

## AN EFFICIENT SPECTRAL METHOD FOR ORDINARY DIFFERENTIAL EQUATIONS WITH RATIONAL FUNCTION COEFFICIENTS

EVANGELOS A. COUTSIAS, THOMAS HAGSTROM, AND DAVID TORRES

ABSTRACT. We present some relations that allow the efficient approximate inversion of linear differential operators with rational function coefficients. We employ expansions in terms of a large class of orthogonal polynomial families, including all the classical orthogonal polynomials. These families obey a simple 3-term recurrence relation for differentiation, which implies that on an appropriately restricted domain the differentiation operator has a unique banded inverse. The inverse is an integration operator for the family, and it is simply the tridiagonal coefficient matrix for the recurrence. Since in these families convolution operators (i.e., matrix representations of multiplication by a function) are banded for polynomials, we are able to obtain a banded representation for linear differential operators with rational coefficients. This leads to a method of solution of initial or boundary value problems that, besides having an operation count that scales linearly with the order of truncation  $N$ , is computationally well conditioned. Among the applications considered is the use of rational maps for the resolution of sharp interior layers.

### 1. INTRODUCTION

The solution of constant-coefficient ordinary differential equations with periodic boundary conditions is especially simple in the Fourier spectral representation, since differentiation of a smooth function is replaced by multiplication of its Fourier coefficient vector by a diagonal matrix. An analogous property is shared by Hermite polynomial expansions in unbounded domains. Other spectral representations give, in general, almost full triangular differentiation matrices. However, for polynomial families such as the Chebyshev and Legendre, the matrices representing some commonly occurring operators, such as the Laplace operator in various separable geometries, are known to be reducible to simple, banded form through the use of appropriate banded preconditioners ([12, Ch. 10], [9, 18]). The origin of most of such simplifications is found in the fact that the matrix operator for integration in any of the classical orthogonal polynomial families is tridiagonal [8].

---

Received by the editor August 9, 1994 and, in revised form, February 12, 1995.

1991 *Mathematics Subject Classification*. Primary 65Q05, 65L60, 65P05, 76M25, 33A45, 33C55, 33C45.

*Key words and phrases*. Spectral methods, orthogonal polynomials, boundary value problems.

Part of the work of the first author was performed at Risø National Laboratory, DK-4000 Roskilde, Denmark. All authors supported in part by DOE Grant DE-FG03-92ER25128.

The work of the second author was partially supported by NSF Grants DMS-9108072, DMS-9304406 and by ICOMP, NASA Lewis Res. Ctr., Cleveland, OH, USA.

In this article we show how to exploit the properties of the operator of integration for arbitrary classical orthogonal polynomial families to arrive at efficient spectral algorithms for the approximate solution of a large class of ordinary differential equations of the form

$$(1) \quad Lu = \sum_{k=0}^n (m_{n-k}(x)D^k)u = f(x) \quad , \quad x \in \Omega = (a, b),$$

subject to the constraints

$$\mathcal{T}u = c \quad ,$$

where  $m_k$  are rational functions of  $x$ ,  $D^k$  denotes  $k$ th-order differentiation with respect to  $x$ ,  $\mathcal{T}$  is a linear functional of rank  $n$ , and  $c \in R_n$ . (Typically, the constraints are boundary or initial conditions, but this is not necessary.)

We must mention that the basic idea of the method presented here was first introduced by Clenshaw [6]. He realized that solving for the highest derivative present in a given ordinary differential equation leads to banded forms for Chebyshev Galerkin discretizations for ODEs with low-order polynomial coefficients, which then he solved by backward recurrence relations. The method is further discussed in the monograph by Fox and Parker [11], again for the Chebyshev polynomials. Among our main contributions are the development of an efficiently implementable algorithm for general, nonsingular problems in arbitrary classical orthogonal polynomial bases, together with its conditioning and convergence analysis, and the application to the resolution of sharp layers through rational maps. We present now the basic description of our method, followed by an outline of the rest of the paper.

The problem of approximating solutions of Ordinary or Partial Differential Equations (O or PDE) by spectral methods, known as Galerkin approximation, involves the projection onto the span of some appropriate set of basis functions, typically arising as the eigenfunctions of a singular Sturm-Liouville (SL) problem. The members of the basis may satisfy automatically the auxiliary conditions imposed on the problem, such as initial, boundary or more general conditions. Alternatively, these conditions may be imposed as constraints on the expansion coefficients, as in the Lanczos  $\tau$ -method [15].

It is well known [5] that the eigenfunctions of certain singular Sturm-Liouville problems allow the approximation of functions in  $C^\infty[a, b]$  whose truncation error approaches zero faster than any negative power of the number of basis functions (modes) used in the approximation, as that number (order of truncation  $N$ ) tends to  $\infty$ . This phenomenon is usually referred to as ‘spectral accuracy’ [12]. The accuracy of derivatives obtained by direct, term-by-term differentiation of such truncated expansions naturally deteriorates [5], but for low-order derivatives and sufficiently high-order truncations this deterioration is negligible, compared to the restrictions in accuracy introduced by typical difference approximations. Since results on the accuracy of spectral methods are well documented in the literature, we shall limit ourselves to the discussion of certain formal properties of orthogonal polynomial families, which allow algorithmic simplifications in their use. Facts about orthogonal polynomials that we shall need can be found in any of the standard references (e.g. [16, 19]).

Throughout, we assume that we are working with a family of polynomials  $\{Q_k\}_0^\infty$  which are orthogonal and complete over the interval  $(a, b)$  (here  $a$  and/or  $b$  can be infinite) with respect to the nonnegative weight  $w(x)$ . In the cases of interest, these

are the eigenfunctions of a Sturm-Liouville problem

$$(2) \quad (p(x)Q'_k)' + \lambda_k w(x)Q_k = 0.$$

Then the  $Q'_k$  form an orthogonal family as well, with nonnegative weight  $p(x)$  which satisfies  $p(x) \rightarrow 0$  as  $x \rightarrow a, b$ . In this paper we focus exclusively on the classical orthogonal polynomials, i.e., the Jacobi (special cases of which are the Chebyshev, Legendre and Gegenbauer polynomials), Laguerre and Hermite polynomials, which are the only polynomial solutions of Sturm-Liouville problems of the form (2) [14]. We will assume that the functions under consideration possess sufficient differentiability properties over  $(a, b)$  and can be expressed as a series involving the  $Q_k$ . See [5] for a discussion of the convergence properties in the relevant function spaces.

We introduce the spaces  $Q_m^n$  by

$$Q_m^n \equiv \text{span}\{Q_k \mid m \leq k \leq n\}.$$

Our method constructs an approximate particular solution of (1) in a subspace of codimension  $n$  (e.g.  $(Q_0^{n-1})^\perp$ ) such that when  $n$ th-order differentiation is restricted to this subspace it has a simple inverse. We also require that  $L$  be invertible when restricted to this subspace and that  $\mathcal{T}$  has full rank when restricted to the space of solutions to the homogeneous problem ((1) with  $f = 0$ ).

Of key importance for our purposes is the requirement that differentiation or its inverse ('integration' in an appropriately restricted domain) must have banded form. For example, the first derivative operator in the Chebyshev representation,  $D$  has elements

$$\frac{1}{2}D_{i,j} = \begin{cases} 0, & i \geq j, \\ 0, & i < j, \quad i+j \text{ even}, \\ j, & 0 < i < j, \quad i+j \text{ odd}, \\ \frac{j}{2}, & i = 0, \quad j \text{ odd}. \end{cases}$$

Its inverse, when respective domains and ranges are appropriately restricted, is given by

$$B = \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & \cdots & 0 \\ 2 & 0 & -1 & \cdots & \cdots & \cdots & 0 \\ 0 & 1/2 & 0 & -1/2 & \cdots & \cdots & 0 \\ 0 & 0 & \ddots & \ddots & \ddots & \cdots & 0 \\ 0 & 0 & 0 & 1/k & 0 & -1/k & \cdots \end{pmatrix}.$$

Now,  $DB = I_{Q_0^\infty}$  while  $BD = I_{Q_1^\infty}$ . Clearly,  $D^k B^k = I_{Q_0^\infty}$  as well. However,  $B^k D^k \neq I$ . If we apply  $k$ -fold differentiation to an arbitrary function, all information about the first  $k$  coefficients in its Chebyshev expansion is lost. If however we restrict the action of  $D^k$  to the space  $Q_k^\infty$ , then  $B^k$  is a left inverse provided its range is restricted to the same space. We introduce the notation  $A_{[k]}$  to denote a matrix  $A$  with its first  $k$  rows set to zero. Thus, we have that  $B_{[k]}^k D^k = I_{Q_k^\infty}$ . We note that these relationships carry over to finite truncations if we replace the last column of  $B$  and the last  $k$  columns of  $B_{[k]}^k$  with zeros, since  $D^k : Q_k^N \rightarrow Q_0^{N-k}$  while  $B_{[k]}^k : Q_0^{N-k} \rightarrow Q_k^N$ . It is easy to see that these simple inversion (integration) operators originate in the recursions

$$(3) \quad \frac{T'_{k+1}}{k+1} - \frac{T'_{k-1}}{k-1} = 2T_k \quad , \quad k = 1, \dots, \\ T'_0 = 0 \quad , \quad T'_1 = T_0,$$

$$(4) \quad \frac{T''_{k+2}}{4(k+2)(k+1)} - \frac{2T''_k}{4(k^2-1)} + \frac{T''_{k-2}}{4(k-1)(k-2)} = T_k \quad , \quad k = 1, \dots ,$$

$$T''_0 = 0 \quad , \quad T''_1 = 0 \quad , \quad T''_2 = 4T_0,$$

and so on for higher derivatives. Clearly,  $B$  and  $B_{[2]}^2$  are the matrices of recursion coefficients for equations (3), (4), respectively. In the discussion we use the same symbol for an infinite-dimensional matrix operator and its finite-dimensional truncation, where the distinction is clear from the context.

More generally, if  $\{Q_k(x)\}_0^\infty$  is a family of orthogonal polynomials, then a three-term recurrence for multiplication by the monomial  $x$

$$(5) \quad \sum_{l=-1}^1 Q_{k+l} a_{k+l,k} = xQ_k \quad , \quad k = 0, 1, \dots$$

follows easily from the orthogonality of the  $Q_k$  [19]. Since the  $Q'_k$  are orthogonal (with weight  $p(x)$ , as is easily seen by integrating (2) by parts), they also satisfy a relation of form (5):

$$(6) \quad \sum_{l=-1}^1 Q'_{k+l} a_{k+l,k}^{(1)} = xQ'_k \quad , \quad k = 0, 1, \dots$$

Therefore, by differentiating (5) and combining with (6), we arrive at [8]

$$(7) \quad \sum_{l=-1}^1 Q'_{k+l} b_{k+l,k} = Q_k \quad , \quad k = 0, 1, \dots$$

which allows the efficient inversion of differentiation to all orders. The coefficients in (7) can be derived from those of the basic recurrence (5), which defines the family.

The method we shall present in §3, explained in detail for 2nd-order operators but not limited to them, relies on restricting the domain of  $D^n$  to the subspace  $Q_n^N = \text{span}\{Q_k\}_n^N$ , thus ensuring the existence of a unique inverse. Throughout, we tacitly assume that the operator  $L_N$ , the  $N$ th-order Galerkin approximation to  $L$ , has rank  $N-n$  when acting on elements of  $Q_0^N$ . Thus, the problem of solving the resulting algebraic system for right-hand sides restricted to  $Q_0^{N-n}$  has a solution containing  $n$  free parameters. We moreover assume that the operator is nonsingular when restricted further to  $Q_n^N$ . Thus, the null space contains no element orthogonal to  $Q_0^{n-1}$ . These assumptions are not as restrictive as one might at first expect. The method is most effective when the above problem needs to be solved repeatedly for several right-hand sides  $f$  and high accuracy is desired. This type of problem arises, e.g., when the Navier-Stokes equations are solved in a geometry in which the Laplace operator is separable, and the boundary conditions are periodic in all directions except one. Common examples are provided by the Laplace operator in various separable curvilinear coordinates, where expansions of smooth functions in terms of eigenfunctions of the Laplacian in the bounded direction do not possess good convergence properties.

In §3 we give some examples of the inversion of the Laplacian in some common geometries, including a disk and an annulus in cylindrical and helical coordinate systems. The use of the method for initial value problems is demonstrated through a study of the Airy equation, while the biharmonic equation, analyzed in §4, provides an example for a higher-order problem. Also considered is the Stokes problem:

here a coupled system of two second-order equations is studied with boundary conditions given for only one. The method is easily extended to cover this case. The Chebyshev polynomials are an especially important family, because of their optimal approximation properties as well as the applicability of the Fast Fourier Transform. Thus, most of our explicit calculations are carried out for Chebyshev-Galerkin matrices. In §4 we carry out a detailed conditioning analysis for typical problems. It is found that if the leading coefficient  $m_0(x)$  does not vanish in the interval under consideration, the method generically produces well-conditioned operators. Finally, in §5 we discuss how to use rational mappings to stretch the coordinate system near points where the solution of a BVP exhibits rapid variation, thus ensuring a more efficient representation of the solution without sacrificing the speed of the method.

2. RECURSIVE DETERMINATION OF DERIVATIVES

Throughout, we assume that  $\{Q_k(x)\}_0^\infty$  is a family of orthogonal polynomials in  $[a, b]$  with weight  $w(x)$ , such that if  $u \in C^\infty[a, b]$  and if we set

$$u_N = \sum_0^N \hat{u}_k^0 Q_k$$

with

$$\hat{u}_k^0 = \frac{1}{h_k} \int_a^b u(x) Q_k(x) w(x) dx \quad , \quad \text{where } h_k = \| Q_k \|_w^2 \quad ,$$

then the error  $\| u - u_N \|_w \xrightarrow{N \rightarrow \infty} 0$  faster than any negative power of  $N$ . This is for example true for the eigenfunctions of certain singular Sturm–Liouville problems [5].

We shall write  $D_N^n$  for the restriction of the  $n$ th-derivative operator with respect to  $x$  on  $Q_0^N = \text{span} \{Q_k\}_0^N$ . We adopt the notation

$$(8) \quad D^n u = D^n \sum_0^\infty \hat{u}_k^0 Q_k \equiv \sum_0^\infty \hat{u}_k^n Q_k,$$

and we write  $\hat{u}_N = \text{col}(\hat{u}_i) \in R^{N+1}$  ( $i = 0, 1, \dots, N$ ). In the sequel we will drop the subscript  $N$  when the distinction between truncated and nontruncated expansions is clear. Also, as stated earlier, we shall write  $A_{[k]}$  for a matrix  $A$  whose first  $k$  rows have been set equal to zero.

We now prove the following theorem, which is a special case of Theorem 2.2, but because of its simplicity serves to explain ideas. In this form, the theorem applies, e.g., to the Legendre polynomials.

**Theorem 2.1.** *If the family  $\{Q_k(x)\}_0^\infty$  satisfies the recurrence*

$$(9) \quad Q'_{k+1} - Q'_{k-1} = f(k)Q_k \quad , \quad k = 0, 1, \dots \quad ,$$

with  $Q_{-1} \equiv 0$ , then

$$(10) \quad \frac{\hat{u}_{k+1}^1}{f(k+1)} - \frac{\hat{u}_{k-1}^1}{f(k-1)} = -\hat{u}_k^0 \quad , \quad k = 1, 2, \dots \quad .$$

*Proof.* Clearly,

$$Q'_{k+1}(x) = \sum_{\substack{m=0 \\ m+k \text{ even}}}^k f(m)Q_m(x) \quad ,$$

so that

$$\begin{aligned} u' &= \sum_{k=0}^{\infty} \hat{u}_k^0 Q'_k \equiv \sum_{k=0}^{\infty} \hat{u}_k^1 Q_k \\ &= \sum_{k=0}^{\infty} \hat{u}_k^0 \sum_{\substack{m=0 \\ m+k \text{ odd}}}^{k-1} f(m) Q_m \\ &= \sum_{m=0}^{\infty} \left\{ f(m) \sum_{\substack{k=m+1 \\ m+k \text{ odd}}}^{\infty} \hat{u}_k^0 \right\} Q_m , \end{aligned}$$

and finally

$$\hat{u}_m^1 = f(m) \sum_{\substack{k=m+1 \\ k+m \text{ odd}}}^{\infty} \hat{u}_k^0 ,$$

resulting in the recurrence claimed above. □

Applying the formula of Theorem 2.1 repeatedly, we can derive similar recursions for the inversion of higher derivatives as well. For example, for  $D^2$  we have

$$(11) \quad \frac{\hat{u}_{k+2}^2}{f(k+1)f(k+2)} - \hat{u}_k^2 \frac{(f(k+1) + f(k-1))}{f(k+1)f(k)f(k-1)} + \frac{\hat{u}_{k-2}^2}{f(k-1)f(k-2)} = \hat{u}_k^0 , \quad k = 2, 3, \dots$$

The above formulae lead to simple algorithms for the computation of derivatives of functions expanded in terms of the  $Q$ 's as well as for the solution of simple Initial (I) or Boundary Value Problems (BVPs). For example, the solution to the problem

$$u_x = g(x) \quad , \quad u(a) = \alpha \quad ,$$

where

$$g(x) = \sum_{m=0}^{\infty} \hat{g}_m Q_m(x) ,$$

can be found in the form

$$\hat{u}_k^0 = \frac{\hat{g}_{k-1}}{f(k-1)} - \frac{\hat{g}_{k+1}}{f(k+1)} \quad , \quad k = 1, 2, \dots$$

while

$$\hat{u}_0^0 = \left( \alpha - \sum_{m=1}^{\infty} \hat{u}_m^0 Q_m(a) \right) / Q_0(a) .$$

Other simple linear BVPs of the form

$$Lu = g \quad , \quad Bu = l \quad ,$$

can be solved efficiently by the inversion of banded matrices if the differential operator  $L$  has constant coefficients. For example, let

$$(12) \quad L = -\frac{d^2}{dx^2} + \lambda^2 .$$

In order to solve the BVP (12) with boundary conditions

$$(13) \quad u(a) = \alpha \quad , \quad u(b) = \beta$$

numerically, by assuming a truncated expansion for  $u(x)$  of order  $M$ , we set

$$\hat{u}_k^2 = \lambda^2 \hat{u}_k^0 - \hat{g}_k$$

in Eq. (11) for  $k = 2, \dots, M$  to get, together with the  $\tau$ -conditions

$$\begin{aligned} \sum_{m=0}^M \hat{u}_m^0 Q_m(a) &= \alpha \quad , \\ \sum_{m=0}^M \hat{u}_m^0 Q_m(b) &= \beta \quad , \end{aligned}$$

an almost pentadiagonal system (except for the first two rows, which are full) for the coefficients  $\hat{u}_m^0$ . This can be easily solved by LU decomposition. Thus the  $\tau$ -conditions are viewed as the first two equations, followed by the first  $M - 1$  recurrence relations for the determination of the  $\hat{u}_k^0$ ,  $k = 2, \dots, M$ , with  $\hat{u}_k^0 = \hat{g}_k = 0$ ,  $k > M$ . This is equivalent to the usual way of stating the  $\tau$ -method [12].

An alternative approach is suggested here. We specifically look for null vectors in the form  $e_k = Q_k + u_k$ ,  $u_k \in (Q_0^{n-1})^\perp$ ,  $k = 1, \dots, n - 1$ . Then, if  $u_p \in (Q_0^{n-1})^\perp$  is a particular solution, the solution to the BVP can be written as  $u = u_p + \sum \alpha_k e_k$ , with  $\alpha_k$  satisfying an  $n \times n$  system. Note that if repeated solution of the system is required with different right-hand sides, the  $u_k$  need only be determined once, and there is a slight reduction in computational overhead of our method when compared, e.g., with an efficient implementation of the  $\tau$ -method. In fact, our method can effectively optimize the conditioning of a problem by restricting attention to the most stable subspace. So, for example, if one is required to solve the Poisson equation  $\Delta u = -g$  in a region  $\Omega$ , where  $\Omega$  is a 2-(3-)dimensional rectangle with one (two) periodic directions and one bounded direction several times by adopting a Fourier-(Fourier)- $Q$  expansion, the problem decomposes to equations of type (12) in the bounded direction for each Fourier mode. The LU decomposition can be performed in a preprocessing stage and the results stored, resulting in only  $\approx (10M)$  operations per solution per Fourier mode at all subsequent stages. The cost is thus comparable to solving the Poisson equation in the pure Fourier case! Similar results can be easily derived for other ordinary differential operators with constant coefficients.

A straightforward generalization of the previous formulas, which is useful in deriving properties for the Chebyshev polynomials, follows [8]:

**Theorem 2.2.** *If the family  $\{Q_k(x)\}_0^\infty$  satisfies the recurrence*

$$(14) \quad \sum_{l=-1}^1 Q'_{k+l} b_{k+l,k} = Q_k \quad , \quad k = 0, 1, \dots \quad ,$$

with  $Q_{-1} \equiv 0$ , then, if  $f(x) = \sum_{k=0}^\infty \hat{f}_k Q_k(x)$  is a sufficiently differentiable function, there holds

$$(15) \quad \hat{f}_k^{(0)} = \sum_{l=-1}^1 b_{k,k+l} \hat{f}_{k+l}^{(1)} \quad , \quad k = 1, 2, \dots \quad ,$$

where the  $l$ th derivative of the function  $f(x)$  has expansion coefficients  $\hat{f}_k^{(l)}$ .

Thus, even in this more general case, the expansion coefficients of a function can be calculated from those of its derivative in  $O(N)$  operations. The more general form in Theorem 2.2 will be useful dealing with Chebyshev polynomials, for which it agrees with the usual normalizations. The proof is straightforward, but we give it for completeness.

*Proof of Theorem 2.2.* We can introduce the vectors

$$\hat{f}^{(l)} = \left( \hat{f}_0^{(l)}, \hat{f}_1^{(l)}, \dots \right)^T$$

and

$$q_x = (Q_0, Q_1, \dots), \quad q'_x = (Q'_0, Q'_1, \dots).$$

Then,  $f(x) = q_x \hat{f}$  and  $\hat{f}' = q'_x \hat{f}$ . Also, by assumption,  $q'_x B = q_x$ , where  $B$  is the coefficient matrix for recurrence (14). Combining, we find

$$q'_x \left( \hat{f} - B \hat{f}^{(1)} \right) = 0.$$

Assuming that the  $Q_i^{(k)}, i = 1, 2, \dots$ , are independent (true for all families that satisfy Eq.(5)), we find the relation claimed.  $\square$

We note that the Chebyshev polynomials in the standard normalization satisfy (14) with  $b_{k,k\pm 1} = (\pm 1)/(2(k \pm 1))$ . Also the Jacobi polynomials in their standard normalization satisfy a relation of type (14). Strictly speaking, Theorem 2.1 applies only to the Legendre polynomials (although we can scale the symmetric Jacobi polynomials so that (9) applies). In any case, (14) shows that integration is always banded, and of a simple form (the recurrence coefficient matrix) for all the classical orthogonal polynomial families. The discussion following Theorem 2.1 was given to clarify ideas, and in principle could have been omitted.

We now focus on the operator  $D$ . This operator has a one-dimensional null space, and if appropriately restricted, it has an inverse. An especially useful restriction involves the subspace  $Q_1^\infty$ . In this space, the operator  $D$  has a well-defined inverse, which we will denote as  $B$ . Although  $D$  has a full upper triangular matrix representation,  $B$  is banded. Indeed, assuming the recursion in the form of Theorem 2.2, we have that  $B$  is the coefficient matrix for the recurrence (14) (note that this matrix had zeros in the first row since  $Q'_0 = 0$ ).

Similarly,  $D^n$  must be restricted to  $Q_n^\infty$ . Indeed,  $\mathcal{N}(D^n) = Q_0^{n-1}$ , so that the operator  $D^n$  is nonsingular on  $Q_n^\infty$ , the orthogonal complement of its null space, and it has a unique inverse, denoted  $B_{[n]}^n : Q_0^\infty \rightarrow Q_n^\infty$ . Any two images of an element  $z \in Q_0^\infty$  under  $n$ -fold integration differ by an element of  $Q_0^{n-1}$ . The specific form of  $B_{[n]}^n$  fixes that element of  $Q_0^{n-1}$  to be the zero element. Thus, the solutions of

$$(16) \quad Lu = \sum_{k=0}^n (m_{n-k}(x) D^k) u = f(x), \quad u \in Q_n^\infty,$$

and

$$(17) \quad Lu = \sum_{k=0}^n (m_{n-k}(x) D^k B_{[n]}^n) z = f(x), \quad z \in Q_0^\infty,$$

are equivalent. Clearly,  $D^k B_{[n]}^n \neq B_{[n-k]}^{n-k}$ . These operators differ since the second matrix has zeros in its first  $n-k$  rows while the first has some nonzero elements there.



**Example .** The operator  $B_{[2]}^2 : Q_0^\infty \rightarrow Q_2^\infty$  (for families that satisfy Theorem 2.1) is

$$B_{[2]}^2 = \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & \cdots & 0 & \cdots & 0 & 0 \\ \frac{1}{f_0 f_1} & 0 & -\frac{f_1+f_3}{f_3 f_2 f_1} & 0 & \frac{1}{f_3 f_4} & \cdots & 0 & 0 \\ 0 & \frac{1}{f_1 f_2} & 0 & -\frac{f_2+f_4}{f_2 f_3 f_4} & 0 & \frac{1}{f_4 f_5} & \cdots & 0 \\ 0 & 0 & \ddots & 0 & \ddots & 0 & \ddots & \cdots \\ 0 & 0 & \cdots & \frac{1}{f_{k-2} f_{k-1}} & 0 & -\frac{f_{k-1}+f_{k+1}}{f_{k-1} f_k f_{k+1}} & 0 & \frac{1}{f_{k+1} f_{k+2}} \end{pmatrix},$$

and using

$$D = \begin{pmatrix} 0 & f_0 & 0 & f_0 & \cdots & 0 & f_0 & \cdots \\ 0 & 0 & f_1 & 0 & \cdots & f_1 & 0 & \cdots \\ \cdots & \cdots & \cdots & \ddots & \ddots & \cdots & \cdots & \cdots \\ 0 & \cdots & \cdots & \cdots & 0 & f_{2k-1} & 0 & \cdots \\ 0 & \cdots & \cdots & \cdots & \cdots & 0 & f_{2k} & \cdots \end{pmatrix},$$

we find that

$$B_{[1]} \neq DB_{[2]}^2 = \begin{pmatrix} 0 & \frac{f_0}{f_1 f_2} & 0 & -f_0 & 0 & \cdots & 0 \\ \frac{1}{f_0} & 0 & -\frac{1}{f_2} & 0 & \cdots & \cdots & 0 \\ 0 & \frac{1}{f_1} & 0 & -\frac{1}{f_3} & \cdots & \cdots & 0 \\ 0 & 0 & \ddots & \ddots & \ddots & \cdots & 0 \\ 0 & \cdots & 0 & \frac{1}{f_{k-1}} & 0 & -\frac{1}{f_{k+1}} & \cdots \end{pmatrix}.$$

In the general case, the operator  $D$  is hard to write explicitly, as it is the ‘inverse’ (in the sense discussed above) of a general tridiagonal matrix. However, all that is needed in our method is the expression for  $D^k B_{[n]}^n$ , which is identical to the matrix  $B_{[n-k]}^{n-k}$  except for the first  $n - k$  rows which, in general, will contain some nonzero elements. These elements are easy to compute however, as they can be expressed in terms of elements of  $B_{[1]}$  and the  $n \times n$  principal submatrix of  $D$ . For example, the operator  $DB_{[2]}^2$  for the general case is identical to  $B_{[1]}$  except for the first row, whose elements are

$$\text{row}_0 \left( DB_{[2]}^2 \right) = -d_{01} (b_{10} b_{11}, b_{11}^2 + b_{12} b_{21}, b_{12} (b_{11} + b_{22}), b_{12} b_{23}, 0, \dots, 0).$$

Here, the elements  $b_{ij}$  for the classical orthogonal polynomials can be found in Table 1, together with the elements of the matrix  $A$  and other relevant quantities using the standard notation [1]. Also,  $d_{01}$  is the corresponding entry of the differentiation matrix  $D$ , which is simply the derivative of  $Q_1$  expressed in terms of  $Q_0$ . For example, for the general Jacobi polynomials,  $d_{01} = (\alpha + \beta + 2)/2$ , etc.

The relations for the Gegenbauer polynomials  $C_n^{(\nu)}$  can be constructed from those of the Jacobi polynomials since

$$C_n^{(\nu)} = \frac{\Gamma(\alpha + 1)\Gamma(2\alpha + n + 1)}{\Gamma(\alpha + n + 1)\Gamma(2\alpha + 1)} P_n^{(\alpha, \beta)},$$

where  $\alpha = \beta = \nu - 1/2$ . Arrays are indexed from 0 to  $N$ , the maximum order of truncation.

TABLE 1. Recursions for common families ( $a_{0,k} = 0$ ; for the  $T_k$ ,  $a_{1,0} = b_{1,0} = 1$ )

Family	Hermite $H_k$	Laguerre $L_k^{(\alpha)}$	Chebyshev $T_k$	Legendre $P_k$	Jacobi $P_k^{(\alpha,\beta)}$
$Q_0$	1	1	1	1	1
$Q_1$	$2x$	$1 + \alpha - x$	$x$	$x$	$\frac{1}{2}((\alpha - \beta) + (\alpha + \beta + 2)x)$
$a_{k-1,k}$	$k$	$-(k + \alpha)$	$\frac{1}{2}$	$\frac{k}{2k+1}$	$\frac{2(k + \alpha)(k + \beta)}{(2k + \alpha + \beta + 1)(2k + \alpha + \beta)}$
$a_{k,k}$	0	$2k + \alpha + 1$	0	0	$-\frac{(2k + \alpha + \beta + 2)(2k + \alpha + \beta)}{2(k + 1)(k + \alpha + \beta + 1)}$
$a_{k+1,k}$	$\frac{1}{2}$	$-(k + 1)$	$\frac{1}{2}$	$\frac{k+1}{2k+1}$	$\frac{(2k + \alpha + \beta + 2)(2k + \alpha + \beta + 1)}{(2k + \alpha + \beta + 2)(2k + \alpha + \beta + 1)}$
$b_{k-1,k}$	0	0	$-\frac{1}{2(k-1)}$	$-\frac{1}{2k+1}$	$-\frac{2(k + \alpha)(k + \beta)}{(k + \alpha + \beta)(2k + \alpha + \beta)(2k + \alpha + \beta + 1)}$
$b_{k,k}$	0	1	0	0	$\frac{(2k + \alpha + \beta)(2k + \alpha + \beta + 2)}{2(k + \alpha + \beta + 1)}$
$b_{k+1,k}$	$\frac{1}{2(k+1)}$	-1	$\frac{1}{2(k+1)}$	$\frac{1}{2k+1}$	$\frac{(2k + \alpha + \beta + 1)(2k + \alpha + \beta + 2)}{(2k + \alpha + \beta + 1)(2k + \alpha + \beta + 2)}$
$w(x)$	$e^{-x^2}$	$x^\alpha e^{-x}$	$(1 - x^2)^{-1/2}$	1	$(1 - x)^\alpha (1 + x)^\beta$
$p(x)$	$e^{-x^2}$	$x^{\alpha+1} e^{-x}$	$(1 - x^2)^{1/2}$	$(1 - x^2)$	$(1 - x^2)w(x)$
$(a, b)$	$(-\infty, \infty)$	$(0, \infty)$	$(-1, 1)$	$(-1, 1)$	$(-1, 1)$
$h_k$	$\sqrt{\pi} 2^k k!$	$\frac{\Gamma(\alpha + k + 1)}{k!}$	$\pi/2(\pi, k = 0)$	$\frac{2}{2k+1}$	$\frac{2^{\alpha+\beta+1}}{2k + \alpha + \beta + 1} \frac{\Gamma(k + \alpha + 1)\Gamma(k + \beta + 1)}{k!\Gamma(k + \alpha + \beta + 1)}$
$\lambda_k$	$2k$	$k$	$k^2$	$k(k+1)$	$k(k + \alpha + \beta + 1)$

## 3. THE METHOD

We present now a method for the efficient inversion of operators resulting from the spectral solution of the ODE

$$(18) \quad Lu = \sum_{k=0}^n (m_{n-k}(x)D^k)u = f(x) \quad , \quad x \in \Omega = (a, b),$$

subject to the constraints

$$\mathcal{T}u = c.$$

The constraints are represented by a linear functional  $\mathcal{T}$  of rank  $n$ .

We assume now that the matrices  $M_k$ , representing multiplication by  $m_k$ , are banded. This is for example the case if the original DE had rational coefficients. After multiplying out the denominators, we are left with low-order polynomial coefficients, and as a result of the simple recursion operator  $A$  of multiplication by the monomial  $x$  (5), these have banded representations as convolution operators. The simplest form is found if we expand the resulting polynomial coefficients in terms of the  $Q_k$ , then exploit the properties of the banded operators of multiplication by  $Q_k$ . In constructing an approximate solution, we look for a solution of

$$L_N u_N = f_N \quad ; \quad u_N \in Q_0^N, \quad f_N \in Q_0^{N-n},$$

with  $L_N$  the Galerkin approximation to  $L$ , as usual [12]. The main result can be expressed as follows:

**Theorem 3.1.** *Assume that the  $M_k$  are banded. Also assume that  $Q_0^N = \mathcal{N}(L_N) \oplus Q_n^N$ . Then, if there is a solution  $u_N$ , it can be written as a combination of an element  $w \in \mathcal{N}(L_N)$  and an element  $u_p \in Q_n^N$  such that  $L_N u_p = f$ . The solution of the latter problem can be performed in  $\mathcal{O}(N)$  operations.*

To construct the particular solution, let  $z = D^n u_p \in Q_0^{N-n}$  so that  $u_p = B_{[n]}^n z \in Q_n^N$  is uniquely defined. Then the equation can be rewritten

$$(19) \quad \sum_{k=0}^n (M_{n-k} D^k) \hat{U}_p = \sum_{k=0}^n M_{n-k} D^k B_{[n]}^n \hat{Z} = \hat{F},$$

where we have introduced  $\hat{U}_p$ ,  $\hat{Z}$ ,  $\hat{F}$  to represent the vectors of expansion coefficients. (We use the same notation, however, for the differential and integral operators as for their matrix representations.) Since a solution  $u_p$  of this problem was guaranteed to exist,  $z$ , its  $n$ th derivative, exists as well. Here we must note that in the case of weak solutions the highest derivative must be handled carefully, but in this case convergence would be slow and the method would be impractical. Now we address our main question: Is the new system any easier to solve than the original? This is clearly true: the integration operators are all banded, and to find  $u$  from  $z$  we perform one more banded matrix multiplication. To simplify the notation, in the rest of the paper, unless otherwise stated, we will write  $D^{j-n} \equiv D^j B_{[n]}^n$ ,  $j = 0, \dots, n-1$ , where  $n$  is the order of the differential operator under investigation.

Finally, we need to determine a convenient basis for the nullspace of the operator  $L_N$ . We define

$$e_k = Q_k + w_k \quad , \quad w_k \in Q_n^N,$$

with

$$L_N e_k = 0 \Rightarrow L_N w_k = -L_N Q_k.$$

Then

$$\mathcal{T}u = \mathcal{T}u_p + \sum_{k=0}^{n-1} \alpha_k \mathcal{T}e_k = c \quad ,$$

so that, when the numbers  $T_{kl} = (\mathcal{T}e_k)_l$  are found, we have

$$\sum_{k=0}^{n-1} \alpha_k T_{kl} = c_l - (\mathcal{T}u_p)_l.$$

In other words, for every new right-hand side we simply need to solve the standard BVP for  $u_p$ , then evaluate the quantities  $(\mathcal{T}u_p)_l$  and solve an  $n \times n$  system for the  $\alpha_k$ . The condition of this system is known in advance.

**Example .** The radial Laplace equation for the  $n$ th Fourier mode is

$$(x+a) \frac{\partial}{\partial x} \left( (x+a) \frac{\partial u}{\partial x} \right) - n^2 u = f \quad , \quad u \in Q_2^N.$$

This leads to the pentadiagonal matrix

$$(x+a)^2 I + (x+a) D^{-1} - n^2 D^{-2}.$$

**Example .** The helical Laplace equation for the  $n$ th Fourier mode is

$$\frac{1}{r} \partial_r \left( \frac{r}{1+r^2 \alpha^2} \partial_r u \right) - \frac{n^2}{r^2} u = f.$$

This is transformed to the (almost nine-diagonal) matrix

$$r^2(1+r^2 \alpha^2) I + r(1-\alpha^2 r^2) D^{-1} - n^2(1+\alpha^2 r^2) D^{-2}.$$

**Example .** The initial value problem for the Airy equation,

$$y'' - \alpha^3(x-x_0)y = 0 \quad , \quad y(0) = .355029403792807 \quad , \quad y'(0) = .258819403792807\alpha \quad ,$$

has the solution

$$y(x) = \text{Ai}(\alpha(x-x_0)),$$

the Airy function of the first kind. Here we include the parameters  $\alpha$ ,  $x_0 \in [-1, 1]$  in order to scale the interval over which the problem is solved, since Chebyshev expansions apply naturally over the interval  $[-1, 1]$ . For  $x > 0$  the solutions decay exponentially and the numerical algorithm converges rapidly. However, for  $x < 0$  the solutions exhibit oscillatory behavior with ever increasing frequency, and convergence can only be achieved if sufficient modes are included to resolve the most rapid oscillations present. In Figure 1 we show the solution to the problem with  $x_0 = -1$ ,  $\alpha = 10$  with  $N = 30, 40$ , respectively. The first case is underresolved, and the maximum absolute error is  $O(10^{-2})$ , while the second case is barely resolved, and the error is  $O(10^{-5})$ . A slight increase in the order of truncation improves the solution dramatically. With  $N = 64$ , the error is already less than  $10^{-11}$ .

**Example .** The two-dimensional Stokes problem is expressed by the system

$$(20) \quad \Delta_N \psi = -\omega,$$

$$(21) \quad H\omega = f \quad ,$$

where  $\Delta_N$  is the nonperiodic part of the Laplacian for the  $m$ th Fourier mode in a two-dimensional geometry with one nonperiodic and one periodic direction; likewise,  $H$  is a second-order linear operator with rational coefficients. No conditions are given on  $\omega$  while  $\psi$  and  $\psi_x$  are specified at  $x = \pm 1$ . Such a system results, e.g., from the time discretization of the Stokes equations in appropriate two-dimensional

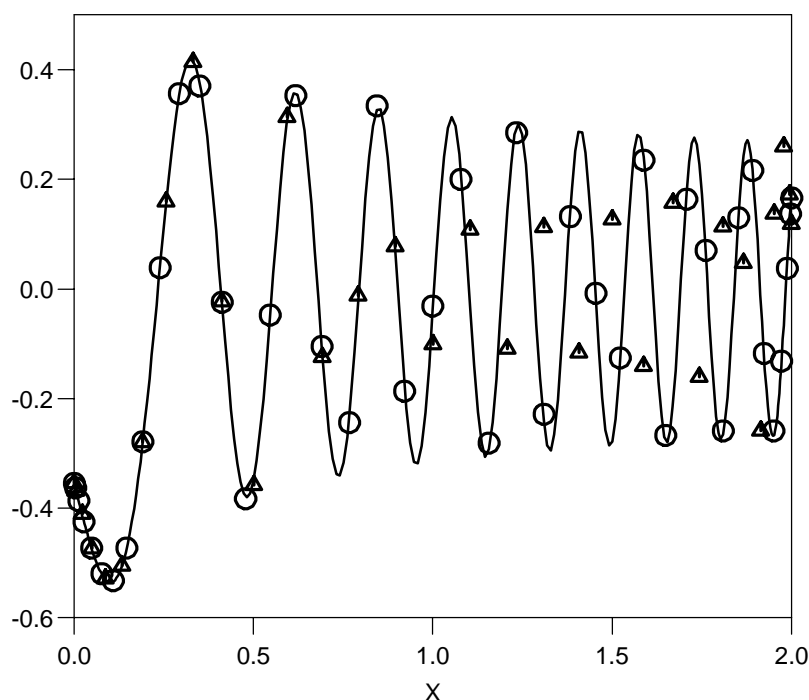


FIGURE 1. Computed solutions with 30 modes (triangles) and 40 modes (circles) are plotted versus the exact solution of the Airy equation for  $\alpha = 10$

domains. We consider projections  $f \rightarrow f_{N-4} \in Q_0^{N-4}$ ,  $\omega \rightarrow \omega_{N-2} \in Q_0^{N-2}$  and  $\psi \rightarrow \psi_N \in Q_0^N$ . We determine a particular solution for (21) as  $\omega_p \in Q_2^{N-2}$  and homogeneous solutions  $\omega_k = Q_k + \Omega_k$ ,  $k = 0, 1$ , with  $\Omega_k \in Q_2^{N-2}$ . The general solution for (21) is then

$$(22) \quad \omega_{N-2} = \omega_p + \alpha_0 \omega_0 + \alpha_1 \omega_1.$$

The general solution of (20) can now be written as  $\psi_N = \psi_p + \beta_0 \psi_0 + \beta_1 \psi_1$ , where  $\psi_k = Q_k + \Psi_{k+2}$ ,  $k = 0, 1$  (with  $\Psi_{k+2} \in Q_2^N$ ) are the homogeneous solutions and  $\psi_p = \Psi_p + \alpha_0 \Psi_0 + \alpha_1 \Psi_1 \in Q_2^N$  is a particular solution with  $\Delta_N \Psi_p = -\omega_p$  and  $\Delta_N \Psi_k = -\omega_k$ ,  $k = 0, 1$ . The boundary conditions can now be applied to  $\psi_N$  to produce a  $4 \times 4$  system in the  $\alpha_k, \beta_k$ ,  $k = 0, 1$ :

$$\mathcal{A}\sigma = \phi,$$

with

$$\begin{aligned} \mathcal{A}_{1,j+1} &= \Psi_j(1), & \mathcal{A}_{2,j+1} &= \Psi_j(-1), \\ \mathcal{A}_{3,j+1} &= \Psi_{j,x}(1), & \mathcal{A}_{4,j+1} &= \Psi_{j,x}(-1) \end{aligned}$$

( $j = 1, 2, 3, 4$ ), so that  $\mathcal{A}$  need only be evaluated once,  $\sigma_{k+1} = \alpha_k$ ,  $k = 0, 1$ , and

$\sigma_{k+3} = \beta_k$ ,  $k = 0, 1$ , and  $\phi_1 = \psi(1) - \Psi_p(1)$  etc. If this problem is solved as part of a time integration of the Navier Stokes equations in a two-dimensional region, the only functions that need to be computed repeatedly are  $\Psi_p \in Q_2^N$  and  $\omega_p \in Q_2^{N-2}$ .

Finally, we comment on the solution of the BVP with arbitrary BC. In this case, we need to determine the null space, and add an arbitrary combination of nullvectors to the particular solution to satisfy any desired BC. An alternative form of our method can also be considered. It is based on commuting the polynomials with the differential operators and multiplying on the left by the integration operator  $D_n$ , so that the differential operator matrix becomes banded. This approach is discussed in [8]. This is in fact the  $\tau$ -method, and various instances that have been worked out (e.g. [12, 9]) are of this type. Theorem 2.2 establishes the success of this approach as a consequence of the basic recurrence relation (14).

However, other preconditioners may be available, depending on the special structure of the matrix operator  $L_N$  [18]. The present method's main appeal, besides the fact that a convergence analysis is available (see §4) is its simplicity and generality. Indeed, the complicated expressions for the differential operators are entirely avoided. Of course, it should be expected that in special cases simpler forms and preconditioners might be possible. An example of this is provided by the Laplace equation in a circle, for which the integration preconditioner leads to pentadiagonal forms while a simpler tridiagonal form is in fact possible [18]. This is related to the special properties of the operator  $(x+1)\frac{d}{dx}$ . We have not yet explored the flexibility of the choice of spaces, which may sometimes lead to more efficient algorithms.

#### 4. STABILITY AND CONVERGENCE

As differential operators are unbounded (in the usual norms), so are their difference and spectral analogues under mesh refinement. Numerical studies have shown that the spectral radii of Chebyshev and Legendre differentiation matrices are  $O(N)$ , where  $N$  is the dimension of the subspace. Moreover, these matrices are far from normal, so that their norms can grow even faster. For example, the maximum norm of the differentiation matrix,  $D$ , for a family satisfying the simple recursion (14), such as the Legendre polynomials, satisfies the lower bound

$$(23) \quad \|D\|_\infty \geq \max_{m \leq N} f(m) \frac{N-m}{2}.$$

From Table 1 we see that  $f(m) = O(m)$ , so that the lower bound above is  $O(N^2)$ . This lower bound grows accordingly with the order of the derivative being approximated. (For example the second derivative can behave as  $O(N^4)$ .)

The poor conditioning of the matrices arising from spectral discretizations both limits the accuracy of solutions, owing to roundoff errors, and imposes severe limitations on the time step for explicit solutions of dynamic problems [17] or on the convergence rate of iterative solvers. It is well known that the reformulation of differential equations as integral equations often leads to bounded operators and well-conditioned problems. As our formulation of the discrete equations is based on integral operators, we also expect to obtain well-conditioned linear systems. For some constant-coefficient problems, Greengard [13] has directly analyzed spectral approximations to equivalent integral equations and demonstrated the gains in accuracy which can be attained. In this section we generalize and expand on

Greengard's results to cover the algorithms we have proposed. We use the stability estimates we obtain to prove convergence of the method.

**4.1. Estimates of the condition number.** We assume the differential equation takes the form (18) where the coefficients  $m_j(x)$  are polynomials and  $m_0$  is bounded away from zero on  $[a, b]$ . Of course, this last condition is required to avoid singularities in the solution, where the spectral approximation itself may not be well behaved. We concentrate on the system (19), used to determine a particular solution

$$(24) \quad \left( M_0 + \sum_{j=1}^n M_j D^{-j} \right) \hat{Z} \equiv A\hat{Z} = \hat{F},$$

where the  $M_j$  are the Galerkin approximations to multiplication by the polynomial coefficients and  $D^{-j} = D^{n-j} B_{[n]}^n$ . The additional problem to be solved, involving the boundary conditions, will be of much lower dimension, and its conditioning will depend on the specific constraint conditions. We begin by estimating various norms of the integration operators. We view these, now, as operators on  $l_2$  and  $l_\infty$ . The bounds we derive obviously extend to finite truncations. We make the following assumptions about the orthogonal family, which are satisfied for proper normalizations of any of the Jacobi polynomials, as well as the Hermite polynomials (see Table 1):

**Assumption 4.1.** The orthogonal family  $Q_k$  satisfies

$$(25) \quad \frac{\sup_k \langle Q_k, Q_k \rangle_\omega}{\inf_k \langle Q_k, Q_k \rangle_\omega} = \kappa_Q^2 < \infty,$$

$$(26) \quad |b_{kj}| \leq \frac{\alpha}{k^p}, \quad p > 0.$$

Here,  $\langle \cdot, \cdot \rangle_\omega$  denotes the weighted inner product defining the family. The best exponent,  $p$ , is 1 for the Jacobi family and 1/2 for the Hermite family.

(To use Table 1 to verify the statement concerning  $p$ , form the normalized families by dividing  $Q_k$  by  $\sqrt{h_k}$  and note that  $b_{kj}$  is transformed to  $b_{kj} \sqrt{h_k/h_j}$ .) We have the following lemma, describing the structure of the integration matrices:

**Lemma 4.1.** *The matrices  $D^{-j}$ ,  $j = 1, \dots, n$ , are banded with bandwidth  $j$ , with the possible exception of a finite number of elements in the first  $j$  rows, and there exists a constant  $\bar{B}_j$  such that  $|(D^{-j})_{kl}| \leq \bar{B}_j k^{-jp}$ .*

*Proof.* We first note that the integration operator  $B_{[n]}^n = D^{-n}$  coincides with  $(B_{[1]}^1)^n$  except for the first  $n$  rows which are zero. Since  $B_{[1]}^1$  is tridiagonal with elements satisfying (26), the result is immediate. For the other terms we use the fact that  $DB_{[1]}^1 = I$  to write

$$(27) \quad D^{-j} = D^{n-j} B_{[n]}^n = D^{n-j} ((B_{[1]}^1)^n + C) = (B_{[1]}^1)^j + D^{n-j} C,$$

where the nonzero elements of  $C$  are simply the negatives of the nonzero elements of  $(B_{[1]}^1)^n$  in the first  $n$  rows. Since  $D^{n-j}$  is upper triangular with nonzero elements only in superdiagonals  $n-j$  to  $\infty$ , we see that the nonzero elements of  $D^{(n-j)}C$  are restricted to the first  $j$  rows and  $2n+1$  columns. Since the result holds for  $(B_{[1]}^1)^j$ , the lemma is proved.  $\square$

This leads immediately to the following theorem:

**Theorem 4.1.** *For  $r = 2$  or  $r = \infty$ , the operators  $D^{-j} : l_r \rightarrow l_r$  are compact.*

*Proof.* The boundedness of the infinity norm follows directly from Lemma 4.1 and the fact that the norm is equal to the maximum absolute row sum. For the 2-norm we have

$$\begin{aligned}
 \|D^{-j}y\|_2 &= \left( \sum_{i=1}^{\infty} \left| \sum_{l=\max(1,i-j)}^{\max(2n+1,i+j)} (D^{-j})_{il} y_l \right|^2 \right)^{1/2} \\
 (28) \quad &\leq \sqrt{2n+1} \bar{B}_j \left( \sum_{i=1}^{\infty} \sum_{l=\max(1,i-j)}^{\max(2n+1,i+j)} |y_l|^2 \right)^{1/2} \\
 &\leq (2n+1) \bar{B}_j \|y\|_2.
 \end{aligned}$$

To prove compactness it is sufficient to show that  $D^{-j}$  can be approximated by a sequence of bounded operators of finite rank,  $\{D_M^{-j}\}$ . For these we simply take the operators defined by setting all rows below  $M$  to zero. Repeating the arguments above we have, for  $r = 2$  or  $r = \infty$ ,

$$(29) \quad \|D^{-j} - D_M^{-j}\|_r \leq (2j+1) \bar{B}_j M^{-jp} \rightarrow 0, \quad M \rightarrow \infty,$$

completing the proof.  $\square$

We would now like to bound the norms and condition numbers of the Galerkin polynomial multiplication matrices,  $M_k$ . We begin by showing that the Galerkin matrix representing multiplication by an arbitrary polynomial,  $\phi(x)$ , is nonsingular if the zeros of  $\phi$  lie outside  $[a, b]$ .

**Theorem 4.2.** *Let  $\Phi$  be the matrix representation of the Galerkin approximation to multiplication by the degree- $q$  polynomial  $\phi(x)$  relative to the orthogonal system  $\{Q_j(x)\}_0^N$  on  $[a, b]$ . If the zeros of  $\phi$  lie outside  $[a, b]$ , then  $\Phi$  is nonsingular.*

*Proof.* Suppose the contrary. Then there exists a nonzero polynomial,  $\mu$ , of degree  $N$  such that  $\phi\mu = \sum_{j=N+1}^{N+q} c_k Q_k(x)$ . We then have that  $\phi\mu$  is orthogonal to all polynomials of degree less than or equal to  $N$  and has at least  $q$  zeros (counting multiplicities) outside  $[a, b]$ . This implies that  $\phi\mu$  has at most  $N$  zeros of odd multiplicity in  $(a, b)$ . Let  $r_i, i = 1, \dots, s$ , denote these zeros. Then  $\psi(x) = \prod_{i=1}^s (x-r_i)$  is a polynomial of degree less than or equal to  $N$  such that  $\phi\mu\psi$  is of one sign on  $[a, b]$ . However, we also have  $\int_a^b \omega \phi\mu\psi dx = 0$  by the orthogonality of  $\phi\mu$  to polynomials of degree not more than  $N$ . This is a contradiction, so  $\mu$  cannot exist.  $\square$

An immediate corollary of this theorem is:

**Corollary 4.1.** *The spectrum of  $\Phi$  is contained within  $\{y = \phi(x), x \in [a, b]\}$ .*

*Proof.* Suppose  $\Phi - \lambda I$  is singular. Since  $\Phi - \lambda I$  is the Galerkin approximation to multiplication by  $\phi - \lambda$ , we conclude  $\phi - \lambda$  must have a zero in  $[a, b]$ , completing the proof.  $\square$

From the eigenvalues we can easily bound the norms:



**Theorem 4.3.** *The matrix,  $\Phi$ , satisfies the following bounds:*

$$(30) \quad \|\Phi\|_2 \leq \kappa_Q \max_{x \in [a,b]} |\phi(x)|, \quad \|\Phi^{-1}\|_2 \leq \kappa_Q \max_{x \in [a,b]} |(1/\phi(x))|.$$

*Proof.* Let  $\tilde{\Phi}$  denote the matrix representing the Galerkin approximation to multiplication by  $\phi$  relative to the orthonormal basis obtained by normalizing the  $Q$ 's. Since  $\tilde{\Phi}$  is symmetric, its 2-norm and the 2-norm of its reciprocal are bounded, respectively, by the largest and the inverse of the smallest eigenvalues (in absolute value). These are in turn bounded by  $\max_{x \in [a,b]} |\phi(x)|$  and  $\max_{x \in [a,b]} |(1/\phi(x))|$  from Corollary 4.1. Let  $R = \text{diag}(\sqrt{\langle Q_i, Q_i \rangle_\omega})$ . Then  $\Phi = R^{-1} \tilde{\Phi} R$ . Taking norms yields the final result.  $\square$

We have shown that the system defining the particular solution has the form  $M_0 + K$ , where, for regular problems,  $M_0$  has a bounded condition number uniformly in  $N$ , and  $K$  approaches a compact operator. To complete our analysis, we must develop lower bounds on  $M_0 + K$ , which can only be expected if the homogeneous differential equation admits no nontrivial solutions in  $(Q_0^{n-1})^\perp$ . We make this explicit in the following assumption.

**Assumption 4.2.** If  $w$  is a solution of the homogeneous problem ((18) with  $f = 0$ ) satisfying  $\langle w, Q_k \rangle_\omega = 0$  for all  $k = 0, \dots, n-1$ , then  $w = 0$ .

We remark that the existence of a nontrivial solution of the homogeneous differential equation which is orthogonal to all polynomials of degree less than  $n$  is clearly not generic. If it holds, the difficulties with the method can be remedied by looking for particular solutions in a different subspace. In the future we plan to consider the case of singular problems in more detail, particularly in cases where the lead coefficient is zero somewhere in  $[a, b]$ .

We now define the operator  $\tilde{K} : l_2 \rightarrow l_2$  by

$$(31) \quad (\tilde{K}y)_k = h_k^{-1} \left\langle Q_k, (m_0)^{-1} \sum_{j=1}^n m_j \left( \sum_{i=0}^{\infty} (D^{-j}y)_i Q_i \right) \right\rangle_\omega.$$

We then have:

**Lemma 4.2.** *The operator  $\tilde{K}$  is compact and, if Assumption 4.2 holds,  $(I + \tilde{K})^{-1}$  is bounded.*

*Proof.* The proof of compactness again follows by approximating  $D^{-j}$  by  $D_M^{-j}$ . If  $\tilde{K}_M$  denotes the resulting approximation to  $\tilde{K}$ , it is clear that  $\tilde{K}_M$  is bounded and has finite rank. Moreover, as  $M \rightarrow \infty$ ,

$$(32) \quad \|(\tilde{K} - \tilde{K}_M)y\|_2 \leq n\kappa_Q \max_{j=1, \dots, n} \left( \max_{x \in [a,b]} |m_0^{-1}(x)m_j(x)| \right) \|(D^{-j} - D_M^{-j})\|_2 \|y\|_2 \rightarrow 0.$$

Therefore,  $\tilde{K}$  is compact. By the Riesz-Schauder theory the boundedness of  $(I + \tilde{K})^{-1}$  holds if and only if  $(I + \tilde{K})z = 0$  has no nontrivial solution in  $l_2$ . Suppose such a solution exists. Let  $w = D^{-n}z$ ,  $w = \sum_{k=n}^{\infty} W_k Q_k$ . Then we can write

$$(33) \quad Lw = m_0(1 + m_0^{-1} \sum_{j=1}^n m_j D^{-j}) D^n w = m_0 \sum_k \left( Q_k((I + \tilde{K})z)_k \right) = 0.$$

That is,  $w$  is a weak solution of the homogeneous problem in  $(Q_0^{n-1})^\perp$ , with  $n$  derivatives in  $L_\omega^2$ . By repeated differentiation and use of the fact that  $m_0$  is bounded away from zero, we establish that arbitrary derivatives are in  $L_\omega^2$ . Since  $\omega$  is bounded above and below in arbitrary closed subintervals of  $(a, b)$ , Sobolev's inequality implies that  $w$  is a classical solution, violating Assumption 4.2. This completes the proof.  $\square$

We are now in a position to uniformly bound the condition number of  $A$ .

**Theorem 4.4.** *Suppose Assumption 4.2 holds. Then there exist constants  $C_0$  and  $C_1$  and an integer  $N_0$  such that for all  $N > N_0$  and vectors  $y$  with Euclidean norm 1,*

$$(34) \quad C_0 \leq \|Ay\|_2 \leq C_1.$$

*Proof.* The finite system can be written in the form  $A_N = M_{0,N}(I + K_N)$ , where

$$(35) \quad K_N = M_{0,N}^{-1} \sum_{j=1}^n M_{j,N} D_N^{-j}.$$

Here, the subscript  $N$  indicates that the degree- $N$  Galerkin approximation is being considered. The existence of the uniform upper bound,  $C_1$ , follows from Lemmas 4.1 and 4.3. To deduce the existence of the lower bound, we define  $\bar{K}_N : l_2 \rightarrow l_2$  by

$$(36) \quad (\bar{K}_N y)_k = \begin{cases} (K_N y_N)_k, & k = 0, \dots, N, \\ 0, & k > N. \end{cases}$$

Here,  $y_N$  is the  $(N+1)$ -vector formed from the first  $N+1$  components of  $y$ . Let  $\epsilon > 0$  be given. We will find  $N(\epsilon)$  such that  $\|\bar{K}_N - \tilde{K}\| < \epsilon$  for  $N > N(\epsilon)$ . Given any element,  $y$ , of  $l_2$  with norm 1 and any  $M > 0$ , set  $y = y_M + x_M$ , where only the first  $M+1$  components of  $y_M$  are nonzero and the first  $M+1$  components of  $x_M$  equal 0. For  $M$  sufficiently large, independent of  $y$  and  $N$ , the estimates in the proof of Lemma 4.1 imply that  $\|\bar{K}_N x_M\|, \|\tilde{K} x_M\| < \epsilon/4$ . Moreover, for  $N > M + n + q$ , where  $q$  is the maximum degree of  $m_j$ ,  $j = 1, \dots, n$ , we have

$$(37) \quad s(x) = \sum_{j=1}^n m_j(x) \left( \sum_{i=0}^{M+j} (D^{-j} y_M)_i Q_i(x) \right) = \sum_{j=1}^n \sum_{i=0}^{M+q+j} (M_{j,N} D_N^{-j} y_M)_i Q_i(x).$$

Set  $e = (\bar{K}_N - \tilde{K})y_M$  and let  $e = \bar{e} + \tilde{e}$ , where  $\bar{e}$  is nonzero only in the first  $N+1$  components and the first  $N+1$  components of  $\tilde{e}$  are zero. Note that

$$(38) \quad \tilde{e}_k = -h_k^{-1} \langle Q_k, m_0^{-1} s \rangle_\omega, \quad k \geq N+1.$$

Denote by  $S$  the vector of expansion coefficients of  $s$ , recalling that all components,  $S_k$ , are zero for  $k > M + n + q$ . Then, treating vectors whose nonzero components have index no greater than  $N$  as  $N+1$  vectors, we obtain

$$(39) \quad M_{0,N} \bar{e} = S - M_{0,N} (\tilde{K} y_M)_N.$$

However, by the bandedness of  $M_0$ , we have

$$(40) \quad S = M_{0,N,N+q} (\tilde{K} y_M)_{N+q},$$

where  $M_{0,N,N+q} = (M_{0,N} \ E_{N,q})$  is the rectangular matrix formed from the first  $N + 1$  rows of  $M_{0,N+q}$ . Therefore,

$$(41) \quad M_{0,N}\bar{e} = (M_{0,N} \ E_{N,q})(\tilde{K}y_M)_{N+q} - (M_{0,N} \ O)(\tilde{K}y_M)_{N+q} = (O \ E_{N,q})(\tilde{K}y_M)_{N+q}.$$

Hence,

$$(42) \quad \|\bar{e}\| = \|(O \ M_{0,N}^{-1}E_{N,q})(\tilde{K}y_M)_{N+q}\| \leq \|M_{0,N}^{-1}\| \cdot \|M_{0,N+q}\| \cdot \|\bar{e}\|.$$

Therefore, we have, for some constant  $\bar{C}$ ,

$$(43) \quad \|(\bar{K}_N - \tilde{K})y_M\| \leq \bar{C}\|\bar{e}\|.$$

For fixed  $M$  the functions  $m_0^{-1}s$  as well as any of their derivatives may be bounded independent of  $y_M$ ,  $\|y_M\| \leq 1$ . Therefore, for any integer  $\mu > 0$ , standard approximation results (e.g. [5, Ch. 9]) imply the existence of constants  $C(\mu, M)$  such that the right-hand side of (43) is bounded above by  $C(\mu, M)N^{-\mu}$ . We may then choose  $N(\epsilon, M)$  sufficiently large that

$$(44) \quad \max_{\|y_M\| \leq 1} \|(\bar{K}_N - \tilde{K})y_M\| < \frac{\epsilon}{2}, \quad N > N(\epsilon, M).$$

We finally have, for  $M = M(\epsilon)$ ,  $N > N(\epsilon, M(\epsilon))$  and  $\|y\| = 1$ ,

$$(45) \quad \|(\bar{K}_N - \tilde{K})y\| \leq \|(\bar{K}_N - \tilde{K})y_M\| + \|\bar{K}_N x_M\| + \|\tilde{K}x_M\| < \epsilon.$$

By the Banach lemma,

$$(46) \quad \|(I + \bar{K}_N)^{-1}\| \leq \|(I + \tilde{K})^{-1}\|(1 - \epsilon\|(I + \tilde{K})^{-1}\|)^{-1},$$

for  $N > N(\epsilon)$  and  $\epsilon < (\|(I + \tilde{K})^{-1}\|)^{-1}$ . Since  $(I + K_N)^{-1}$  is a block diagonal submatrix of  $(I + \bar{K}_N)^{-1}$  it follows that  $\|(I + K_N)^{-1}\| \leq \|(I + \bar{K}_N)^{-1}\|$ . As we have uniform lower bounds on  $M_0$ , the existence of  $C_0$  follows, completing the proof.  $\square$

**4.2. Error estimates.** Given these bounds on the condition number of the linear system, a convergence result is easily proved. We restrict attention to symmetric Jacobi (Gegenbauer) polynomials, where good results for interpolation have been obtained by Bernardi and Maday [4]. We explicitly assume that the original problem has the following properties:

- Assumption 4.3.** (a) The constraint/boundary operators  $\mathcal{T}$  satisfy an inequality of the form  $|\mathcal{T}w| \leq \|w\|_{\omega,n}$ .  
 (b) The forcing function,  $f(x)$ , is in  $C^r([a, b])$ ,  $r \geq 1$ .  
 (c) If  $w$  is a solution of the homogeneous problem ((18) with  $f = 0$ ) satisfying  $\mathcal{T}w = 0$  or  $\langle w, Q_k \rangle_{\omega} = 0$  for all  $k = 0, \dots, n-1$ , then  $w = 0$ .

We now prove a sequence of estimates of various parts of the error. For the continuous problem, Assumption 4.3 implies the following (e.g. [7]):

- (i) There exists a basis,  $\{u_j\}_{j=0}^{n-1} \in C_{\infty}([a, b])$ , for the space of solutions to the homogeneous problem taking the form  $u_j = Q_j + \tilde{u}_j$ , with  $\tilde{u}_j$  orthogonal to  $Q_0^{n-1}$ .  
 (ii) There exists a unique solution,  $u \in C_{n+r}([a, b])$ , which can be written  $u(x) = u_s(x) + \sum_{j=0}^{n-1} c_j u_j(x)$  with  $u_s$  orthogonal to  $Q_0^{n-1}$ .

(iii) The  $n \times n$  matrix

$$T_e = [ \mathcal{T}u_0 \quad \mathcal{T}u_1 \quad \dots \quad \mathcal{T}u_{n-1} ]$$

is nonsingular.

Let  $v_s(x)$  denote the approximate particular solution, that is, the polynomial whose expansion coefficients are given by  $D^{-n}z$ , and  $v_j(x)$  denote the approximate solution of the homogeneous problem taking the form  $Q_j + \tilde{v}_j$  with  $\tilde{v}_j$  orthogonal to  $Q_0^{n-1}$ . We assume that the right-hand side of the inhomogeneous equation is obtained via interpolation at the relevant Gauss or Gauss-Lobatto points. Let

$$(47) \quad T_a = [ \mathcal{T}v_0 \quad \mathcal{T}v_1 \quad \dots \quad \mathcal{T}v_{n-1} ].$$

We then have:

**Lemma 4.3.** *There exists  $N_0$  such that for  $N > N_0$ :*

(i) *There exist constants  $G_{l,\mu}$  such that*

$$\|u_j^{(l)} - v_j^{(l)}\|_\omega \leq G_{l,\mu} N^{-\mu}, \quad 0 \leq l, \mu < \infty, \quad j = 0, \dots, n-1.$$

(ii) *There exist constants  $R_\mu$  such that*

$$\|T_e - T_a\| \leq R_\mu N^{-\mu}, \quad 0 \leq \mu < \infty.$$

(iii) *There exist constants  $D_l$  such that*

$$\|u_s^{(l)} - v_s^{(l)}\|_\omega \leq D_l N^{-r} \|f\|_{\omega,r}, \quad 0 \leq l \leq n.$$

*Proof.* We rely extensively on the approximation results listed in [4] and [5, Ch. 9]. Now  $u_j - v_j = \tilde{u}_j - \tilde{v}_j$ . Let  $\tilde{u}_j = \tilde{u}_{j,N} + \tilde{w}_j$ , where  $\tilde{u}_{j,N}^{(n)} \in Q_0^{N-n}$  and  $\tilde{w}_j^{(n)} \in Q_{N-n+1}^\infty$ . Then

$$(48) \quad \tilde{u}_j - \tilde{v}_j = (\tilde{u}_{j,N} - \tilde{v}_j) + \tilde{w}_j.$$

Estimates of the last term and its derivatives follow directly from results on approximation by singular Sturm-Liouville eigenfunctions and the smoothness of  $\tilde{u}_j$ . For the first, we rewrite the expansion coefficients as  $B_{[n]}^n \hat{Z}_{u,j,N}$  and  $B_{[n]}^n \hat{Z}_{v,j}$  and introduce  $\hat{Z}_{e,j} = \hat{Z}_{u,j,N} - \hat{Z}_{v,j}$ . Let  $\bar{Z}_{u,j,N} = \hat{Z}_{u,j,N+Q} - E_{N+Q} \hat{Z}_{u,j,N}$ , where  $Q$  is the bandwidth of the matrices  $A$  and  $E$  represents extension by 0 of a vector to a longer vector. Denoting explicitly by  $A_m$  the matrix  $A$  associated with degree- $(m+n-1)$  truncations and by  $P_m$  the restriction of a vector of order larger than  $m$  to the  $m$ -vector containing its first  $m$  components, we have

$$(49) \quad A_{N-n} \hat{Z}_{e,j} = P_{N-n} A_{N-n+Q} \bar{Z}_{u,j,N}.$$

By the properties of  $A$  we have

$$(50) \quad \|\hat{Z}_{e,j}\|_2 \leq C \|\bar{Z}_{u,j,N}\|_2.$$

Now  $\bar{Z}_{u,j,N}$  can be estimated by derivatives of  $\tilde{w}_j$ . Therefore, we have estimates of the  $n$ th derivative of the error in terms of the difference between  $\tilde{u}_j^{(n)}$  and its projection into  $Q_0^{N-n}$ . From [5] we directly obtain the estimate in (i). For lower derivatives we simply apply the bounded operators  $D^{-j}$  to  $\hat{Z}_{e,j}$ . For higher derivatives we apply the derivative operators, which, though unbounded, still contribute only polynomial growth. To derive (ii), we use the estimates in (i) and assumption (a) on the constraint operators.

Estimates of the particular solution follow the same pattern. Introduce  $\bar{f}_N$ , the polynomial approximation to  $f$  used to compute  $v_s$ , and  $f_N$ , the orthogonal projection of  $f$  into  $Q_0^{N-n}$ . Write  $u_s^{(n)} = u_{s,N}^{(n)} + w_s^{(n)}$ , where  $u_{s,N}^{(n)} \in Q_0^{N-n}$  and  $w_s^{(n)} \in Q_{N-n+1}^\infty$ . Let the expansion coefficients of  $u_{s,N}$  and  $v_s$  be given by  $B_{[n]}^n \hat{Z}_{s,u,N}$  and  $B_{[n]}^n \hat{Z}_{s,v}$  respectively. Let  $\hat{R}_s = \hat{Z}_{s,v} - \hat{Z}_{s,u,N}$  and  $\bar{Z}_{s,u,N} = \hat{Z}_{s,u,N+Q} - E_{N+Q} \hat{Z}_{s,u,N}$ . Then we have

$$(51) \quad A_{N-n} \hat{R}_s = P_{N-n} A_{N-n+Q} \bar{Z}_{s,u,N} + \bar{F}_N - \hat{F}_N.$$

From the boundedness of the  $A$ 's, the first term can be estimated in terms of  $\|w_s^{(n)}\|_\omega = O(N^{-r}) \cdot \|w_s^{(n)}\|_{\omega,r}$ . (This holds because  $u_{s,N}$  was constructed by projecting  $u_s^{(n)}$  into  $Q_0^{N-n}$  and applying integration operators.) The second is bounded by  $\|f - \bar{f}_N\|_\omega + \|f - f_N\|_\omega = O(N^{-r}) \cdot \|f\|_{\omega,r}$ . The boundedness of  $A^{-1}$  then implies (iii) for the  $n$ th derivative. The bounds for the lower derivatives then follow by application of the bounded integration operators.  $\square$

We are now in a position to prove:

**Theorem 4.5.** *For some  $N_0 < \infty$  there exist constants  $H_l$  such that, for all  $N > N_0$ , the difference between the true solution,  $u$ , and the approximate solution,  $v$ , satisfies*

$$\|u^{(l)} - v^{(l)}\|_\omega \leq H_l N^{-r} \|f\|_{\omega,r}, \quad l = 0, \dots, n.$$

*Proof.* We have

$$(52) \quad u = u_s + \sum_{j=0}^{n-1} \gamma_j u_j, \quad T_e \gamma = c - T u_s,$$

$$(53) \quad v = v_s + \sum_{j=0}^{n-1} \delta_j v_j, \quad T_a \delta = c - T v_s.$$

Introducing  $e = u - v$ ,  $\nu = \gamma - \delta$  and taking the difference of the equations above, we obtain

$$(54) \quad e = u_s - v_s + \sum_{j=0}^{n-1} (\gamma_j (u_j - v_j) - \nu_j v_j), \quad T_a \nu = (T_a - T_e) \gamma - T(u_s - v_s).$$

Applying estimates (ii) and (iii) of Lemma 4.3 and the Banach lemma to the second equation, we obtain  $|\nu| = O(N^{-r}) \|f\|_{\omega,r}$ . Substituting this into the first equation and again using parts (i) and (iii) of Lemma 4.3, we obtain the desired result.  $\square$

We note that the estimate for the  $n$ th derivative is of optimal order for finite  $r$ . Of course, for  $f \in C^\infty([a, b])$ , we have convergence at a rate faster than any negative power of  $N$ .

**4.3. Direct computations of the condition number.** Finally, we illustrate the conditioning results by computing the singular values of the matrix used in the numerical example in §3, namely Airy's equation with a Chebyshev discretization,

$$(55) \quad A = I + \alpha^3 (x+1) D^{-2}.$$

The singular values for various  $N$  and  $\alpha$  were computed using the lapack routine, dgesvd. The results are presented in Table 2.

TABLE 2. Extreme singular values for  $I + \alpha^3(x+1)D^{-2}$ 

$N$	$\alpha = 5$			$\alpha = 10$			$\alpha = 20$		
	$\sigma_1$	$\sigma_{N-1}$	$\kappa_2$	$\sigma_1$	$\sigma_{N-1}$	$\kappa_2$	$\sigma_1$	$\sigma_{N-1}$	$\kappa_2$
32	46.3	.077	605	374	.007	53517	2992	.100	29891
64	46.3	.077	605	374	.023	16015	2992	.042	71295
128	46.3	.077	605	374	.023	16015	2992	.008	378611
256	46.3	.077	605	374	.023	16015	2992	.008	378611
512	46.3	.077	605	374	.023	16015	2992	.008	378611
1024	46.3	.077	605	374	.023	16015	2992	.008	378611

TABLE 3. Extreme singular values for  $I - \alpha D^{-4}$ 

$N$	$\alpha = 1$			$\alpha = 100$			$\alpha = 10000$		
	$\sigma_1$	$\sigma_{N-1}$	$\kappa_2$	$\sigma_1$	$\sigma_{N-1}$	$\kappa_2$	$\sigma_1$	$\sigma_{N-1}$	$\kappa_2$
32	1.00	.995	1.01	1.31	.602	2.17	69.9	.070	1004
64	1.00	.995	1.01	1.31	.602	2.17	69.9	.070	1004
128	1.00	.995	1.01	1.31	.602	2.17	69.9	.070	1004
256	1.00	.995	1.01	1.31	.602	2.17	69.9	.070	1004
512	1.00	.995	1.01	1.31	.602	2.17	69.9	.070	1004
1024	1.00	.995	1.01	1.31	.602	2.17	69.9	.070	1004

We see that the extreme singular values and, hence, the condition number of the system matrix are independent of  $N$ , once  $N$  is taken large enough to resolve the problem. (The large condition number for large  $\alpha$  simply reflects the large but bounded condition number of the integral equation.) To illustrate the insensitivity of this result to the order of the underlying differential equation, we have carried out the same computation for the biharmonic; that is for  $A = I - \alpha D^{-4}$ , with the results tabulated in Table 3.

Here, the results are quite independent of the truncation, as the extreme singular values are resolved with  $N = 32$ , so the only growth in the condition number is associated with the growth of  $\alpha$ .

## 5. RATIONAL MAPS FOR LAYER RESOLUTION

The Chebyshev approximation to a function with a region or regions of very rapid variation may exhibit Gibbs-type phenomena, that is large amplitude oscillations of the error, unless many basis functions are used. Therefore, adaptive computations using coordinate mappings to stretch these regions have been proposed. Bayliss and Turkel [3] have made a comparative study of various functional forms, all of which were transcendental functions.

In order to be able to use our fast solvers, we consider rational maps. That is, we directly solve (1) in  $y$ -space, where

$$(56) \quad x = \frac{P(y; \eta)}{Q(y; \eta)};$$

the polynomials  $P$  and  $Q$  are of low degree, and  $\eta$  is a parameter vector. In an adaptive procedure,  $\eta$  would be chosen to minimize some measure of the error, for

example the error functional proposed by Bayliss and coworkers [2]. A very simple construction of an appropriate map can be motivated in the following way: Let  $g(y) = P/Q$ . The convergence of the Chebyshev expansion (in  $y$ ) depends on the behavior of

$$(57) \quad \frac{d^k u}{dy^k} = \left( \frac{dg}{dy} \right)^k \frac{d^k u}{dx^k} + \dots$$

This suggests that improved convergence will follow from making  $dg/dy$  small where  $d^k u/dx^k$  is large. We imagine an underlying linear map (so that the limits of the computational region will be  $[-1, 1]$ ) stretched near a finite number of points,  $x_j$ . This can be accomplished by subtracting scaled and shifted multiples of the function

$$(58) \quad h_s(y; \alpha, \beta, \gamma) = \frac{\alpha y + \gamma}{1 + \beta y^2}.$$

Note that  $h'_s(0; \alpha, \beta, \gamma) = \alpha$  and that the derivative approaches zero as  $\beta y^2 \rightarrow \infty$ . We then propose

$$(59) \quad g(y) = Sy + C - \sum_j h_s(y - y_j; \alpha_j, \beta_j, \gamma_j).$$

The number of terms in the sum, and, hence, the degree of the map and bandwidth of the resulting matrices, depends on the number of layers present. Then, if  $S - \alpha_j$  is small and  $\beta_j$  is large, enhanced resolution at  $g(y_j)$  will be obtained.

We demonstrate the idea on the following simple boundary value problem:

$$(60) \quad \epsilon \frac{d^2 \phi}{dx^2} + x \frac{d\phi}{dx} = 0, \quad \phi(\pm 1) = \pm 1.$$

This problem has the exact solution

$$\phi(x) = -1 + \int_{-1}^x e^{-x^2/(2\epsilon)} dx,$$

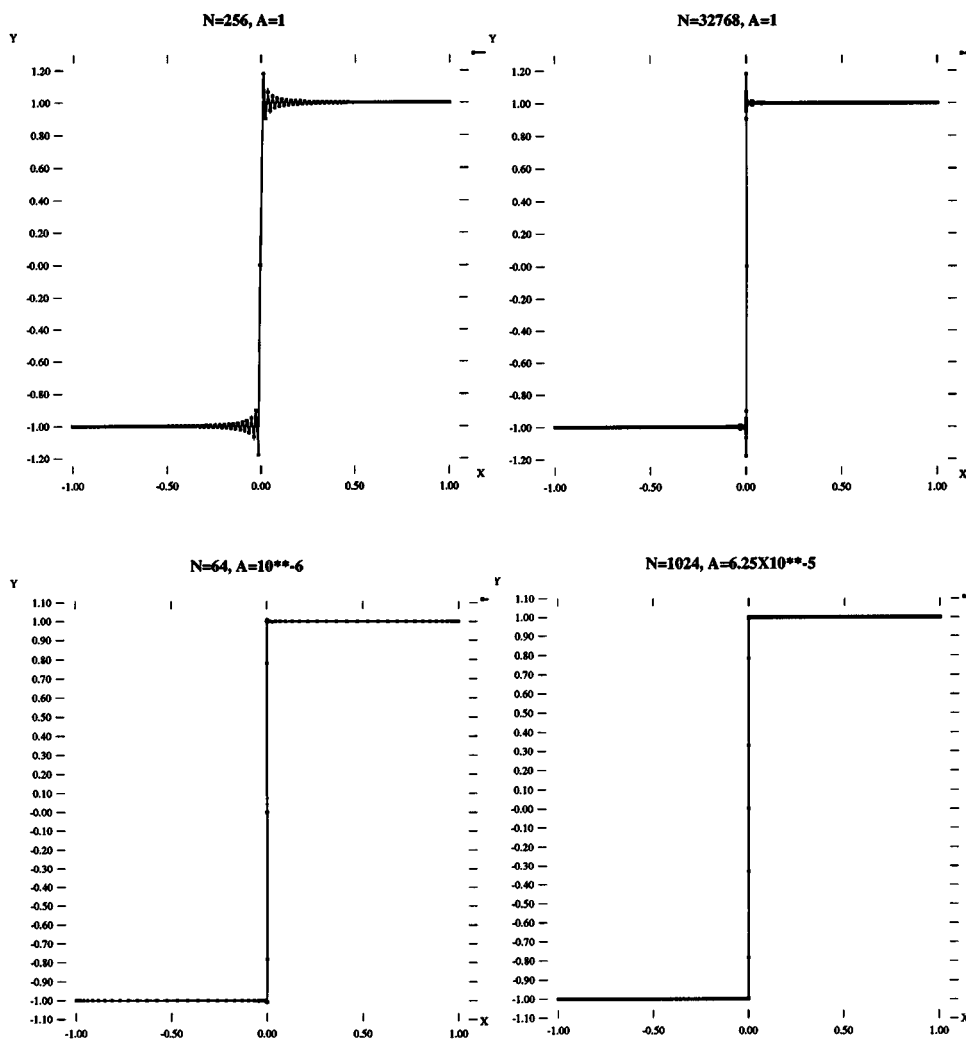
which, for  $\epsilon$  small, exhibits a region of rapid transition near  $x = 0$  whose width is  $O(\epsilon^{-1/2})$ . We see below that for  $\epsilon$  small, and even a large number of modes, there is a very strong Gibbs-like behavior. We consider the rational map

$$(61) \quad x = \frac{2}{A+1} y \frac{A+y^2}{1+y^2},$$

which is derived from the general expression above, making use of the symmetry. In particular, the derivative is minimized at  $y = 0$  where its value is  $2A/(A+1)$ . Under the change of variables,  $\psi(y) = \phi(x(y))$ , the equation becomes

$$(62) \quad \epsilon \frac{d^2 \psi}{dy^2} + \left( x \left( \frac{dx}{dy} \right) - \epsilon \left( \frac{d^2 x}{dy^2} \right) / \left( \frac{dx}{dy} \right) \right) \frac{d\psi}{dy} = 0, \quad \psi(\pm 1) = \pm 1.$$

We studied this equation for various parameter values. We found that with  $\epsilon = 10^{-12}$ , a value of the parameter  $A$  of the map of the order of  $10^{-6}$  yielded the best results. This is reasonable, as one might expect  $A = O(\sqrt{\epsilon})$  to match the scaling in the layer. We did not systematically search for the optimal value. In Figure 2 we present the solutions obtained for various numbers of modes,  $N$ , and mapping parameters,  $A$ . Note that  $A = 1$  yields the identity map, i.e., the case of standard Chebyshev approximation. With  $A = 1$  and  $N = 256$  we see oscillations near the layer with an overshoot of about 18%. Increasing  $N$  to 32768, which was the largest value considered, had *no* effect on the amplitude of the overshoot, but

FIGURE 2. Solutions for  $\varepsilon = 10^{-12}$ 

did contract the region of oscillation. With  $A \neq 1$ , on the other hand, we obtain reasonable results with many fewer nodes. For example, with  $A = 10^{-6}$  and  $N = 64$  the overshoot is less than 1%. Increasing  $N$  to 1024 and spending a bit more effort optimizing the parameter reduces this to  $3 \times 10^{-4}$ .

It is clear that there is no change in the essential ( $O(N)$ ) amount of work needed to solve the problem, although the bandwidth does increase by approximately a factor of 5. If an iteration were used for the minimization of some error functional, by shifting the position of the shock and changing the magnification factors, each step would require the recomputation of the operator coefficients and the solution of the problem. These procedures are of comparable numerical cost, so that the desirable features of the method are essentially preserved under the change of variables.

It is worth noting that there are limits to the capability of this method to concentrate a large fraction of the mesh in a small region. A simple calculation indicates



that  $g'(y) = O(\epsilon)$  in a  $y$ -interval of width  $\sqrt{\epsilon}$ . To achieve a greater magnification, one must use rational maps of higher degree, which results in larger bandwidths. A more detailed study of the properties of rational coordinate mappings is planned for the future.

*Note added in proof.* The three-term recurrence relation for the derivatives of the Jacobi polynomials appears also in the recent review by Fornberg [10]. We would like to thank one of the referees for pointing out the paper by Bernardi and Maday [4], which allowed us to improve our convergence estimate for Gegenbauer polynomials.

## REFERENCES

1. M. Abramowitz and I. Stegun, *Handbook of Mathematical Functions*, Dover, New York, (1965). MR **29**:4914
2. A. Bayliss, D. Gottlieb, B. Matkowsky and M. Minkoff, An adaptive pseudo-spectral method for reaction diffusion problems, *J. Comp. Phys.*, 81, (1989), 421-443.
3. A. Bayliss and E. Turkel, Mappings and accuracy for Chebyshev pseudo-spectral computations, *J. Comp. Phys.*, 101, (1992), 349-359. CMP 92:15
4. C. Bernardi and Y. Maday, Polynomial interpolation results in Sobolev spaces, *J. Comput. and Appl. Math.*, 43, 53-82, 1992. MR **93k**:65010
5. C. Canuto, M. Hussaini, A. Quarteroni and T. Zang, *Spectral Methods in Fluid Dynamics*, Springer Verlag, Berlin, (1988). MR **89m**:76004
6. C.W. Clenshaw, The numerical solution of linear differential equations in Chebyshev series, *Proc. Camb. Phil. Soc.*, 53, 134-149, (1957). MR **18**:516a
7. E. Coddington and N. Levinson, *Theory of Ordinary Differential Equations*, McGraw-Hill, New York, (1955). MR **16**:1022b
8. E.A. Coutsias, T. Hagstrom, J.S. Hesthaven and D. Torres, Integration preconditioners for differential operators in spectral  $\tau$ -methods, 1995 (preprint).
9. E.A. Coutsias, F.R. Hansen, T. Huld, G. Knorr and J.P. Lynov, Spectral Methods for Numerical Plasma Simulations, *Phys. Scripta*, 40, 270-279, 1989.
10. B. Fornberg, A review of pseudospectral methods for solving partial differential equations, *Acta Numerica* (1994), 203-267. CMP 94:16
11. L. Fox and I.B. Parker, *Chebyshev Polynomials in Numerical Analysis*, Oxford Univ. Press, London, (1968). MR **37**:3733
12. D. Gottlieb and S. Orszag, *Numerical Analysis of Spectral Methods: Theory and Applications*, SIAM, Philadelphia, (1977). MR **58**:24983
13. L. Greengard, Spectral integration and two-point boundary value problems, *SIAM J. Numer. Anal.*, 28, (1991), 1071-1080. MR **92H**:65033
14. H. Hochstadt, *The Functions of Mathematical Physics*, Dover, New York (1986). MR **88b**:33001
15. C. Lanczos, *Applied Analysis*, Prentice-Hall, London, (1956). MR **18**:823c
16. W. Magnus and F. Oberhettinger, *Formulas and Theorems for the Special Functions of Mathematical Physics*, Chelsea, New York (1949). MR **8**:532b
17. L. Trefethen and M. Trummer, An instability phenomenon in spectral methods, *SIAM J. Numer. Anal.*, 24, (1987), 1008-1023. MR **89a**:65139
18. L.S. Tuckerman, Transformations of matrices into banded form, *J. Comp. Phys.*, 84, 360-376 (1989). MR **91g**:65082
19. N. N. Lebedev, *Special Functions and their Applications*, Dover, New York, 1972. MR **50**:2568

DEPARTMENT OF MATHEMATICS AND STATISTICS, UNIVERSITY OF NEW MEXICO, ALBUQUERQUE, NEW MEXICO 87131

*E-mail address:* [vageli@math.unm.edu](mailto:vageli@math.unm.edu)

*E-mail address:* [hagstrom@math.unm.edu](mailto:hagstrom@math.unm.edu)

*E-mail address:* [dtorres@math.unm.edu](mailto:dtorres@math.unm.edu)