.15)–(3.16), the errors are computed according to

$$Err(dm) := \left\{\frac{1}{dm+1}\sum_{i=0}^{dm}\frac{|u_x(x_i) - (y_m)_x(x_i) - (z_{dm})_x(x_i)|^2}{\sup_{0\leq i\leq dm}|u_x(x_i)|^2}\right\}^{\frac{1}{2}}.$$

In Table 2, the same type of errors are listed for the case $\gamma = 3$.

In Table 3, we have tabulated $Err(dm)/Err(2dm)$ for the Galerkin scheme (3.14) with $\gamma = 2$ and $\gamma = 3$. The results clearly indicate that for $\gamma = 2$, the scheme (3.14) is first order convergent and for $\gamma = 3$, it becomes second order convergent. The results led us to conclude that in practice (3.14) has a convergence rate of order $\gamma - 1$ for the given problem, which is essentially half order better than the theory predicted.

Therefore, to keep the scheme (3.15)–(3.16) as accurate as (3.14), the following condition should be satisfied:

$$(3.30) \qquad\qquad\qquad d^{\gamma-1} \lesssim m^2.$$

The results in Tables 1 and 2 indicate that as long as (3.30) is satisfied, the solutions of (3.15)–(3.16) are essentially the same as that given by (3.14). However, when (3.30) is violated, the solutions of (3.15)–(3.16) are less accurate than that of (3.14), which are sometimes substantially less. Hence (3.30) is indeed the correct criterion for choosing $d$.

TABLE 1
*Error in $H_0^1(I)$ Norm:* $\gamma = 2$.

| dm | GAL $(d = 1)$ | NLG $d = 2$ | NLG $d = 4$ | NLG $d = 8$ | NLG $d = 16$ | NLG $d = 32$ |
|---|---|---|---|---|---|---|
| 16 | 6.23E−2 | 6.24E−2 | 6.25E−2 | | | |
| 32 | 3.02E−2 | 3.02E−2 | 3.02E−2 | 3.04E−2 | | |
| 64 | 1.47E−2 | 1.47E−2 | 1.47E−2 | 1.47E−2 | 1.52E−2 | |
| 128 | 7.22E−3 | 7.22E−3 | 7.22E−3 | 7.22E−3 | 7.22E−3 | 8.12E−3* |
| 256 | 3.56E−3 | 3.56E−3 | 3.56E−3 | 3.56E−3 | 3.56E−3 | 3.56E−3 |

\* Condition (3.30) is violated!

TABLE 2
*Error in $H_0^1(I)$ Norm:* $\gamma = 3$.

| dm | GAL $(d = 1)$ | NLG $d = 2$ | NLG $d = 4$ | NLG $d = 8$ | NLG $d = 16$ | NLG $d = 32$ |
|---|---|---|---|---|---|---|
| 16 | 9.25E−3 | 9.26E−3 | 1.10E−2 | | | |
| 32 | 2.23E−3 | 2.23E−3 | 2.24E−2 | | | |
| 64 | 5.50E−4 | 5.50E−4 | 5.50E−4 | 5.55E−4 | 4.60E−3* | |
| 128 | 1.36E−4 | 1.36E−4 | 1.36E−4 | 1.36E−4 | 1.41E−4* | |
| 256 | 3.38E−5 | 3.38E−5 | 3.38E−5 | 3.38E−5 | 3.39E−5 | 3.88E−5* |

\* Condition (3.30) is violated!

TABLE 3
*$Err(dm)/Err(2dm)$ of the Galerkin approximation.*

| | 16/32 | 32/64 | 64/128 | 128/256 |
|---|---|---|---|---|
| $\gamma = 2$ | 2.06 | 2.05 | 2.04 | 2.03 |
| $\gamma = 3$ | 4.15 | 4.05 | 4.04 | 4.02 |

## 4. Nonlinear Galerkin method.

We consider a class of nonlinear evolution equations of the form

$$(4.1) \qquad \frac{\partial u}{dt} - \nu u_{xx} + B(u) = f \quad \text{in } I \times [0, T].$$

Equation (4.1) is supplemented with initial condition $u(x, 0) = u_0(x)$ and with the homogeneous Dirichlet boundary conditions $u(\pm 1, t) = 0$. We assume that the nonlinear term can be written as $B(u) = B(u, u)$, where $B(\cdot, \cdot)$ is a bilinear form.

To simplify our presentation, we will restrict ourselves to the Legendre–Galerkin case. Hence $\omega \equiv 1$ and we write $\| \cdot \|_{m,\omega} = \| \cdot \|_m$. We set $V = H_0^1(I)$ and $b(u, v, w) = (B(u, v), w)$ and assume that the trilinear form $b$ satisfies

$$(4.2) \qquad b(u, v, v) = 0 \quad \forall u, v \in V,$$

and

$$(4.3) \qquad b(u, v, w) \leq C \|u\|_V \|v\|_V \|w\| \quad \forall u, v, w \in V.$$

The variational formulation of (4.1) supplemented by the initial and boundary conditions is

$$(4.4) \qquad \frac{d}{dt}(u, v) + \nu a(u, v) + b(u, u, v) = (f, v) \quad \forall v \in V, \quad t \in [0, T],$$
$$u(\cdot, 0) = u_0(\cdot).$$

It is standard to show that there exists a unique solution $u$ of (4.4) such that $u \in C([0,T]; V)$.

In the sequel, we will use $M$ to denote a generic constant that depends on $\nu$, $T$, $u_0$, and on the solution $u$ through $M_1 = \sup_{t \in [0,T]} \|u(t)\|_V$.

*Example* 1 (1-D Burgers' equation).

$$B(u) = uu_x, \quad B(u,v) = \frac{2}{3}uv_x + \frac{1}{3}vu_x.$$

It is an easy matter to verify by using the Sobolev and Young inequalities that (4.3) holds.

*Example* 2 (The Kuramoto–Sivashinsky equation). This equation reads as follows:

$$u_t - \nu u_{xxxx} + u_{xx} + uu_x = f, \quad \text{in } I \times [0,T],$$

with appropriate initial and boundary conditions. This equation is not exactly of the type (4.1) since the principal linear operator is of fourth order. However, as noted in Remark 3.2(c), the fourth order operator can also be decomposed accordingly; note that the fourth order term is dissipative while the second order one is antidissipative and drives the flow. As for the nonlinear term, we can define $B(u)$ and $B(u,v)$ similar to those for the Burgers' equation, and relations corresponding to (4.2) and (4.3) can also be established accordingly (with $\|\cdot\|_V = \|\cdot\|_2$).

The Galerkin approximation for (4.1) in $V_{dm}$ is as follows.

Find $u_{dm} \in V_{dm}$ such that

$$(4.5) \qquad \frac{d}{dt}(u_{dm}, v) + \nu a(u_{dm}, v) + b(u_{dm}, u_{dm}, v) = (f, v) \quad \forall v \in V_{dm},$$

with $u_{dm}(\cdot, 0) = \Pi_{dm} u_0(\cdot)$.

**4.1. A first scheme.** We assume here $f$ is independent of time $t$; the cases with a time dependent force $f$ will be treated later in this section. We propose to approximate (4.4) by the following scheme of nonlinear Galerkin type.

Find $y_m \in V_m$ and $z_{dm} \in W_{dm}$ such that

$$(4.6) \qquad \frac{d}{dt}(y_m, v) + \nu a(y_m, v) + b(y_m + z_{dm}, y_m + z_{dm}, v) = (f, v) \quad \forall v \in V_m,$$

$$(4.7) \qquad \nu a(z_{dm}, w) + b(y_m + z_{dm}, y_m, v) = (f, w) \quad \forall w \in W_{dm},$$

with $y_m(\cdot, 0) = \Pi_m u_0(\cdot)$.

*Remark* 4.1. The scheme (4.6)–(4.7) is a variation of the original nonlinear Galerkin scheme proposed in [14]. The treatment of the nonlinear term here is computationally more efficient.

THEOREM 4.1. *Under the assumptions* (4.2) *and* (4.3), *and for $m$ sufficiently large, we have*

$$(4.8) \qquad \left( \int_0^T \|(u - y_m - z_{dm})(\rho)\|_1^2 d\rho \right)^{\frac{1}{2}} \leq M(dm)^{1-s} \|u\|_{L^2(0,T;H^s)}$$

$$+ Mm^{-s}\|u\|_{L^2(0,T;H^s)} + Mm^{1-\gamma}\|u_t\|_{L^2(0,T;H^{\gamma-2})}.$$

*Proof.* The existence and uniqueness of $y_m(t)$ and $z_{dm}(t)$ on some interval $[0, T_m]$ with $T_m \leq T$ follow directly from standard results on the Cauchy problem for a system of ordinary differential equations. From the a priori estimate given below, one can conclude that $T_m = T$.

Let us first derive an a priori estimate for $y_m$ and $z_{dm}$. Using (4.2), we find

$$
\begin{aligned}
b(y_m + z_{dm}, & y_m + z_{dm}, y_m) + b(y_m + z_{dm}, y_m, z_{dm}) \\
&= b(y_m + z_{dm}, y_m + z_{dm}, y_m) + b(y_m + z_{dm}, y_m + z_{dm}, z_{dm}) \\
&= b(y_m + z_{dm}, y_m + z_{dm}, y_m + z_{dm}) = 0.
\end{aligned}
$$

Setting $v = y_m$ in (4.6) and $w = z_{dm}$ in (4.7), summing up the two relations and using the above identity, we obtain

$$
\frac{1}{2}\frac{d}{dt}\|y_m\|^2 + \nu\|y_m\|_1^2 + \nu\|z_{dm}\|_1^2 = (f, y_m + z_{dm}) \le C\|f\|^2 + \frac{\nu}{2}(\|y_m\|_1^2 + \|z_{dm}\|_1^2).
$$

Integrating the above inequality from 0 to $t$, since $\|\Pi_m u_0\| \le C\|u_0\|_1$, we derive

$$
(4.9) \quad \|y_m(t)\|^2 + \nu \int_0^t (\|y_m\|_1^2 + \|z_{dm}\|_1^2)ds \le C\left(\|u_0\|_1^2 + \int_0^T \|f\|^2 ds\right) \quad \forall t \in [0, T].
$$

We now turn our attention to the error estimates. From the definition of $\Pi_{dm}$, we find that the solution $u$ satisfies
(4.10)
$$
\frac{d}{dt}(u, v + w) + \nu((\Pi_{dm}u)_x, (v + w)_x) + b(u, u, v + w) = (f, v + w) \quad \forall v \in V_m, w \in W_{dm}.
$$

Set $e = \Pi_{dm}u - (y_m + z_{dm})$, $\xi = \Pi_m u - y_m$, and $\eta = (\Pi_{dm} - \Pi_m)u - z_{dm}$. We note that $e = \xi + \eta$, $\xi \in V_m$, $\eta \in W_{dm}$, and since we have set $\|u\|_1^2 = a(u, u)$, we have, in particular,

$$
(4.11) \qquad\qquad \|e\|_1^2 = \|\xi\|_1^2 + \|\eta\|_1^2.
$$

Subtract (4.6) and (4.7) from (4.10), since $(y_x, z_x) = 0$, for all $y \in V_m, z \in W_{dm}$, we find

$$
\begin{aligned}
(\xi_t, v) + \nu a(e, v + w) =& ((\Pi_m u - u)_t, v) - (u_t, w) \\
& + [b(y_m + z_{dm}, y_m + z_{dm}, v) - b(u, u, v)] \\
& + [b(y_m + z_{dm}, y_m, w) - b(u, u, w)].
\end{aligned}
$$

Now set $v = \xi$ and $w = \eta$ in the above relation; we obtain

$$
\begin{aligned}
(4.12) \qquad \frac{1}{2}\frac{d}{dt}\|\xi\|^2 + \nu\|e\|_1^2 =& [((\Pi_m u - u)_t, \xi) - (u_t, \eta)] \\
& + [b(y_m + z_{dm}, y_m + z_{dm}, \xi) - b(u, u, \xi)] \\
& + [b(y_m + z_{dm}, y_m, \eta) - b(u, u, \eta)] = I_1 + I_2 + I_3.
\end{aligned}
$$

In the following, we will use frequently the Cauchy–Schwarz inequality, the Young inequality, the Poincaré inequality, (4.2), and (4.3) to bound the three terms on the right-hand side of (4.12).

First, since $\eta \in W_{dm}$, we get easily

$$
\begin{aligned}
(4.13) \qquad I_1 =& [((\Pi_m u - u)_t, \xi) - ((I - P_{m-2})u_t, \eta)] \\
\le& \frac{\nu}{8}(\|\xi\|_1^2 + \|\eta\|_1^2) + C\|(\Pi_{dm}u - u)_t\|_{-1}^2 + C\|(I - P_{m-2})u_t\|_{-1}^2 \\
=& \frac{\nu}{8}\|e\|_1^2 + C\|(\Pi_{dm}u - u)_t\|_{-1}^2 + C\|(I - P_{m-2})u_t\|_{-1}^2.
\end{aligned}
$$

We rewrite $I_2$ as follows:

$$\begin{aligned}
I_2 &= b(y_m + z_{dm}, y_m + z_{dm}, \xi) - b(u, u, \xi) \\
&= b(y_m + z_{dm}, y_m + z_{dm}, \xi) - b(y_m + z_{dm}, u, \xi) + b(y_m + z_{dm}, u, \xi) - b(u, u, \xi) \\
&= b(y_m + z_{dm}, y_m + z_{dm} - u, \xi) + b(y_m + z_{dm} - u, u, \xi) \\
&= -b(y_m + z_{dm}, u - \Pi_{dm}u + e, \xi) - b(u - \Pi_{dm}u + e, u, \xi) \\
&= -b(y_m + z_{dm}, u - \Pi_{dm}u, \xi) - b(y_m + z_{dm}, e, \xi) - b(u - \Pi_{dm}u + e, u, \xi) \\
&= I_{21} + I_{22} + I_{23}.
\end{aligned}$$

Similarly,

$$\begin{aligned}
I_3 &= b(y_m + z_{dm}, y_m, \eta) - b(u, u, \eta) \\
&= b(y_m + z_{dm}, y_m - u, \eta) + b(y_m + z_{dm} - u, u, \eta) \\
&= -b(y_m + z_{dm}, u - \Pi_m u + \xi, \eta) - b(u - \Pi_{dm}u + e, u, \eta) \\
&= -b(y_m + z_{dm}, u - \Pi_m u, \eta) - b(y_m + z_{dm}, \xi, \eta) - b(u - \Pi_{dm}u + e, u, \eta) \\
&= I_{31} + I_{32} + I_{33}.
\end{aligned}$$

Thanks to (4.2), we find

$$\begin{aligned}
I_{22} + I_{32} &= -b(y_m + z_{dm}, e, \xi) - b(y_m + z_{dm}, \xi, \eta) \\
&= -b(y_m + z_{dm}, e, \xi) - b(y_m + z_{dm}, e, \eta) = -b(y_m + z_{dm}, e, e) = 0.
\end{aligned}$$

Using (4.3) and the Young inequality, we find

$$\begin{aligned}
I_{23} + I_{33} &= -b(u - \Pi_{dm}u + e, u, e) \\
&\leq C(\|u - \Pi_{dm}u\|_1 + \|e\|_1)\|u\|_1\|e\| \\
&\leq M(\|u - \Pi_{dm}u\|_1\|e\| + \|e\|\|e\|_1) \\
&\leq \frac{\nu}{8}\|e\|_1^2 + M(\|u - \Pi_{dm}u\|_1^2 + \|e\|^2).
\end{aligned}$$

Using (4.2), we find

$$\begin{aligned}
I_{21} + I_{31} &= -b(y_m + z_{dm}, u - \Pi_{dm}u, \xi) - b(y_m + z_{dm}, u - \Pi_m u, \eta) \\
&= -b(y_m + z_{dm}, u - \Pi_{dm}u, e) - b(y_m + z_{dm}, \Pi_{dm}u - \Pi_m u, \eta) \\
&= b(e - \Pi_{dm}u, u - \Pi_{dm}u, e) + b(e - \Pi_{dm}u, \Pi_{dm}u - \Pi_m u, \eta) \\
&= E_1 + E_2.
\end{aligned}$$

We derive from (2.18) that $\|\Pi_m u\|_1 \leq C\|u\|_1$. Then thanks to (4.3) and the Young inequality, we get

$$\begin{aligned}
E_1 &\leq C(\|\Pi_{dm}u\|_1 + \|e\|_1)\|u - \Pi_{dm}u\|_1\|e\| \\
&\leq M\|u - \Pi_{dm}u\|_1\|e\| + C\|u - \Pi_{dm}u\|_1\|e\|\|e\|_1 \\
&\leq M(\|u - \Pi_{dm}u\|_1^2 + \|e\|^2) + \frac{\nu}{8}\|e\|_1^2 + C\|u - \Pi_{dm}u\|_1^2\|e\|^2 \\
&\leq M(\|u - \Pi_{dm}u\|_1^2 + \|e\|^2) + \frac{\nu}{8}\|e\|_1^2,
\end{aligned}$$

where we have used the fact that $\|u - \Pi_{dm}u\|_1 \leq C\|u\|_1 \leq M$.

To estimate $E_2$, we use the following enhanced Poincaré inequality, which will be proved below:

$$\text{(4.14)} \qquad \|z\|_\omega \leq Cm^{-1}\|z\|_{1,\omega} \quad \forall z \in W_{dm}.$$

Using (4.3), (4.14), and the fact that $\|\eta\|_1 \leq \|e\|_1$ (see (4.11)), we find for $m$ sufficiently large,

$$
\begin{aligned}
E_2 &\leq C(\|\Pi_{dm}u\|_1 + \|e\|_1)\|\Pi_m u - \Pi_{dm}u\|_1\|\eta\| \\
&\leq Cm^{-1}(\|\Pi_{dm}u\|_1 + \|e\|_1)\|\Pi_m u - \Pi_{dm}u\|_1\|\eta\|_1 \\
&\leq Mm^{-1}\|\Pi_m u - \Pi_{dm}u\|_1\|e\|_1 + Cm^{-1}\|\Pi_m u - \Pi_{dm}u\|_1\|e\|_1^2 \\
&\leq \frac{\nu}{16}\|e\|_1^2 + Mm^{-2}\|\Pi_m u - \Pi_{dm}u\|_1^2 + Mm^{-1}\|e\|_1^2 \\
&\leq \frac{\nu}{8}\|e\|_1^2 + Mm^{-2}\|\Pi_m u - \Pi_{dm}u\|_1^2.
\end{aligned}
$$

Therefore, combining the above inequalities, we derive

$$
\begin{aligned}
(4.15) \qquad \frac{d}{dt}\|\xi\|^2 + \nu\|e\|_1^2 &\leq M\|e\|^2 + Mm^{-2}\|\Pi_m u - \Pi_{dm}u\|_1^2 \\
&\quad + C(\|u - \Pi_{dm}u\|_1^2 + \|(u - \Pi_m)u_t\|_{-1}^2 + \|(I - P_{m-2})u_t\|_{-1}^2).
\end{aligned}
$$

By virtue of the enhanced Poincaré inequality (4.14), we find

$$
\begin{aligned}
\|e\|^2 &= \|\xi + \eta\|^2 \leq 2\|\xi\|^2 + 2\|\eta\|^2 \\
&\leq 2\|\xi\|^2 + Cm^{-2}\|\eta\|_1^2 \leq 2\|\xi\|^2 + Cm^{-2}\|e\|_1^2.
\end{aligned}
$$

Hence for $m$ sufficiently large, we can rewrite (4.15) as

$$
\begin{aligned}
\frac{d}{dt}\|\xi\|^2 + \nu\|e\|_1^2 &\leq M\|\xi\|^2 + Mm^{-2}\|\Pi_m u - \Pi_{dm}u\|_1^2 \\
&\quad + C(\|u - \Pi_{dm}u\|_1^2 + \|(u - \Pi_m u)_t\|_{-1}^2 + \|(I - P_{m-2})u_t\|_{-1}^2).
\end{aligned}
$$

Applying the Gronwall lemma to the above inequality and using Lemma 2.1, (3.22), and noticing that $\xi(0) = 0$, we derive that for all $t \in [0, T]$, we have

$$
\begin{aligned}
\|\xi(t)\|^2 + \int_0^t \|e(\rho)\|_1^2 d\rho &\leq Mm^{-2}\int_0^t (\|\Pi_m u - u\|_1^2 + \|u - \Pi_{dm}u\|_1^2)d\rho \\
&\quad + M\int_0^t (\|u - \Pi_{dm}u\|_1^2 + \|(u - \Pi_m u)_t\|_{-1}^2 + \|(I - P_{m-2})u_t\|_{-1}^2)d\rho \\
&\leq M(m^{-2s} + (dm)^{2-2s})\int_0^T \|u\|_s^2 d\rho + Mm^{2-2\gamma}\int_0^T \|u_t\|_{\gamma-2}^2 d\rho.
\end{aligned}
$$

We then conclude from the above inequality and Lemma 2.2 that

$$
\begin{aligned}
\int_0^T \|(u - y_m - z_{dm})(\rho)\|_1^2 d\rho &\leq \int_0^T \|u - \Pi_{dm}u\|_1^2 d\rho \\
&\quad + M(m^{-2s} + (dm)^{2-2s})\int_0^T \|u\|_s^2 d\rho + Mm^{2-2\gamma}\int_0^T \|u_t\|_{\gamma-2}^2 d\rho \\
&\leq M(m^{-2s} + (dm)^{2-2s})\int_0^T \|u\|_s^2 d\rho + Mm^{2-2\gamma}\int_0^T \|u_t\|_{\gamma-2}^2 d\rho.
\end{aligned}
$$

It remains only to prove (4.14). For any $z \in W_{dm}$, let $w$ be the solution of the Poisson equation $-w_{xx} = z$, $w(\pm 1) = 0$. Then by using (2.4) and (2.5), we find that

$$
\begin{aligned}
\|z\|^2 &= (z, -w_{xx}) = \inf_{v \in V_m}(z, -(w - v)_{xx}) \\
&\leq \alpha\|z\|_1 \inf_{v \in V_m}\|w - v\|_1 \leq C\|z\|_1 m^{-1}\|w\|_2 \leq C\|z\|_1 m^{-1}\|z\|. \qquad \square
\end{aligned}
$$

*Remark* 4.2. (a) Error estimates for the classical Galerkin scheme can be easily recovered from the above process. In fact, taking $d = 1$, we have $W_{dm} = \{0\}$ and therefore $z_{dm} = 0$. In this case, (4.6)–(4.7) reduce to the classical Galerkin scheme in $V_m$, and $y_m$ becomes the solution of the classical Galerkin scheme in $V_m$. Repeating the proof of Theorem 4.1 with $z_{dm} = 0$, we can obtain the following error estimates for the solution $u_{dm}$ of (4.5):

$$(4.16) \qquad \left( \int_0^T \|(u - u_{dm})(\rho)\|_1^2 d\rho \right)^{\frac{1}{2}} \leq M(dm)^{1-s} \|u\|_{L^2(0,T;H^s)}$$
$$+ M(dm)^{1-\gamma} \|u_t\|_{L^2(0,T;H^{\gamma-2})},$$

and

$$(4.17) \qquad \|u(t) - u_{dm}(t)\| \leq M(dm)^{1-s} \left\{ \|u\|_{L^\infty(0,T;H^{s-1})} + \|u\|_{L^2(0,T;H^s)} \right\}$$
$$+ M(dm)^{1-\gamma} \|u_t\|_{L^2(0,T;H^{\gamma-2})} \quad \forall t \in [0,T].$$

(b) Comparing (4.8) with (4.16), it is clear that the first term of the error estimate in (4.8) is inherited from the classical Galerkin approximation, while the second term of the error estimate is due to our nonlinear Galerkin treatment. However, since we assume $f$ is independent of time, we can show (see [9]) that $u$ is analytic in time with value in $V$. Using Cauchy's formula for analytic functions, we can prove $\|u_t\|_s \sim \|u\|_s$. Therefore by taking $\gamma = s + 2$ in (4.8), we find that as long as $d^{s-1} \lesssim m$, (4.6)–(4.7) is as accurate as the classical Galerkin scheme (4.5) in $V_{dm}$. Note that (4.16)–(4.17) as well as the results in Theorem 4.1 are only valid for finite time intervals since the constant $M$ may depend (exponentially) on $T$. It would be nice if this last result could be extended in some way to unbounded time intervals.

(c) Error estimates in $L^\infty(0,T;H^1(\Omega))$ norm can also be derived. We shall not pursue this direction since the main purpose of our error analysis is to determine a quantitative guideline for the proper choice of $d$.

We note that recently Marion and Xu [15] has derived similar error estimates for the nonlinear Galerkin method in the context of two-grid finite elements.

**4.2. A second scheme.** For a general function $f \in L^2(0,T;L^2(I))$, the solution $u$ will not be analytic in time, and therefore the term $dz_{dm}/dt$ may not be negligible by comparison with $\nu(z_{dm})_{xx}$ and hence cannot be neglected in the approximation. Consequently, we propose to approximate (4.4) by the following multilevel scheme of nonlinear Galerkin type.

Find $y_m \in V_m$ and $z_{dm} \in W_{dm}$ such that

$$(4.18) \quad \frac{d}{dt}(y_m + z_{dm}, v) + \nu a(y_m, v) + b(y_m + z_{dm}, y_m + z_{dm}) = (f, v) \quad \forall v \in V_m,$$

$$(4.19) \quad \frac{d}{dt}(y_m + z_{dm}, v) + \nu a(z_{dm}, w) + b(y_m + z_{dm}, y_m, v) = (f, w) \quad \forall w \in W_{dm},$$

with $y_m(0) = \Pi_m u_0$ and $z_{dm}(0) = (\Pi_{dm} - \Pi_m)u_0$.

The following results can be established for the above scheme.

THEOREM 4.2. *Under assumptions* (4.2) *and* (4.3), *we have for m sufficiently*

*large,*

$$\left( \int_0^T \|(u - y_m - z_{dm})(\rho)\|_1^2 d\rho \right)^{\frac{1}{2}} \leq M(dm)^{1-s} \|u\|_{L^2(0,T;H^s)}$$

$$+ Mm^{-s} \|u\|_{L^2(0,T;H^s)} + M(dm)^{1-\gamma} \|u_t\|_{L^2(0,T;H^{\gamma-2})};$$

$$\|u(t) - y_m(t) - z_{dm}(t)\| \leq M(dm)^{1-s} \left\{ \|u\|_{L^\infty(0,T;H^{s-1})} + \|u\|_{L^2(0,T;H^s)} \right\}$$

$$+ Mm^{-s} \|u\|_{L^2(0,T;H^s)} + M(dm)^{1-\gamma} \|u_t\|_{L^2(0,T;H^{\gamma-2})} \quad \forall t \in [0, T].$$

The proof of this result is very similar to, and in fact a little easier than, that of Theorem 4.1. Hence we leave it to the interested readers.

*Remark* 4.3. (a) The linear part of the system (4.18)–(4.19) is not totally decoupled but it is indeed quasidecoupled since there are only two modes of coupling in each direction.

(b) Comparing the result in Theorem 4.2 with (4.16)–(4.17), we find that again as long as $d^{s-1} \lesssim m$, (4.18)–(4.19) is as accurate as the classical Galerkin scheme in $V_{dm}$.

*Remark* 4.4. (a) It has been shown (cf. [19], [21]) that the stability conditions for time discretized nonlinear Galerkin schemes of implicit-explicit or explicit-type only depend on the number of low modes, while the stability conditions of accordingly time discretized classical Galerkin schemes depend on the number of total (low and high) modes. Hence the larger $d$ is, the better the stability condition of the nonlinear Galerkin scheme becomes.

(b) It is now transparent that the stronger the nonlinearity is, the more restrictive the condition on $d$ becomes. The error estimates in Theorems 4.1 and 4.2 and the condition $d^{s-1} \lesssim m$ are optimal with respect to assumption (4.3) on the nonlinear term. However for equations with weaker nonlinearities, such as the reaction-diffusion equation and the Kuramoto–Sivashinsky equation, the condition can be relaxed accordingly to $d^{s-1} \lesssim m^\alpha$ for some $\alpha \in (1, 2]$.

*Remark* 4.5. The condition $d^{s-1} \lesssim m$ should be understood properly. At first sight, it seems that for highly smooth functions ($s >> 1$) we are forced to choose $d \approx 1$, which means no gain for the nonlinear Galerkin approach. However we should keep in mind that the spectral accuracy can only be achieved when the structure of the solution is fully resolved, otherwise the convergence rate of the spectral method is only algebraic in $m^{-1}$, even if the solution is infinitely differentiable. Thus for intermediate realistic values of $m$, the proposed algorithm with the appropriate value of $d$ as indicated above may be most efficient. This remark applies as well for the previous linear schemes since these schemes also produce a coupling between low and high modes.

## REFERENCES

[1] C. CANUTO, M. Y. HUSSAINI, A. QUARTERONI, AND T. A. ZANG, *Spectral Methods in Fluid Dynamics*, Springer-Verlag, New York, Heidelberg, Berlin, 1987.

[2] C. CANUTO AND A. QUARTERONI, *Spectral and pseudo-spectral methods for parabolic problems with nonperiodic boundary conditions*, Calcolo, 18 (1981), pp. 197–218.

[3] C. CANUTO AND A. QUARTERONI, *Approximation results for orthogonal polynomials in Sobolev spaces*, Math. Comp., 38 (1982), pp. 67–86.

[4] P. CONSTANTIN, C. FOIAS, B. NICOLAENKO, AND R. TEMAM, *Integral manifolds and inertial manifolds for dissipative partial differential equations*, Appl. Math. Sci. Series, Vol. 70,

New York, Heidelberg, Berlin, 1988.

[5] C. DEVULDER, M. MARION, AND E. TITI, *On the convergence rate of the nonlinear Galerkin methods*, Math. Comp., 60 (1993), pp. 495–514.

[6] T. DUBOIS, F. JAUBERTEAU, AND R. TEMAM, *Solution of the incompressible Navier–Stokes equations by the nonlinear Galerkin method*, J. Sci. Comput., 8 (1993), pp. 167–194.

[7] C. FOIAS, G. R. SELL, AND R. TEMAM, *Variétés inertielles des équations différentielles dissipative*, C. R. Acad. Sci. Paris Ser. I, 301 (1985), pp. 139–141.

[8] ———, *Inertial manifolds for nonlinear evolutionary equations*, J. Differential Equations, 73 (1988), pp. 308–353.

[9] C. FOIAS AND R. TEMAM, *Some analytic and geometric properties of the evolution Navier–Stokes equations*, J. Math. Pures Appl., 58 (1979), pp. 339–368.

[10] D. GOTTLIEB AND S. A. ORSZAG, *Numerical Analysis of Spectral Methods: Theory and Applications*, Society for Industrial and Applied Mathematics, Philadelphia, PA, 1977.

[11] F. JAUBERTEAU, C. ROSIER, AND R. TEMAM, *The nonlinear Galerkin method in computational fluid dynamics*, Appl. Numer. Math., 6 (1989/90), pp. 361–370.

[12] M. S. JOLLY, I. G. KEVREKIDIS, AND E. S. TITI, *Approximate inertial manifolds for the Kuramoto–Sivashinsky equation: Analysis and computations*, Phys. D, 44 (1990), pp. 38–60.

[13] Y. MADAY AND A. QUARTERONI, *Legendre and Chebyshev spectral approximation of Burgers' equation*, Numer. Math., 37 (1981), pp. 321–332.

[14] M. MARION AND R. TEMAM, *Nonlinear Galerkin methods*, SIAM J. Numer. Anal., 26 (1989), pp. 1139–1157.

[15] M. MARION AND J. C. XU, *Error estimates on a new nonlinear Galerkin method based on two-grid finite elements*, SIAM J. Numer. Anal., to appear.

[16] G. SACCHI-LANDRIANI, *Spectral tau approximation of the two-dimensional Stokes problem*, Numer. Math., 23 (1988), pp. 683–699.

[17] J. SHEN, *Efficient spectral-Galerkin method I. Direct solvers for second and fourth order equations by using Legendre polynomials*, SIAM J. Sci. Comput., 15 (1994), pp. 1489–1505; *Efficient spectral-Galerkin methods II. Direct solvers for second and fourth order equation by using Chebyshev polynomials*, SIAM J. Sci. Comput. 16 (1995), pp. 74–87.

[18] ———, *A spectral-tau approximation for the Stokes and Navier–Stokes equations*, Math. Model. Numer. Anal., 22 (1988), pp. 677–693.

[19] ———, *Long time stability and convergence for fully discrete nonlinear Galerkin methods*, Appl. Anal., 38 (1990), pp. 201–229.

[20] R. TEMAM, *Infinite-Dimensional Dynamical Systems in Mechanics and Physics*, Appl. Math. Sci. Series, Vol. 68, New York, Heidelberg, Berlin, 1988.

[21] ———, *Stability analysis of the nonlinear Galerkin method*, Math. Comp., 57 (1991), pp. 477–505.

[22] E. S. TITI, *On approximate inertial manifolds to the 2-D Navier–Stokes equations*, Math. Anal. Appl., 149 (1990), pp. 540–557.