

# A Critical Evaluation of the Resolution Properties of B-Spline and Compact Finite Difference Methods

Wai Yip Kwok,\* Robert D. Moser,\* and Javier Jiménez†

\**Department of Theoretical and Applied Mechanics, University of Illinois at Urbana–Champaign, Urbana, Illinois 61801; and †School of Aeronautics, Universidad Politécnica, 28040 Madrid, Spain*

Received July 2, 2000; revised March 7, 2001

---

Resolution properties of B-spline and compact finite difference schemes are compared using Fourier analysis in periodic domains, and tests based on solution of the wave and heat equations in finite domains, with uniform and nonuniform grids. Results show that compact finite difference schemes have a higher convergence rate and in some cases better resolution. However, B-spline schemes have a more straightforward and robust formulation, particularly near boundaries on nonuniform meshes. © 2001 Elsevier Science

*Key Words:* resolution; B-spline methods; compact finite difference methods.

---

## 1. INTRODUCTION

Many physical phenomena involve a broad range of spatial scales. One example is turbulent fluid flows, which have a wide and continuous spectrum of length scales describing its composition of eddies of different sizes [2]. Simulation of these physical phenomena requires spatial discretization schemes with high resolution, or in other words, schemes that can produce accurate numerical results over as broad a range of length scales as possible for a given discretization.

In numerical simulation of turbulent fluid flows, spectral methods are attractive spatial discretization schemes due to their very good resolution properties. As a result, many direct numerical simulations (DNS) have been performed with spectral methods in Cartesian coordinates with various boundary conditions [4, 14]. These include simulations of simple fundamental flows such as isotropic turbulence, turbulent channel flows [20], and turbulent boundary layers [35]. One distinctive feature of spectral methods is that they use infinitely differentiable global basis functions [4]. Two common choices are Fourier series expansions and polynomial basis functions, with the first being applied to simulations with periodic boundary conditions and the second to simulations in finite intervals [20, 30]. However, the

global character of the basis functions also limits spectral methods to simple geometries and boundary conditions [24], and there is a great need for simulations in complex geometries. This is very important if turbulence simulations are to contribute to many engineering applications such as external aerodynamics and propulsion systems. Such simulations would require spatial discretization schemes that not only retain the good resolution properties of spectral methods, but also provide flexibility with respect to geometries and grid distribution.

Local numerical representations, such as finite difference and finite element schemes, have much greater flexibility in discretizing complex geometries, so high resolution schemes of these types would be of great interest. For example, Lele has studied compact finite difference schemes for use in problems with a broad range of spatial scales [24], using Fourier analysis to investigate how well the schemes represent a range of wavenumbers. There has also been a trend to combine local discretization algorithms and spectral methods. A typical example of such a confluence of numerical algorithms is the spectral element method, which is based on finite element and spectral methods [18, 19, 29].

Another choice for local numerical representation is to use splines. Unlike finite difference methods, spline methods are functional expansion methods that make use of a set of local basis functions. This property provides us with a straightforward way to implement boundary conditions. Spline methods are similar to finite element methods as they both use piecewise polynomial representations. However, spline methods use basis functions that retain a higher degree of continuity. In short, spline methods have much of the flexibility afforded by the use of local expansions, as in finite elements, and have the resolution advantage afforded by highly continuous expansions, as in spectral methods.

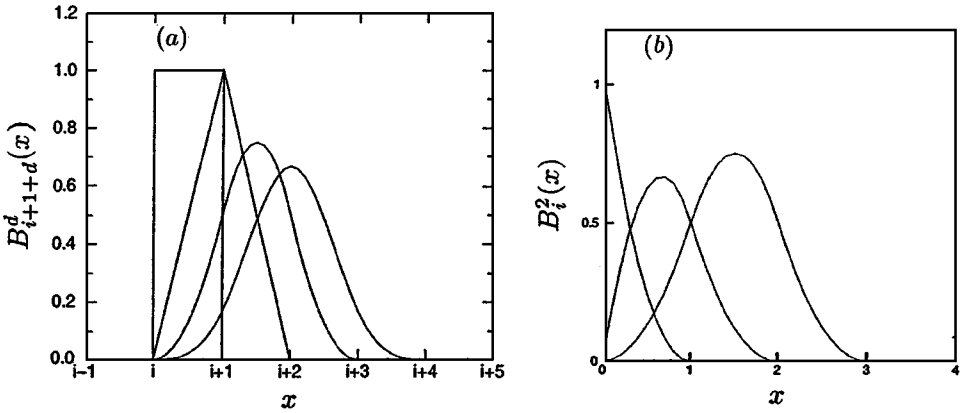
In the research reported here, we investigate the properties of spline methods, in particular spline collocation methods, and their relation to finite difference and finite element methods. Section 2 introduces the basic properties of spline, compact finite difference, finite element methods, and their different formulations. The basic resolution properties of these spatial discretization schemes are presented in Section 3 using Fourier analysis in periodic domains. Of particular interest are the approximations to the first and second derivative operator, which are common in equations describing many physical phenomena. In Sections 4 and 5, the first-order wave equation and heat equation are solved with spline collocation and compact finite difference schemes in bounded domains, in both uniform and nonuniform grids. Concluding remarks are given in Section 6.

## 2. NUMERICAL REPRESENTATIONS

The resolution properties of the numerical methods discussed here are most easily understood in one spatial dimension. Thus, the methods to be evaluated are introduced here in their one-dimensional form. Spline methods, compact finite difference methods, and finite element methods will be discussed.

### 2.1. Spline Methods

Consider a domain divided into  $N$  intervals, a one-dimensional spline is defined to be a polynomial of degree  $d$  in each interval that is continuously differentiable  $d - 1$  times at the interval boundaries. The boundaries of the intervals are called knots.



**FIG. 1.** B-splines  $B_i^d(x)$  on uniform knots with knot spacing  $\Delta x = 1$ , with (a)  $d = 0, 1, 2$ , and  $3$ ; and with (b)  $d = 2$  and  $i = 1, 2$  and  $3$ .

Spline methods have been used before to solve differential equations and fluid mechanics problems [13, 31, 32]. The work of Kasi Viswanadham and Koneru [39] and Davies [6, 7] used B-splines as basis functions and the Galerkin formulation. Most of the research, however, is confined to cubic splines ( $d = 3$ ). More recently, Kravchenko *et al.* [22] and Shariff and Moser [34] used the basis functions of splines to solve partial differential equations and simulate turbulent fluid flows. In particular, mesh embedding techniques are developed to make basis spline methods very effective in solving physical problems in complex geometries.

To use splines as a representation for the solution of a partial differential equation, it is necessary to have a convenient basis for the space of spline functions under consideration. Here the so-called basis splines or “B-splines” as described in [8] and [16] are used. A B-spline is defined as a normalized spline which has support over the minimum possible number of intervals. In fact, it has support on only  $d + 1$  intervals. As an example, B-splines for uniformly spaced knots are plotted in Fig. 1a for  $d$  up to 3. By using a basis with support on the minimum possible number of intervals, minimum bandwidth of the resulting matrices is ensured.

Near a boundary, the basis splines are different than those in Fig. 1a since the presence of the boundary removes the constraint that the B-splines have  $d - 1$  zero derivatives at the edge of its interval of supports. An example of the quadratic B-splines near the boundary is shown in Fig. 1b.

To use the B-splines in a practical computation, one needs to evaluate them and their derivatives at points in the domain. This will be sufficient to compute the various matrices representing different linear operators. An efficient and stable technique to evaluate the B-splines and their derivatives is the recurrence relation described in [8] (see Appendix A). Both interior and boundary splines are generated this way by formally introducing a multiplicity of knots at the boundary (see [8] and Appendix A).

Consider the B-spline representation of a possibly nonlinear spatial operator  $\mathcal{F}$  operating on  $\phi$ . We first postulate an expansion for  $\phi$  in terms of B-splines of order  $d$  on a selected knot set:

$$\phi(x) \approx \tilde{\phi}(x) = \sum_i \alpha_i B_i^d(x). \quad (1)$$

An approximation  $\tilde{\mathcal{F}}$  to the operator  $\mathcal{F}$  is sought that maps splines in  $S_d$  (i.e.,  $\tilde{\phi}$ ) to splines in  $S_d$ , where  $S_d$  is the space of splines of order  $d$  for the selected knot set. That is

$$\mathcal{F}(\tilde{\phi}) \approx \tilde{\mathcal{F}}(\tilde{\phi}) = \tilde{\gamma} = \sum_i \beta_i B_i^d(x). \quad (2)$$

There are several ways to generate such an approximation. Two will be considered here, namely Galerkin and collocation methods.

### 2.1.1. B-spline Galerkin Methods

In the Galerkin formulation, the approximation of the linear differential operator  $\mathcal{D}$  on  $\tilde{\phi}$  is given by

$$(B_j^d, \tilde{\gamma}) = (B_j^d, \mathcal{D}\tilde{\phi}), \quad j = 1, 2, \dots, N_\zeta, \quad (3)$$

where  $(f, g)$  denotes the  $L_2$  inner product  $\int fg \, dx$  in the domain and  $N_\zeta$  is the number of B-splines. This forces the error in  $\tilde{\gamma}$  to be orthogonal to  $S_d$ , thus minimizing the  $L_2$  error in this space. Given the linearity of  $\mathcal{D}$  and the representations of  $\tilde{D}$ ,  $\tilde{\phi}$ , and  $\tilde{\gamma}$ , the above equation can be written

$$\sum_{i=1}^{N_\zeta} \beta_i (B_j^d, B_i^d) = \sum_{i=1}^{N_\zeta} \alpha_i (B_j^d, \mathcal{D}(B_i^d)), \quad j = 1, 2, \dots, N_\zeta. \quad (4)$$

The inner products of Eq. (4) are the elements of matrices  $M$  and  $D$ , with  $M_{ij} = (B_j^d, B_i^d)$  and  $D_{ij} = (B_j^d, \mathcal{D}(B_i^d))$ . The matrix  $M$  is called the ‘‘mass’’ matrix and  $D$  the operator matrix. To obtain  $\tilde{\gamma}$  given  $\tilde{\phi}$ , one solves the linear system  $M\beta = D\alpha$ . Note that both  $M$  and  $D$  are banded matrices since individual B-splines have only local support. The bandwidth  $w$  of the matrices is given by  $w = 2d + 1$ .

### 2.1.2. B-Spline Collocation Methods

The collocation formulation imposes different requirements to obtain the coefficients  $\beta_i$ . Here the approximation  $\tilde{\gamma} = \sum_i \beta_i B_i^d$  of the operator  $\mathcal{D}$  on  $\tilde{\phi}$  must satisfy

$$\tilde{\gamma} = \mathcal{D}\tilde{\phi} \quad \text{at } x = \zeta_j, \quad j = 1, 2, \dots, N_\zeta, \quad (5)$$

which implies

$$\sum_{i=1}^{N_\zeta} \beta_i B_i^d = \sum_{i=1}^{N_\zeta} \alpha_i \mathcal{D}(B_i^d) \quad \text{at } x = \zeta_j, \quad j = 1, 2, \dots, N_\zeta. \quad (6)$$

The values of the B-splines and their derivatives are the elements of the matrices  $M$  and  $D$ , respectively, with  $M_{ij} = B_i^d(\zeta_j)$  and  $D_{ij} = [\mathcal{D}(B_i^d)](\zeta_j)$ . Again, given  $\tilde{\phi}$ ,  $\tilde{\gamma}$  is found by solving the linear system  $M\beta = D\alpha$ . Using the collocation formulation, the matrix bandwidth  $w$  is given by  $w = d$  for odd  $d$ .

### 2.1.3. Selection of Knots and Collocation Points

To use B-splines in a computation, one first needs to determine the location of the knot points and for the collocation method the collocation points. In a periodic domain with  $N$  uniform width intervals, there are  $N$  knots and  $N$  splines spanning the spline space. Therefore,  $N$  collocation points are needed in a collocation scheme. There are only two locations for the collocation points that preserve the spatial symmetry of the operators: collocation points at the knots and collocation points at the center of the intervals. The former is appropriate for odd-order splines, the latter for even.

In a nonperiodic domain, it is more complicated. There are  $N$  intervals,  $N + 1$  distinct knots and  $N + d$  collocation points are needed. There are two basic ways to select knots and collocation points in a finite domain. The first is that  $N + d$  collocation points can be selected by whatever resolution criteria are appropriate and then  $N + 1$  of these points can be chosen to be the knots. Generally, those collocation points that are not knots are near the boundary, though the knot at the boundary is retained. This is referred to as a “not-a-knot” condition, which is commonly used in spline interpolation.

The alternative is to start by selecting the knots according to some resolution criteria. This is more natural since the knots directly determine the spline space and therefore are more closely related to resolution than the collocation points. Furthermore, in a Galerkin scheme all one selects are the knots, so direct comparison of Galerkin and collocation is only possible if one starts by selecting the knots. Selecting the collocation points can then be done in several ways, but there are two choices that seem particularly appropriate: place a collocation point at the maximum of each B-spline function or place it at the centroid of each B-spline function. These prescriptions have the advantage that they are applicable throughout the domain (nothing special about the boundary), and they associate a collocation point directly with each B-spline function. This latter property is useful for applications in multidimensional embedded grids of the type described by Shariff and Moser [34]. Note that with uniform knots away from the boundary, the symmetry of the B-splines places the maxima and centroid at the same location: at the knots or at the center of the intervals for odd and even splines, respectively. In the current paper, collocation points at the B-spline maxima are selected, because this naturally places a collocation point at the boundary, which is useful for imposing boundary conditions. Two knot distributions are used: uniformly spaced knots and nonuniform knots distribution according to

$$x = 0.5 \left\{ 1 - \frac{\cos \left[ \pi \frac{(N-1)\xi + 1}{N+1} \right]}{\cos \left[ \pi \frac{1}{N+1} \right]} \right\} \quad (7)$$

where  $\xi = j/N$  for  $j = 0, 1, \dots, N$ . This nonuniform grid is basically a Chebyshev grid with the boundary singularities removed. It is denser near the boundary.

## 2.2. Compact Finite Difference Methods

Compact finite difference schemes have long been applied to fluid mechanics and other physics problems [17, 23, 33]. Recently, higher order compact finite difference schemes have seen increasing use in the direct numerical simulation of complex fluid flows [12, 28]. Lele presented a comprehensive study on the compact finite difference methods [24]. Consider a uniform mesh where the nodes are indexed by  $i$ . The independent variable

at the nodes is  $x_i$  and the function values at the nodes  $v_i = v(x_i)$  are given. The compact schemes are derived by writing approximations of the form:

$$\beta v'_{i-2} + \alpha v'_{i-1} + v'_i + \alpha v'_{i+1} + \beta v'_{i+2} = c \frac{v_{i+3} - v_{i-3}}{6\Delta x} + b \frac{v_{i+2} - v_{i-2}}{4\Delta x} + a \frac{v_{i+1} - v_{i-1}}{2\Delta x}. \quad (8)$$

Similarly, approximations to the second derivative operator are derived by the following relationship:

$$\begin{aligned} & \beta v''_{i-2} + \alpha v''_{i-1} + v''_i + \alpha v''_{i+1} + \beta v''_{i+2} \\ & = c \frac{v_{i+3} - 2v_i + v_{i-3}}{9(\Delta x)^2} + b \frac{v_{i+2} - 2v_i + v_{i-2}}{4(\Delta x)^2} + a \frac{v_{i+1} - 2v_i + v_{i-1}}{(\Delta x)^2}. \end{aligned} \quad (9)$$

The relations between the coefficients  $a$ ,  $b$ ,  $c$ , and  $\alpha$ ,  $\beta$  are obtained by matching the Taylor series coefficients of various orders. Higher orders can be obtained by including more nodes in the above two equations.

In this study, compact schemes with the same stencil size on both sides of the equations are selected ( $c = 0$  in Eqs. (8) and (9) for example). This is because mass and operator matrices with the same bandwidth is a property shared by B-spline methods. All the coefficients then are used to match the Taylor series to as high an order as possible. The value of the coefficients are listed in Tables I and II for schemes with matrix bandwidth  $w$  up to 11. Schemes in which convergence order is sacrificed to improve resolution have also been proposed (see [24] and Table I). Note that since no restriction is imposed on the coefficients other than those from Taylor series matching, mass matrices associated with the first, second, and higher derivatives are all different. This issue will be addressed in more detail in Section 6.

### 2.3. Finite Element Methods

Most of the finite element applications in fluid dynamics use the Galerkin finite element formulation [11]. The application of finite element method to fluid mechanics is treated by Thomasset [37] and Baker [1].

In this study, one-dimensional finite elements with  $C_0$  and  $C_{(d-1)/2}$  continuity are used, where  $d$  is the degree of polynomials.  $C_{(d-1)/2}$  continuity is the highest that can be imposed while preserving the iso-parametric property of the elements. These are commonly called Hermite finite elements. As with B-splines, the finite elements are polynomials on a series of knots (element boundaries). However, because a lower order of continuity is imposed, there are many more degrees of freedom per interval (element). If there are  $N$  intervals, then there would be  $dN$  and  $\frac{d+1}{2}N$  degrees of freedom for  $C_0$  and  $C_{(d-1)/2}$  finite elements, respectively. In this paper, only finite element Galerkin methods are considered, though collocation methods are also possible. Note that this method of increasing the local degree of the polynomial shape-function is very similar to the “ $p$ ” finite element method [10], in which an element may neighbor an element having different polynomial order. The main advantage of finite element methods with low order of continuity is flexibility with respect to geometry. In most applications of finite element methods, elements are typically chosen to be at most quadratic [3, 9], and consequently, a high order of convergence is not achieved. This is exactly opposite to the characteristics of spectral methods. The intention to

**TABLE I**  
**Table of Coefficients for Discretized First Derivative Operators Using Compact Finite Difference Schemes**

Band-width	Order	$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_5$
3	2	2.822510141559E-01	0	0	0	0
3	4	$\frac{1}{4}$	0	0	0	0
5	6	4.907480792180E-01	3.935368647117E-02	0	0	0
5	8	$\frac{4}{9}$	$\frac{1}{36}$	0	0	0
7	12	$\frac{9}{16}$	$\frac{9}{100}$	$\frac{1}{400}$	0	0
9	16	$\frac{16}{25}$	$\frac{4}{25}$	$\frac{16}{1225}$	$\frac{1}{4900}$	0
11	20	$\frac{25}{36}$	$\frac{100}{441}$	$\frac{25}{784}$	$\frac{25}{15876}$	$\frac{1}{63504}$
Band-width	Order	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$
3	2	1.564502028312E+00	0	0	0	0
3	4	$\frac{3}{2}$	0	0	0	0
5	6	1.450612391632E+00	6.09591139745939E-01	0	0	0
5	8	$\frac{40}{27}$	$\frac{25}{54}$	0	0	0
7	12	$\frac{21}{16}$	$\frac{231}{250}$	$\frac{147}{2000}$	0	0
9	16	$\frac{144}{125}$	$\frac{152}{125}$	$\frac{10704}{42875}$	$\frac{761}{85750}$	0
11	20	$\frac{55}{54}$	$\frac{12760}{9261}$	$\frac{5115}{10976}$	$\frac{23045}{500094}$	$\frac{7381}{8001504}$

*Note.* The approximations have the form  $\sum_j \alpha_j (f'_{i+j} + f'_{i-j}) + f'_i = \sum_j a_j \frac{f_{i+j} - f_{i-j}}{2j\Delta x}$ . For the second-order tri-diagonal and sixth-order pentadiagonal schemes, coefficients are chosen to increase 1% resolution (see Section 3.2).

combine these two methods comprehensively leads to the development of spectral element methods [19, 29]. Spectral element methods are basically variational domain decomposition techniques. The computational domain is broken up into macro-elements within which variables are represented as high-order polynomial expansions [18]. The work of Patera [29], Karniadakis [18], and their co-workers illustrates the application of spectral element methods in partial differential equations and fluid mechanics problems.

#### 2.4. Relationship to Other Approximations

There are a variety of other formulations for the first and second derivative (or equivalent) that have not been covered here. Two that are of particular interest, due to their close relationship to the compact finite difference methods, are the Padé finite volume methods discussed by Kobayashi [21] and the coupled derivative formulation of Mahesh [25]. In the Padé finite volume scheme, a compact reconstruction operator is used to reconstruct the values of the function being approximated on the volume boundaries, given the cell averaged values of the functions. These are then the fluxes in a finite volume representation of the convection equation. A similar reconstruction of the derivative at the volume

**TABLE II**  
**Table of Coefficients for Discretized Second Derivative Operators Using Compact Finite Difference Schemes**

Band-width	Order	$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_5$
3	4	1	0	0	0	0
5	8	$\frac{344}{1179}$	$\frac{23}{2358}$	0	0	0
7	12	$\frac{329913}{725308}$	$\frac{18387}{362654}$	$\frac{619}{725308}$	0	0
9	16	5.701357323754E-1	1.127421221433E-1	6.744355726188E-3	6.859460384736E-5	0
11	20	6.483654835488E-1	1.803832051742E-1	2.035363636527E-2	7.637676152837E-4	5.211703267310E-6
Band-width	Order	$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$\alpha_5$
3	4	$\frac{6}{5}$	0	0	0	0
5	8	$\frac{320}{393}$	$\frac{310}{393}$	0	0	0
7	12	$\frac{263655}{725308}$	$\frac{261954}{181327}$	148449	0	0
9	16	7.963758598008E-2	1.60567370122	6.589603271777E-1	3.511632641751E-2	0
11	20	-7.028592427964E-2	1.472309342499	1.120030052713	1.729190106287E-1	4.770127252614E-3

Note. The approximations have the form  $\sum_j \alpha_j (f_{i+j}'' + f_{i-j}'') + f_i'' = \sum_j \alpha_j \frac{i_{i+j}^{-2} f_{i+j}}{j^2 h^2}$ .



boundary yields a compact finite volume representation of the diffusion equation. However, when the relevant reconstruction operator is integrated into the convection or diffusion equation, the same overall operator as the corresponding compact finite difference approximation is obtained. Thus, the error properties of the compact finite differences described above are equally applicable to the compact finite volume schemes of Kobayashi [21]. This is true, however, only for the infinite or periodic domain problem. The near-boundary schemes appropriate for the finite volume representation are different. Kobayashi proposes a fourth-order boundary scheme for use with the fourth-order Padé finite volume method, but no boundary treatments for use with higher order finite volume representations are reported.

In the coupled derivative (CD) formulation of Mahesh [25], one takes advantage of the fact that in many problems one has both a first and second derivative (convection and diffusion). By computing them together, one is able to obtain a higher order approximation to each (with the same stencil size) than would be possible by computing them separately. In addition, even when comparing schemes of the same order, the CD methods have somewhat better resolution properties than the standard compact finite difference methods. For example the sixth-order CD first derivative approximation has an error approximately a factor of 3 smaller than sixth-order Padé approximation for wavenumbers less than approximately  $k_{\max}/2$ . For the second derivative, the effective wavenumber has smaller error at large wavenumber ( $k > k_{\max}/2$ ). The matrices needed to implement the CD method have larger bandwidth, with the result that the computational cost is slightly higher than the standard compact finite difference methods of the same order [25]. In finite domains, the stable boundary schemes investigated by Mahesh were third order for the first derivative and fourth order or less for the second derivative. The effects of the reduced order boundary schemes on resolution properties of finite domain problems is of concern.

## 2.5. Basis for Comparison

To compare the resolution properties of the several spatial discretization schemes discussed above, it is necessary to define the basis of comparison. The question is: comparing B-spline, finite element, and finite difference methods, what characteristics of these methods (i.e., what degree polynomials, or what stencil size) should be compared. In this paper we take the view that comparison should be done between schemes with matrices that have the same bandwidth. The bandwidth of the matrices is an indicator of the computational cost of performing the linear algebra associated with the scheme, so methods with similar linear algebra costs are compared. This is related to the common practice of characterizing finite difference methods by their stencil size.

There are several reasons that a comparison based on bandwidth is appropriate in the current context. First, cost is an important consideration and the linear algebra cost for which bandwidth is an indicator often dominates the computational cost in the numerical solution of PDEs. Second, commonly used bases of comparison, such as polynomial degree, are not defined for all methods, or, as with order of accuracy or convergence, may have different common interpretations in different methods (e.g., order of derivative approximation in finite difference versus order of function approximation in finite elements). Finally, it is not always clear which of several accuracy indicators are of most interest, and so should form the basis of comparison. By making bandwidth (cost) the basis of comparison, one can more conveniently assess the relative merit of disparate schemes by a variety of measures.

Of course, there are many computational costs that are associated with the numerical solution of any given problem, not all of which scale with the bandwidth. These will vary with the details of the problem being solved. So, any general cost based comparison like this is inherently imperfect. Nonetheless, when comparing disparate schemes, some basis for comparison based on cost is appropriate, and matrix bandwidth is the best indicator of the relative computational complexity of these schemes that we were able to devise.

### 3. FOURIER ANALYSIS

In this section, a Fourier analysis of the errors associated with the approximation of differential operators by the several spatial discretization schemes discussed in Section 2 is presented. The resolution properties of the numerical schemes are most directly investigated using a Fourier analysis [24, 26, 27, 36, 38], in which the approximations of the operators in a periodic or infinite domain with a uniform grid are compared.

#### 3.1. Effective Wavenumber and Eigenfunctions

One common measure of how well a differential operator is approximated is the effective wavenumber. In a periodic or infinite domain, the eigenfunctions of derivative operators are the complex exponentials, and the eigenvalues of the  $n$ th derivative are  $(ik)^n$ , where  $k$  is the wavenumber of the complex exponential and  $i = \sqrt{-1}$ . The effective wavenumbers  $\tilde{k}$  are obtained from the eigenvalues of the approximate derivative operators  $M^{-1}D$  as  $\tilde{k}_j = \sqrt[n]{\frac{\lambda_j}{i^n}}$ , where  $\lambda_j$  is the  $j$ th eigenvalue of the approximate operators. For central schemes such as those studied in this section,  $\tilde{k}$  is real. Perhaps more important than the effective wavenumber is the error in the eigenvalue  $|\lambda - (ik)^n|$ . Also of interest is how closely the eigenfunctions of the approximate operator correspond with the exact eigenfunctions (the complex exponentials).

While the effective wavenumber has been widely studied as an indicator of the accuracy and resolution of approximate derivative operators, the accuracy with which the eigenfunctions of the operators (the complex exponentials) are represented has not generally been considered. The accuracy of the eigenfunctions of the approximate operators is important because they are a necessary part of the description of the operators (accurate eigenvalues is not enough). The accuracy of the eigenfunctions is clearly related to the ability to approximate the complex exponential, which is also important. In finite element and B-spline methods, the error in a numerical solution of a problem is related to and is certainly limited by this approximation error.

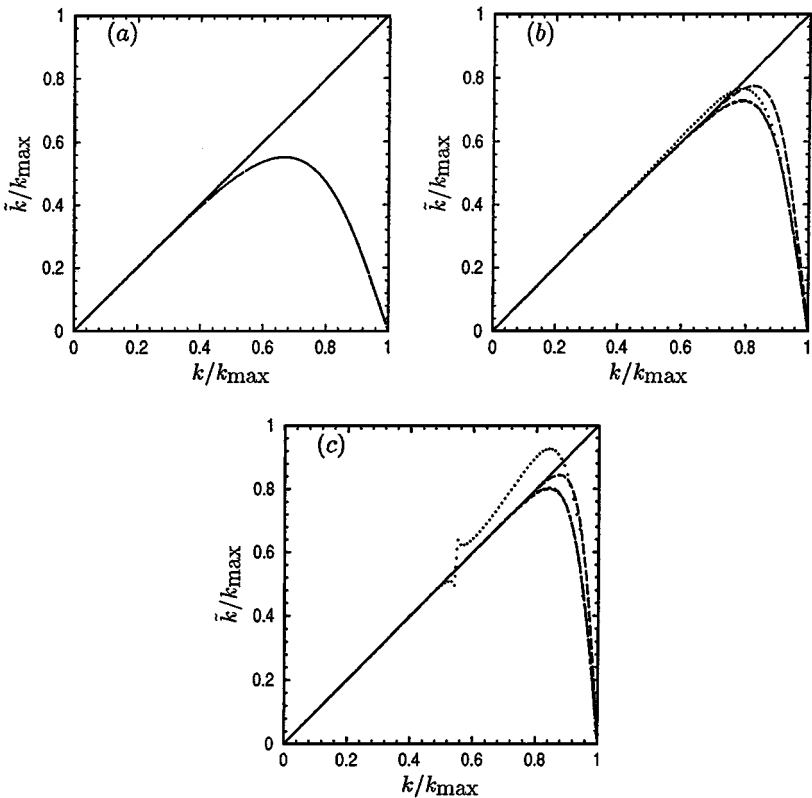
One reason that eigenfunctions have been less often examined is that in finite difference methods, the circulant nature of the operator matrices ensures that the eigenfunctions of the operators exactly represent the values of  $e^{ikx}$  at the finite difference grid points. However, with methods based on functional representations, one can measure the  $L_2$  errors  $\|e^{ikx} - \psi_j(x)\|$ , where  $\psi_j(x)$  are the approximate eigenfunctions. For the B-spline schemes, the matrices  $M$  and  $D$  are circulant. Therefore, the approximate operator has the same eigenfunctions for all derivatives. These approximate eigenfunctions are also the same as those obtained by directly approximating the complex exponential, using the method under considerations (Galerkin or collocation).

The B-spline matrices are circulant because with uniform knots, the basis functions are all the same, only differing by a spatial shift. For high order finite elements, however,

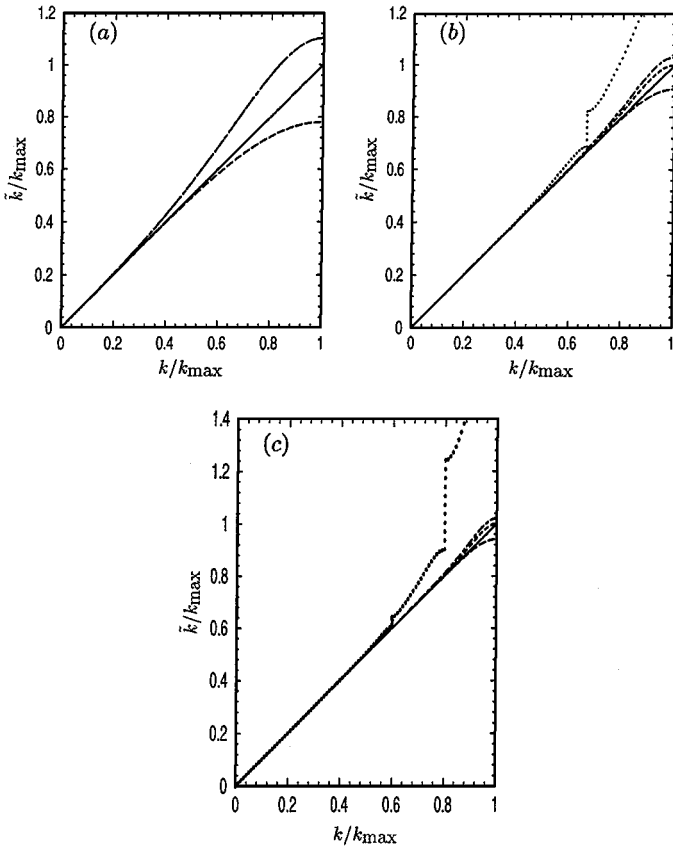
there are several different basis functions, so the matrices are not circulant. But they are block circulant, which allows the eigenvalues and eigenfunctions to be easily determined (see Appendix B). Since the matrices are not circulant, the eigenfunctions of different derivatives are not the same, and they are different from the direct approximation of the complex exponential. However, there is very little difference between the eigenfunctions and representation of the complex exponential, which has a slightly lower error. Therefore, the error in the direct Galerkin finite element approximation will be presented in the following sections.

### 3.2. Comparison of Accuracy and Resolution

The numerical schemes tested using Fourier analysis include B-spline collocation and Galerkin formulations, finite element Galerkin formulations, and compact finite difference methods. Effective wavenumbers associated with the first and second derivatives for the four methods discussed here are shown in Figs. 2 and 3, respectively. Notice that the wavenumber is normalized by the maximum wavenumber  $k_{\max}$ , representable with the numerical method. For the spline and finite difference methods,  $k_{\max} = \frac{\pi}{\Delta x}$ , where  $\Delta x$  is the grid or knot spacing. For the  $C_0$  or  $C_{(d-1)/2}$  finite element schemes,  $k_{\max} = \frac{d\pi}{\Delta x}$  or  $\frac{(d+1)\pi}{2\Delta x}$ , respectively,



**FIG. 2.** Effective wavenumber  $\bar{k}$  of the first derivative operators for matrix bandwidth (a) 3, (b) 7, (c) 11: ----, B-spline; ---, compact finite difference; ·····,  $C_0$  finite element Galerkin; -·-·-,  $C_{(d-1)/2}$  finite element Galerkin; —, exact differentiation. For bandwidth equals 3, the three schemes yield the same result. For bandwidth equals 7 and 11, the difference between B-spline and  $C_{(d-1)/2}$  finite element Galerkin schemes is indistinguishable.

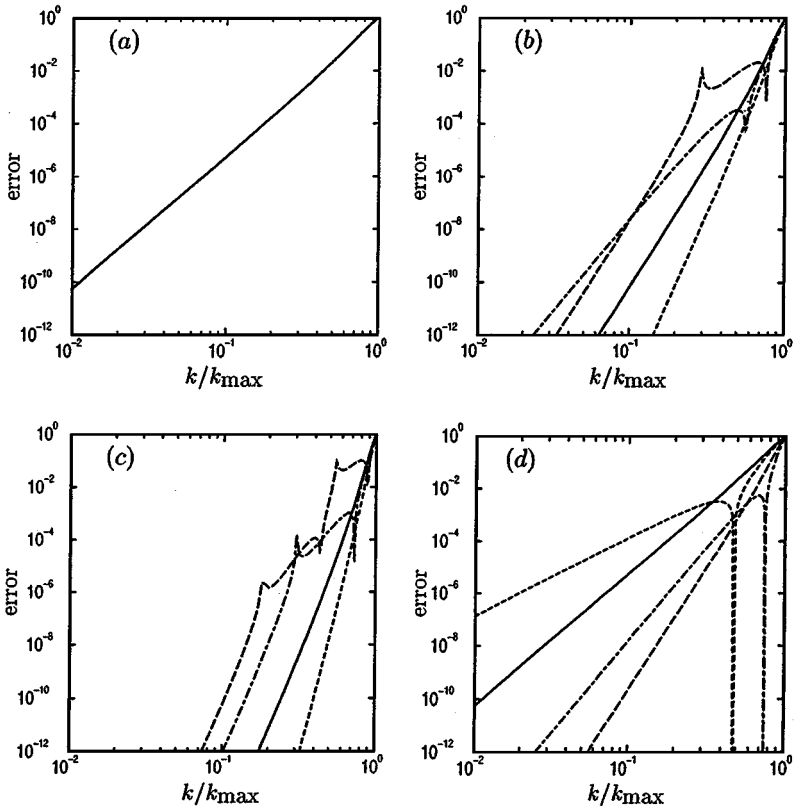


**FIG. 3.** Effective wavenumber  $\tilde{k}$  of the second derivative operators for matrix bandwidth (a) 3, (b) 7, (c) 11: ----, B-spline; ---, compact finite difference; ·····,  $C_0$  finite element Galerkin; -·-·-,  $C_{(d-1)/2}$  finite element Galerkin; —, exact differentiation. For bandwidth equals 3, B-spline and finite element yields the same result.

since there are  $d$  or  $(d + 1)/2$  degrees of freedom per element. This definition of  $k_{\max}$  is appropriate for finite elements, since the number of Fourier modes that can be represented and the size of the calculation are determined by the number of degrees of freedom in the representation, not by the size of the elements.

There are several things to note about the effective wavenumbers. First, for a given matrix bandwidth  $w$ ,  $\tilde{k}$  is identical for B-spline collocation and Galerkin methods. This is despite the fact that for collocation, the order of the splines is higher ( $d = w$ ) than for the Galerkin ( $d = \frac{w-1}{2}$ ). This identity was noted by Swartz and Wendroff [36]. Second, for a tridiagonal matrix, the finite element scheme (linear elements) is identical to the spline Galerkin method (linear splines). For first derivatives, the effective wavenumber is also the same as that for compact finite difference, which is the fourth-order Padé scheme. For the second derivative, however, they are different. Finally, the high-order (large bandwidth) finite element effective wavenumbers depart suddenly from the exact result, effectively limiting the range of wavenumbers for which  $\tilde{k}$  is a good approximation of  $k$ .

The errors in the eigenvalues  $|\lambda - (ik)^n|$  for the first and second derivatives, and errors in representing the complex exponential are plotted in Figs. 4, 5, and 6, respectively. In comparing the different methods, the most obvious difference is the rate of convergence at



**FIG. 4.** Error in the eigenvalue of first derivative operators for matrix bandwidth (a) 3, (b) 7, (c) 11: —, B-spline; - - -, compact finite difference; - · - ·,  $C_0$  finite element Galerkin. - · - ·,  $C_{(d-1)/2}$  finite element Galerkin. For bandwidth equals 3, all schemes yield the same result. In (d), nonmaximum order high-resolution schemes are shown: —, bandwidth = 3, fourth(maximum)-order; - - -, bandwidth = 3, second-order; - · - ·, bandwidth = 5, eighth(maximum)-order; - · - ·, bandwidth = 5, sixth-order.

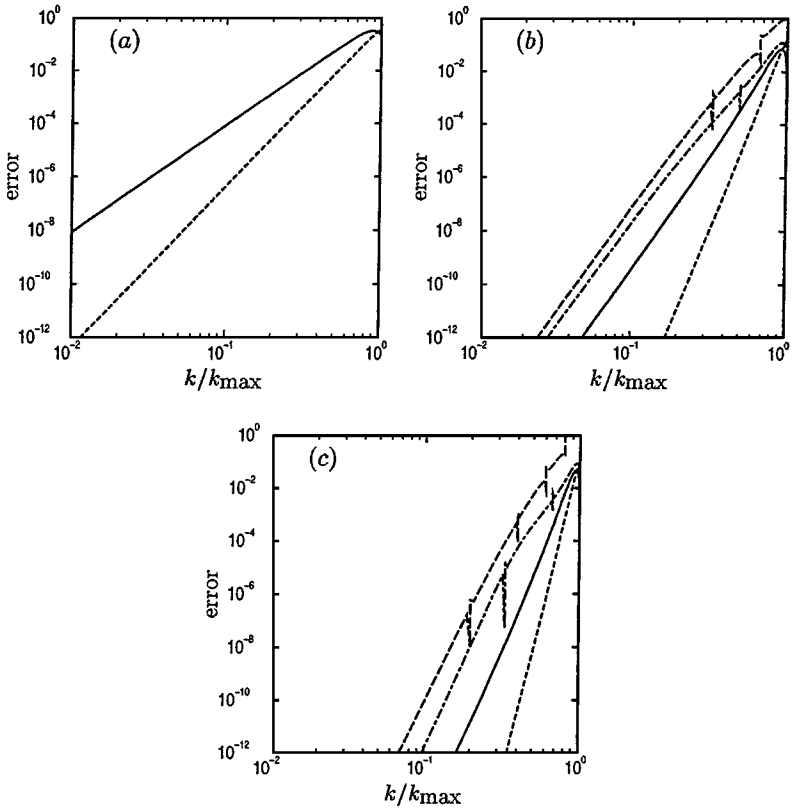
small  $k$ : these curves asymptotically approach zero according to their theoretical convergence rate as shown in Table III.

Note that the compact finite difference convergence rate is significantly faster for large  $w$ . This is possible because in the finite difference method, the “mass” matrix can be different

**TABLE III**  
**Order of Convergence of the Errors of Eigenvalues**  
**and Representation of the Complex Exponential**

Numerical scheme	Eigenvalue of the first derivative operator	Eigenvalue of the second derivative operator	Complex exponential
Finite element Galerkin	$k^{w+2}$	$k^{w+1}$	$k^{\frac{w+1}{2}}$
B-spline Galerkin	$k^{w+2}$	$k^{w+1}$	$k^{\frac{w+1}{2}}$
B-spline collocation	$k^{w+2}$	$k^{w+1}$	$k^{w+1}$
Compact finite difference	$k^{2w-1}$	$k^{2w}$	NA

Note.  $w$  is the matrix bandwidth.

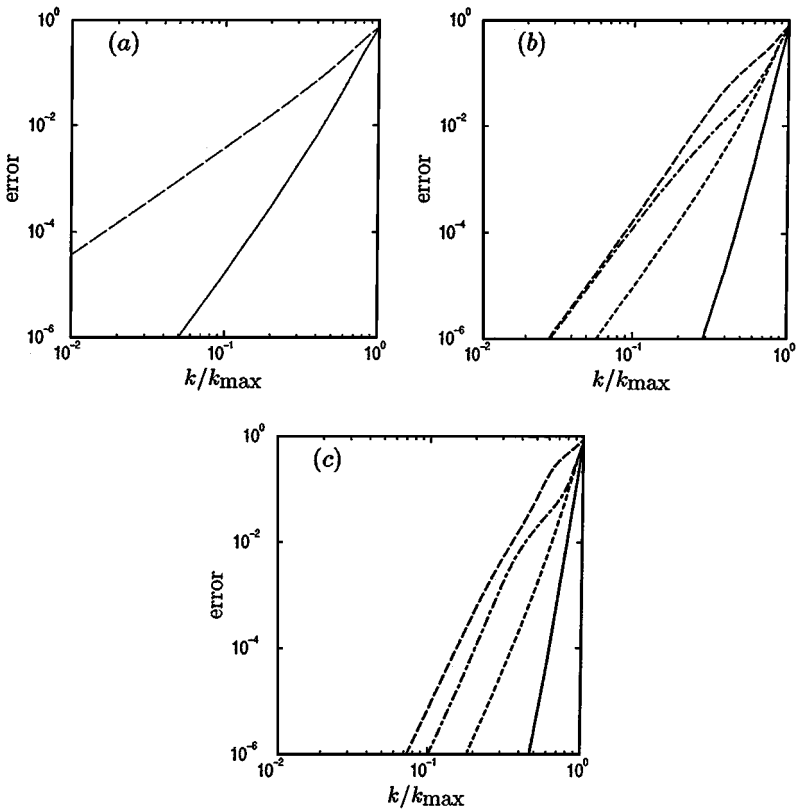


**FIG. 5.** Error in the eigenvalue of second derivative operators for matrix bandwidth (a) 3, (b) 7, (c) 11: —, B-spline; - - -, compact finite difference; - · -,  $C_0$  finite element Galerkin. · · · ·,  $C_{(d-1)/2}$  finite element Galerkin. For bandwidth equals 3, B-spline and finite element yields the same result.

for each order derivative. In contrast, by the nature of functional expansion methods, the mass matrix is the same for all derivatives that can be determined from the representation. If this restriction were imposed on the compact finite difference methods, the same order of convergence as the spline and finite element methods would be obtained.

Another property of the approximate operators is the behavior of the error at large  $k$ . This is important because it determines the range of spatial scales that can be resolved by the numerical method. There is no universally used measure of this resolution property of numerical methods. One measure proposed by Lele [24] is the lowest wavenumber ( $k/k_{\max}$ ) at which the error crosses some arbitrary threshold (say 0.1), giving the fraction of the maximum wavenumber range that is represented to this accuracy or better. In Table IV, this resolved fraction for 10%, 1%, and 0.1% error in the eigenvalues and eigenfunctions is listed for the numerical schemes discussed.

The results discussed above have included only compact finite difference schemes with maximum possible order for the given bandwidth (or stencil size). However, Lele [24] pointed out that one could attain improved resolution properties by decreasing the order of accuracy for a given bandwidth and using the extra degrees of freedom to improve the resolved fraction. The error in first derivative effective wavenumber for two such schemes (bandwidth 3 and 5) are shown in Fig. 4d, along with that of the corresponding maximum



**FIG. 6.**  $L_2$  error in the representation of the complex exponential with wavenumber  $k$  for matrix bandwidth (a) 3, (b) 7, (c) 11: —, B-spline collocation; ---, B-spline Galerkin; - · -,  $C_0$  finite element Galerkin and · · · ·,  $C_{(d-1)/2}$  finite element Galerkin. For bandwidth equals 3, B-spline Galerkin and finite element Galerkin yield the same result.

order scheme. The coefficients for these schemes are also listed in Table I. Each of the high-resolution schemes has an order of accuracy two lower than the maximum possible. This frees one degree of freedom in the scheme which was used to increase the 1% resolved fraction as much as possible. For the tridiagonal scheme, the 1% resolved fraction is increased from 0.35 to 0.52 and in the pentadiagonal case from 0.61 to 0.77. However, the magnitude of the resolved fraction improvement, as well as its importance necessarily decreases with increasing bandwidth and order of the approximation. Also, these schemes tuned to improve 1% resolved fraction degrade the 0.1% resolved fraction, and improve the 10% resolved fraction only marginally. Thus, in using such methods, one needs to be confident that the error level one is targeting is in fact critical, since the high resolution property of the schemes will not be manifested for other error levels. Finally, the resolution improvements discussed here are only for the periodic domain case. In bounded domains, boundary schemes that preserve these properties would need to be developed. Because of these complications, we will consider only the maximal order compact finite difference schemes for each bandwidth, and consider their resolution properties to be representative of what is possible with such schemes. One should keep in mind that for any particular purpose, it may be possible to use the flexibility of compact finite difference methods to attain somewhat better resolution properties by reducing the order or accuracy.

**TABLE IV**  
**Resolved Fraction for Eigenvalues of First and Second Derivative Operators,**  
**and Eigenfunction Representation**

(a)  $d/dx$

Band-width	B-spline collocation			Compact finite difference			$C_0$ finite element Galerkin			$C_{(d-1)/2}$ finite element Galerkin		
	10%	1%	0.1%	10%	1%	0.1%	10%	1%	0.1%	10%	1%	0.1%
3	0.59	0.35	0.20	0.59	0.35	0.20	0.59	0.35	0.20	0.59	0.35	0.20
7	0.80	0.65	0.52	0.84	0.73	0.62	0.82	0.27	0.24	0.80	0.67	0.59
11	0.87	0.77	0.67	0.90	0.83	0.76	0.54	0.50	0.45	0.86	0.78	0.57

(b)  $d^2/dx^2$

Band-width	B-spline collocation			Compact finite difference			$C_0$ finite element Galerkin			$C_{(d-1)/2}$ finite element Galerkin		
	10%	1%	0.1%	10%	1%	0.1%	10%	1%	0.1%	10%	1%	0.1%
3	0.34	0.11	0.03	0.68	0.39	0.22	0.34	0.11	0.03	0.34	0.11	0.03
7	1.00	0.65	0.48	0.94	0.78	0.66	0.57	0.33	0.23	0.81	0.50	0.32
11	1.00	0.80	0.68	0.99	0.88	0.80	0.59	0.45	0.34	1.00	0.66	0.46

(c) Eigenfunction

Band-width	B-spline collocation			B-spline Galerkin			$C_0$ finite element Galerkin			$C_{(d-1)/2}$ finite element Galerkin		
	10%	1%	0.1%	10%	1%	0.1%	10%	1%	0.1%	10%	1%	0.1%
3	0.68	0.43	0.26	0.46	0.16	0.05	0.46	0.16	0.05	0.46	0.16	0.05
7	0.84	0.70	0.57	0.72	0.47	0.29	0.48	0.25	0.15	0.69	0.35	0.17
11	0.89	0.79	0.70	0.81	0.63	0.48	0.53	0.35	0.21	0.77	0.42	0.27

Despite the fact that the order of convergence for the finite element and spline effective wave numbers is the same, the errors in the spline methods are lower at any given  $k$ . In essence, the spline methods have better resolution. This is indicated by Lele's resolution measure, as shown in Table IV. The reason for the lower resolution of the finite elements is the low continuity at the element boundary. One way to understand this (for the first derivative) is to imagine a high-order finite element function  $u$  evolving according to the scalar wave equation:  $\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0$ . At the initial time there are discontinuities in first derivative at the element boundaries. The exact solution would have these discontinuities propagate into the middle of the element, where they cannot be well represented, leading to relatively large errors. This scenario suggests that maximum possible continuity at the knots, that is, splines, is desirable.

The uniform grid periodic analysis is informative, but it does not address two key issues commonly encountered in numerical simulations, that is, nonuniform grids and boundaries. The behavior of finite difference methods in particular is at issue since the formulation discussed in Section 2.2 does not apply directly in these cases. Also, the result that the eigenfunctions of the derivative operators are recovered exactly (Section 3.1) will not hold. It is thus of interest to consider model problems in finite domains. Two such problems are discussed in Sections 4 and 5, namely the first-order wave equation and heat equation. They will only be applied to the B-spline collocation and compact finite difference schemes, the



two best methods discussed above. A preliminary analysis of the other schemes indicates that their performance relative to B-spline collocation and compact finite difference on the wave and heat equations is consistent with that shown above.

### 3.3. Isotropy Properties

The resolution results of the previous sections can be easily extended to multiple dimensions when a tensor product representation is used. The primary complication is that the representation introduces an anisotropy because of the introduction of the grid directions [24, 38]. This anisotropy can be characterized by considering the isotropy of the approximate gradient operator. Consider the gradient of the two-dimensional complex exponential  $\phi = e^{i\mathbf{k}\cdot\mathbf{x}}$ , where  $\mathbf{k}$  is the wave vector, and  $\mathbf{x}$  is the coordinate vector. The exact gradient is given by  $i\mathbf{k}\phi$ , whereas the approximate gradient is  $i\tilde{\mathbf{k}}\phi$ , where  $\tilde{\mathbf{k}}$  is the effective wavenumber vector with components  $\tilde{k}_x = \tilde{k}(k_x) = \tilde{k}(k \cos(\theta))$  and  $\tilde{k}_y = \tilde{k}(k_y) = \tilde{k}(k \sin(\theta))$ , where  $k$  is the magnitude of the wave vector and  $\theta$  is the angle it makes with the  $x$ -axis. The function  $\tilde{k}$  is the one-dimensional effective wavenumber function, as described in Section 3.1.

The relative integrated square error  $e^2$  in the approximate gradient is given by

$$\frac{e^2(\mathbf{k})}{k^2} = \frac{|\mathbf{k} - \tilde{\mathbf{k}}|^2}{k^2} = \left( \cos \theta - \frac{\tilde{k}(k \cos \theta)}{k} \right)^2 + \left( \sin \theta - \frac{\tilde{k}(k \sin \theta)}{k} \right)^2, \quad (10)$$

so there is clearly a variation of this error with the angle  $\theta$ . Assuming that  $\tilde{k}$  is continuously differentiable, this error is minimum when  $\theta = \frac{\pi}{4} + \frac{n\pi}{2}$  and maximum when  $\theta = \frac{n\pi}{2}$ . That is, the error is maximum when the wave-vector is aligned with the grid. If the error in  $\tilde{k}$  (i.e.,  $k - \tilde{k}$ ) does not increase monotonically with  $k$ , then there can be other extrema as well. As is evident in Section 3.2, over most of the wavenumber range for both compact finite difference and B-spline schemes,  $k - \tilde{k}$  can be modeled as

$$k - \tilde{k} = \alpha k^n, \quad (11)$$

where  $n$  is the order of convergence and  $\alpha$  is just a proportionality constant. For this  $k - \tilde{k}$  dependence, the maximum and minimum error can easily be determined:

$$e_{\max}^2(k) = \frac{(k - \tilde{k}(k))^2}{k^2} \quad (12)$$

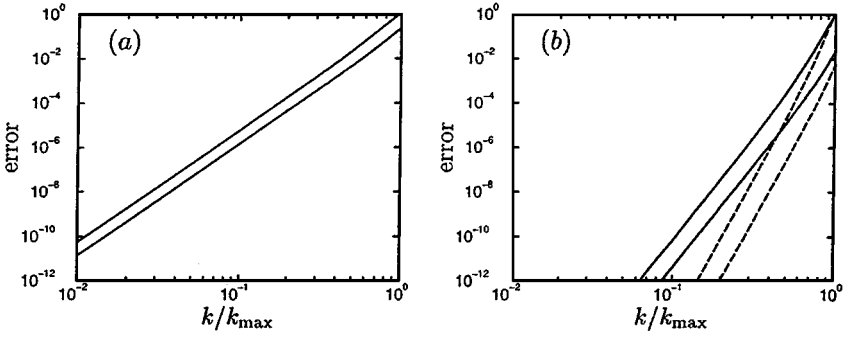
$$e_{\min}^2(k) = \frac{(k - \tilde{k}(k))^2}{k^2} 2^{1-n}. \quad (13)$$

This error variation is just due to the fact that when  $\theta = \pi/4$ , the  $\tilde{k}$  function is evaluated at  $k/\sqrt{2}$  (twice) rather than at  $k$ , so that the error is smaller by a factor of  $2^n$ .

Another quantity of interest is the component of the approximate gradient in the direction of the exact gradient. The ratio of this to the exact gradient is given by

$$c_p(\mathbf{k}) = \frac{\cos \theta \tilde{k}(k \cos \theta) + \sin \theta \tilde{k}(k \sin \theta)}{k}. \quad (14)$$

When solving the two-dimensional wave equation, this is the speed of propagation of the



**FIG. 7.** Maximum and minimum (in  $\theta$ ) error  $(1 - c_p)$  in the phase speed for two-dimensional waves, with matrix bandwidth (a) 3, (b) 7: —, B-spline collocation; - - -, compact finite difference. For bandwidth equals 3, B-spline collocation and compact finite difference yield the same result.

numerical solution relative to the exact speed [24]. This quantity also depends on the angle  $\theta$  and its difference from one is minimum for  $\theta = \frac{\pi}{4} + \frac{n\pi}{2}$ . Again, the anisotropy can be characterized by the maximum and minimum of  $c_p(\text{in } \theta)$ . Using the above model for  $k - \tilde{k}$  we get

$$(1 - c_p(\mathbf{k}))_{\max} = \frac{k - \tilde{k}(k)}{k} \tag{15}$$

$$(1 - c_p(\mathbf{k}))_{\min} = \frac{k - \tilde{k}(k)}{k} 2^{(1-n)/2}. \tag{16}$$

Note that for this simple  $k - \tilde{k}$ ,  $e^2 = (1 - c_p)^2$ ; this is not true in general. As an example,  $(1 - c_p(\mathbf{k}))_{\max}$  and  $(1 - c_p(\mathbf{k}))_{\min}$  are shown in Fig. 7 based on the actual  $\tilde{k}$  for several schemes. Note that maximum and minimum curves are separated a constant ratio of approximately  $2^n$ , consistent with the above analysis.

The above analysis makes it clear that the anisotropy of the approximate gradient operator is governed directly by the errors in  $\tilde{k}$ . Thus, when using tensor product representations of high resolution schemes in which  $k - \tilde{k}$  is small over a wide range of wavenumbers, the anisotropy errors will also be small over the same range of wavenumbers.

#### 4. FIRST-ORDER WAVE EQUATION IN BOUNDED DOMAINS

In this section, B-spline collocation methods and compact finite difference methods are used to solve the first-order wave equation in nonperiodic domains. The problem is defined as

$$\begin{aligned} u_t + u_x &= 0 \quad \text{for } 0 < x < 1, \\ u(0, t) &= \exp(-ikt). \end{aligned} \tag{17}$$

The exact solution assuming periodicity in time is

$$u(x, t) = \exp(ik(x - t)). \tag{18}$$

For numerical solution, it is assumed that  $u(x, t)$  takes the form  $u(x, t) = v(x) \exp(-ikt)$  and solves the following equations for  $v(x)$ :

$$ikv = \frac{dv}{dx}, \quad v(0) = 1. \quad (19)$$

The equation is discretized with B-spline collocation and compact finite difference schemes on both uniform and nonuniform grids.

#### 4.1. B-Spline Collocation Formulation

As mentioned in Section 2.1.3, collocation points at the B-spline maxima are selected. In general, using this “B-spline maxima” collocation formulation with splines of order  $d$ , matrices with  $d + 1$  nonzero diagonals will be obtained. In the case of uniform grids away from the boundary, there are only  $d$  nonzero diagonals as the maxima of splines coincide with the knot points. After discretization, a matrix equation  $i\omega M\alpha = D_1\alpha$  is obtained, where  $M$  and  $D_1$  are the mass and first derivative operator matrix, and  $\alpha$  is the B-spline coefficient vector.

The boundary condition is implemented by replacing the operator at the boundary collocation point with  $v_0 = 1$ .

#### 4.2. Compact Finite Difference Formulation

Lele presents a comprehensive study of high resolution finite difference schemes [24]. In his paper, the effective wavenumber in a periodic domain is investigated. For domains with nonperiodic boundaries, the same analysis is used to obtain the effective wavenumbers both for the interior and the special boundary schemes. The effective wavenumbers for the boundary schemes are in general complex, with the real part associated with the dispersive error and the imaginary part associated with the dissipative error. The conservative formulation, eigenvalue analysis, and stability limits for explicit schemes are also presented. For the details, the reader is directed to [24].

In this section, two issues are addressed. The first is an alternative approach to studying the boundary formulation, instead of the effective wavenumber analysis of Lele. The second is the formulation of schemes with nonuniform grids. The same problem is then solved which offers direct comparison with the B-spline collocation method.

##### 4.2.1. Boundary Formulation

To discretize the hyperbolic equation, the numerical schemes need to resolve the traveling waves in the domain. The boundary formulation is studied using normal modal analysis. Normal modal analysis is also used by Carpenter *et al.* [5] to investigate the stability of boundary treatments for compact finite difference schemes. The similarities of these two analyses will be pointed out after the description of the current boundary formulation.

In the interior, the compact finite difference approximation of the first derivative is derived from Eq. (8), which can be rewritten more generally as

$$v'_i + \sum_{j=1}^m \alpha_j (v'_{i+j} + v'_{i-j}) = \sum_{j=1}^m a_j \frac{v_{i+j} - v_{i-j}}{2j \Delta x}, \quad (20)$$

where  $m$  is related to the matrix bandwidth  $w$  by  $m = \frac{w-1}{2}$ . Knowing that  $v' = ikv$ ,

$$\sum_{j=1}^m c_j v_{i-j} + ikv_i - \sum_{j=1}^m c_j^* v_{i+j} = 0 \quad (21)$$

is obtained, where  $c_j = \frac{a_j}{2j} + i\alpha_j k \Delta x$  and  $c_j^*$  is the conjugate of  $c_j$ . This can be interpreted as a linear recursion relation for  $v_i$ . Such a recursion has solutions  $\Lambda^j$ , where  $\Lambda$  is a function of  $k$ . Substituting  $v_j = \Lambda^j$  into Eq. (21), the characteristic polynomial is obtained,

$$\sum_{j=1}^m c_j \Lambda^{-j} + ik - \sum_{j=1}^m c_j^* \Lambda^j = 0 \quad (22)$$

which has  $2m$  roots. In general, if  $\Lambda_+$  is a root, then  $\Lambda_- = \Lambda_+^{*-1}$  is also a root. These root pairs are denoted as type I root pairs. If  $|\Lambda| = 1$ ,  $\Lambda = \Lambda^{*-1}$ . In this case, there can be two independent roots. These roots are denoted as type II roots. In the limit  $k \rightarrow 0$ , Eq. (22) has the form

$$\sum_{i=1}^m \frac{a_j}{2j} (\Lambda^{-j} - \Lambda^j) = 0, \quad (23)$$

which always has the solutions

$$\Lambda - \Lambda^{-1} = 0 \Rightarrow \Lambda = \pm 1. \quad (24)$$

Changing notation to that of effective wavenumbers,

$$\Lambda = \exp(i\tilde{k}\Delta x) \Rightarrow u(x, t) = \exp\left(i\tilde{k}\left(x - \frac{k}{\tilde{k}}t\right)\right), \quad (25)$$

type I root pairs correspond to conjugate pairs of complex effective wavenumbers  $\tilde{k}$  and  $\tilde{k}^*$ , while type II roots yield real  $\tilde{k}$ . Conjugate pairs of complex effective wavenumbers represent a pair of solutions, one of which grows exponentially in amplitude to the right, the other to the left. Also, for  $k = 0$ , the two solutions yield  $\tilde{k} = 0$  and  $\tilde{k} = k_{\max}$ . Clearly, of the  $2m$  solutions, only one solution with real  $\tilde{k}$  can be a valid approximation to the exact solution. The remainder are spurious. When Eq. (20) is used to solve Eq. (19), the coefficients of the various solutions are determined by the boundary conditions and special differencing schemes used near the boundaries. Clearly, the boundary schemes should be chosen to make the amplitudes of spurious solutions as small as possible.

To see how this works, consider the tridiagonal and pentadiagonal interior scheme (see Table I). For these two cases, the coefficients in the characteristic polynomials and their corresponding roots are given in Table V.  $\Lambda_2$  and  $\Lambda_3$  are complex conjugate pairs while  $\Lambda_0$  and  $\Lambda_1$  have magnitude 1.  $\Lambda_0$  represents the approximation to the exact solution  $\exp(ik\Delta x)$  to the order associated with the scheme and it has a positive group velocity.  $\Lambda_1$  is a spurious wave with a negative group velocity.  $\Lambda_2$  and  $\Lambda_3$  are spurious waves growing exponentially in magnitude to the right and left, respectively. For the spurious waves that grow exponentially to the right, the magnitude of the waves is largest at the right boundary. Thus, by arranging the right boundary schemes to make the  $\Lambda_2$  wave (for example) small at the right boundary,

**TABLE V**  
**Coefficients and Roots of Characteristic Polynomials**

(a) Coefficients of Characteristic Polynomials				
Band-width		$c_1$		$c_2$
3		$\frac{3}{4} + \frac{1}{4}ik\Delta x$		N.A.
5		$\frac{20}{27} + \frac{4}{9}ik\Delta x$		$\frac{25}{216} + \frac{1}{36}ik\Delta x$
(b) Roots of Characteristic Polynomials				
Band-width	$\Lambda_0$	$\Lambda_1$	$\Lambda_2$	$\Lambda_3$
3	$\exp(ik\Delta x) + O((k\Delta x)^5)$	$-1.0000 + 0.3333 ik\Delta x + 0.0555(k\Delta x)^2 + \dots$	N.A.	N.A.
5	$\exp(ik\Delta x) + O((k\Delta x)^9)$	$-1.0000 + 0.1636 ik\Delta x + 0.0134(k\Delta x)^2 + \dots$	$-6.2397 + 1.1118 ik\Delta x - 0.0733(k\Delta x)^2 + \dots$	$-0.1603 + 0.0286 ik\Delta x + 0.0070(k\Delta x)^2 + \dots$

*Note.* The various coefficients in the expressions for  $\Lambda$  are given to four digit accuracy.

the  $\Lambda_2$  solution is small everywhere, regardless of the length of the domain. Similarly, waves growing to the left (e.g.,  $\Lambda_3$ ), should be controlled at the left boundary. For waves with  $|\Lambda| = 1$ , or equivalently real  $\tilde{k}$ , the “group velocity”  $v_g = d\tilde{k}/dk$  determines which boundary should control the wave. With positive group velocity, the left boundary controls the wave because when solving the transient problem (17), information from the boundary will propagate into the domain from the left. Similarly, waves with negative group velocity are controlled at the right boundary. Thus, the spurious solution  $\Lambda_1$  will be controlled by the right boundary scheme, while the physical boundary conditions at the left boundary will control the physical solution  $\Lambda_0$ .

To determine the appropriate inflow boundary schemes, consider the general solution, which near the inflow boundary can be written as

$$v_j = p_0\Lambda_0^j + p_3\Lambda_3^j + O((k\Delta x)^n), \tag{26}$$

where  $n$  is the order of the error in the interior scheme (5 or 9 for tridiagonal and pentadiagonal schemes, respectively). Note that for the tridiagonal scheme,  $\Lambda_2$  and  $\Lambda_3$  can be considered to be zero. The  $O((k\Delta x)^n)$  term is the contribution of the  $\Lambda_1$  and  $\Lambda_2$  waves, which will be this small by construction of the right boundary schemes. Using this expression, the left boundary schemes are constructed to make  $p_0 = 1 + O((k\Delta x)^n)$  and  $p_3 = O((k\Delta x)^n)$  (for the pentadiagonal scheme). This is accomplished using schemes of the form

$$\sum_{j=0}^{m+i} \alpha_{ij} v'_j = \frac{1}{\Delta x} \sum_{j=0}^{3m-i} a_{ij} v_j \quad \text{for } 0 \leq i \leq m-1, \tag{27}$$

for the first  $m = \frac{w-1}{2}$  points, except for the boundary point ( $i = 0$ ), which is replaced by the boundary condition  $v_0 = 1$ . The coefficients for bandwidth 3 and 5 ( $m = 1$  and 2) are shown in Table VI. A Taylor series analysis of these schemes shows them to be of the same

order as the interior scheme, and indeed this is how they were derived. This appears to be a sufficient condition for the suppression of the spurious waves to the desired order. Note, however, that the theory of Gustafsson [5, 15] implies that boundary schemes one order lower than the interior should be adequate to ensure global convergence consistent with the order of the interior scheme.

Near the outflow boundary, the general solution can be written similarly to (26),

$$v_{N-l} = p'_0 \Lambda_0^{-l} + p'_1 \Lambda_1^{-l} + p'_2 \Lambda_2^{-l} + O((k \Delta x)^n), \tag{28}$$

where  $p'_i = p_i \Lambda_i^N$  and  $N$  is the grid number of the right boundary. Boundary schemes that are “mirror images” of the left boundary scheme result in  $p'_1$  and  $p'_2 = O((k \Delta x)^n)$  (pentadiagonal scheme). Thus, we have

$$\sum_{j=0}^{m+N-i} \alpha_{N-i-j} v'_{N-j} = \frac{1}{\Delta x} \sum_{j=0}^{3m-N+i} -a_{N-i-j} v_{N-j} \quad \text{for } N \geq i \geq N - m + 1, \tag{29}$$

where again the coefficients are given in Table VI.

The boundary scheme analysis presented here is similar to the GKS stability analysis of boundary treatments in Carpenter *et al.* [5], in which a similar model problem is used and in which the same spurious waves are treated. However, in Carpenter *et al.*, the assumed temporal form of the solution is more general in that the frequency  $k$  (in our notation, see Eq. (18)) is allowed to be complex. The concern is then whether the time-periodic solutions of Eq. (17) of the form used here are stable. For the fourth-order tridiagonal scheme, Carpenter *et al.* show the combined interior and boundary schemes to be GKS stable, but they do not treat the eighth-order pentadiagonal scheme discussed here. The stability of the solutions to (19) will be discussed in Section 4.4.

#### 4.2.2. Nonuniform Grids

Another issue that needs to be addressed is the formulation of the compact finite difference scheme with nonuniform grids. The approach is to apply a mapping which uses the

**TABLE VI**  
**Coefficients of the Boundary Formulation for the First Derivative**

Band-width	$i$	$\alpha_{i0}$	$\alpha_{i1}$	$\alpha_{i2}$	$\alpha_{i3}$	$\alpha_{i4}$	$\alpha_{i5}$	$\alpha_{i6}$
3	0	1	3	0	0			
5	0	1	12	15	0			
5	1	$\frac{1}{15}$	1	2	$\frac{2}{3}$			
Band-width	$i$	$a_{i0}$	$a_{i1}$	$a_{i2}$	$a_{i3}$	$a_{i4}$	$a_{i5}$	$a_{i6}$
3	0	$-\frac{17}{6}$	$\frac{3}{2}$	$\frac{3}{2}$	$-\frac{1}{6}$	0	0	0
5	0	$-\frac{79}{20}$	$-\frac{77}{5}$	$\frac{55}{4}$	$\frac{20}{3}$	$-\frac{5}{4}$	$\frac{1}{5}$	$-\frac{1}{60}$
5	1	$-\frac{247}{900}$	$-\frac{19}{12}$	$\frac{1}{3}$	$\frac{13}{9}$	$\frac{1}{12}$	$-\frac{1}{300}$	0

uniform mesh scheme in the mapped coordinate. The mesh mapping is given in Eq. (7). The discretization equations (8), (27), and (29) are modified by the mesh mapping

$$\left. \frac{dv}{dx} \right|_i = \left. \frac{dv}{d\xi} \right|_i \left. \frac{d\xi}{dx} \right|_i. \quad (30)$$

Thus, the equation to be solved is  $\frac{dv}{d\xi} \frac{d\xi}{dx} = ikv$ . The same interior and boundary scheme are then used in  $\xi$ .

### 4.3. Comparison

Tests based on the solution of the first-order wave equation were carried out with  $N = 100$ . Before discussing the results, however, it should be noted that different from the effective wavenumbers, the accuracy of the solution of the wave equation is dependent on the number of intervals  $N$  apart from the wavenumber. In this sense, the results here are less general than those of the effective wavenumber. Nevertheless, using the same  $N$  for both schemes allows us to compare their order of convergence and resolution.

The results on uniform grids are discussed first. The  $L_2$  errors in the representation of the solution of the wave problem using B-spline collocation and compact finite difference methods are shown in Fig. 8. For B-spline collocation methods, the errors vary with  $k$  like  $k^{w+2}$ . Notice that in periodic domains, the convergence rates of the eigenvalue of the first derivative operator and the eigenfunction are  $k^{w+2}$  and  $k^{w+1}$ , respectively (see Figs. 4 and 6 and Table III). For compact finite difference schemes, the  $L_2$  error varies with  $k$  like  $k^{2w-1}$ . This is consistent with the theoretical convergence rate (Fig. 4 and Table III), though the curve is not smooth. Apparently, the boundary condition and boundary schemes do not affect the convergence rate of either scheme.

An error of B-spline solution that varies as  $k^{w+2}$  is curious because the  $L_2$  error is bounded from below by the error in representing the complex exponential, which as indicated in Table III, converges like  $k^{w+1}$ . To resolve this apparent inconsistency, note that the equations for the B-spline coefficients have the same form as the compact finite difference equations, and so the analysis in Section 4.2.1 applies to them as well. The solution is characterized

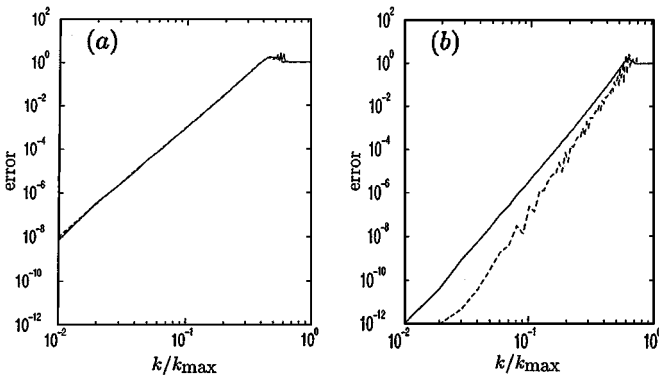
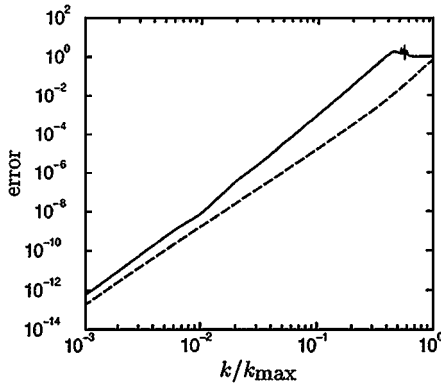


FIG. 8.  $L_2$  error in the representation of the solution of the wave equation with wavenumber  $k$  on uniform grids for matrix bandwidth (a) 3, (b) 5: —, B-spline collocation; ---, compact finite difference.



**FIG. 9.**  $L_2$  representation error with wavenumber  $k$  for matrix bandwidth 3: —, solution of the wave equation (17) in bounded domain; - - -, complex exponential in periodic domain.

by roots of the characteristic polynomial ( $\Lambda_j$ ), only one of which is an approximation to the “exact” solution. In this case the “exact” solution (away from the boundaries) is the dependence of the B-spline coefficients for representation of the complex exponential in an infinite domain, i.e.,  $\Lambda_0 = \exp ik\Delta x$ . With the exact B-spline coefficients, the  $L_2$  error converges like  $k^{w+1}$ . However, in the solution to (17), the analysis in Section 4.2.1 shows that there are errors in the B-spline coefficients of order  $(k\Delta x)^{w+2}$ , due to the error in  $\Lambda_0$  and the errors in  $p_j$ . For large  $k$ , the errors in the B-spline coefficients dominate, resulting in a  $k^{w+2}$  dependence, and for sufficiently small  $k$  the representation error of the complex exponential dominates, resulting in a  $k^{w+1}$  dependence. This is illustrated in Fig. 9, in which the error in the solution to (19) is shown along with the representation error of the complex exponential for the  $w = 3$  case. The change from  $k^5$  to  $k^4$  dependence occurs where  $kL \sim 1$  (or  $k/k_{\max} \sim 1/N\pi$ , where  $N = 100$ ), which is consistent with the fact that the dominant error arising from the error in  $\Lambda_0^N$  is of order  $kL(k\Delta x)^{w+1}$  (see Section 4.2.1). Note that for B-spline Galerkin solutions of this problem, the error in representing the complex exponentials, which goes like  $k^{\frac{w+1}{2}}$  (see Table III), dominates over the errors in  $\Lambda_0$  and  $p_j$  for all  $k$ , resulting in an overall convergence of  $k^{\frac{w+1}{2}}$ . For the compact finite difference methods, the error decreases like  $k^{2w-1}$ , consistent with the analysis of Section 4.2.1 and the convergence rates listed in Table III.

When plotted versus  $k/k_{\max} = k\Delta x/\pi$ , the error curves must depend on  $N$ , since the error behaves as  $kL(k\Delta x)^{n-1} \sim N(k\Delta x)^n$ . One way to obtain a curve that is valid for all  $N$  is to plot  $\text{error}/N$  versus  $k/k_{\max}$ . The resulting curve would not shift as  $N$  changes except when the error is close to 1, and, for B-spline collocation, when  $k/k_{\max} \leq 1/N\pi$  (where the representation error dominates).

Perhaps more important than the order of convergence is the resolution of the two schemes. The well-resolved fraction of the wavenumber range for the solution of the wave equation is shown in Table VII. It can be seen that for tridiagonal schemes, the two have almost the same resolution. For pentadiagonal schemes, compact finite difference has better resolution due to the higher order of convergence.

Another issue is the effect of a nonuniform grid. The  $L_2$  errors in the representation of the solution of the wave problem using the two numerical schemes on nonuniform grids are shown in Fig. 10. Here, the wavenumber is normalized by the maximum wavenumber



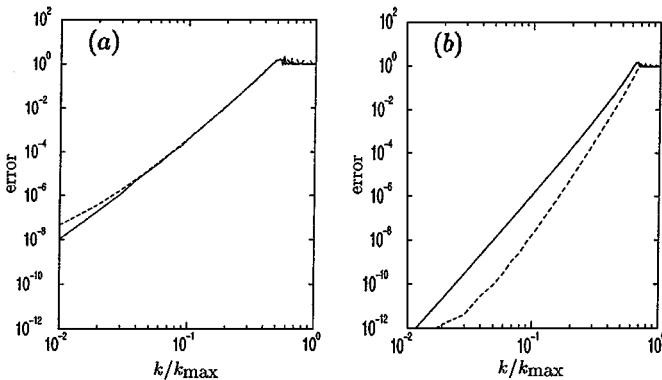
**TABLE VII**  
**Resolved Fraction for the Solution of the Wave Equation**  
**for Uniform Grid Distribution**

Band-width	B-spline collocation			Compact finite difference		
	10%	1%	0.1%	10%	1%	0.1%
3	0.25	0.15	0.10	0.24	0.15	0.09
5	0.40	0.30	0.22	0.45	0.34	0.27

$k_{\max} = \frac{\pi}{\Delta x_{\max}}$ , where  $\Delta x_{\max}$  is the maximum grid spacing. Basically, both schemes maintain the same convergence rate, as in the case of uniform grids. Note that for the compact finite difference schemes, the curves turn up at the lowest wavenumber and the cause is not clear. With regard to the resolved fraction, Table VIII indicates that the two tridiagonal schemes again have the same resolution. (Note however that in nonuniform grids, B-spline collocation has elements outside the three “main” diagonals in the interior.) For pentadiagonal schemes, compact finite difference has better resolution.

The order of convergence of the two schemes suggests that the difference in resolution properties between compact finite difference and B-spline collocation will become bigger as the matrix bandwidth increases. Also, comparing results in periodic and bounded domains (Tables IV and VII), it is found that the resolution in finite domains is substantially lower. In particular, the error reaches 1 at  $k/k_{\max}$  ranging from 0.4 to 0.6 in Fig. 8.

This plateau of the error at 1 for moderate values of  $k/k_{\max}$  is also caused by dominance of the error in  $\Lambda_0^N$ , which is the error (of function values or B-spline coefficients) at the right boundary. In terms of  $k/k_{\max}$ , this error goes like  $N(k/k_{\max})^n$ , where  $n$  is  $w + 2$  for B-spline collocation and  $2w - 1$  for compact finite difference. This error is of order 1 when  $k/k_{\max} \approx N^{-1/n}$ , and thus for larger  $k/k_{\max}$  the overall solution error should be of order 1. With  $N = 100$  (as in this case) and with  $n = 5, 7,$  and  $9$  the value of  $N^{-1/n}$  is 0.4, 0.52, and 0.60, respectively, which is in reasonably good agreement with the start of the plateau for  $w = 3, w = 5$  B-splines and  $w = 5$  compact finite difference.



**FIG. 10.**  $L_2$  error in the representation of the solution of the wave equation with wavenumber  $k$  on nonuniform grids for matrix bandwidth (a) 3, (b) 5: —, B-spline collocation; ---, compact finite difference.

**TABLE VIII**  
**Resolved Fraction for the Solution of the Wave Equation**  
**for Nonuniform Grid Distribution**

Band-width	B-spline collocation			Compact finite difference		
	10%	1%	0.1%	10%	1%	0.1%
3	0.29	0.19	0.12	0.29	0.19	0.12
5	0.47	0.35	0.25	0.55	0.43	0.34

**4.4. Stability of the Time-Harmonic Solutions**

The time-harmonic solutions evaluated in Section 4.3 were forced to be time-harmonic, so there is no guarantee that these solutions are stable when the wave equation is solved as an initial value problem. To evaluate this, we need only examine the eigenvalues of the relevant approximate operator. A perturbation  $\delta u$  from one of the harmonic solutions of (17) is governed by

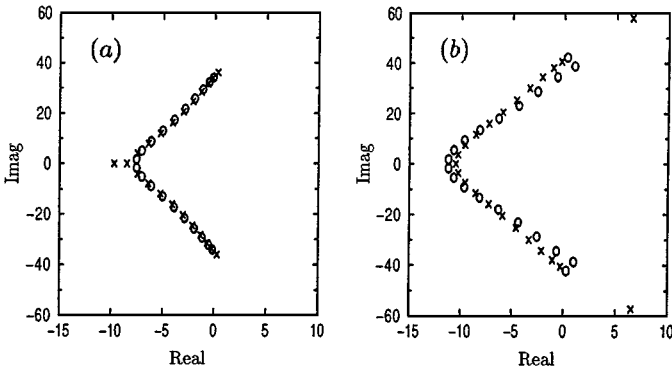
$$\delta u_t + \delta u_x = 0 \quad 0 \leq x \leq 1, \quad \delta u(0, t) = 0. \tag{31}$$

When numerically discretized this equation has the form  $M\delta\alpha_t = D_1\delta\alpha$ , where  $\delta\alpha$  is either the B-spline coefficient vector (B-spline collocation) representing  $\delta u$  or the vector of grid point values of  $\delta u$  (finite difference). Consistent with the implementation for the time-harmonic solutions, boundary conditions are implemented by replacing the equation associated with the point at  $x = 0$  with  $\delta(\alpha_0)_t = 0$ , where  $\delta\alpha_0$  is associated with  $x = 0$ . This has the effect of replacing the first (say) row of  $M$  with all zeros, except for a one in the  $\delta\alpha_0$  column, and replacing the first row of  $D_1$  with all zeros. Call the modified matrices  $\tilde{M}$  and  $\tilde{D}_1$ , respectively, Eq. (31) can be rewritten  $\delta\alpha_t = \tilde{M}^{-1}\tilde{D}_1\delta\alpha$ . All solutions of this equation will decay to zero provided all the eigenvalues of  $\tilde{M}^{-1}\tilde{D}_1$  have negative real parts.

These eigenvalues have been computed for the bandwidth 3 and 5 schemes examined in this section with  $N = 20$ . For tridiagonal compact finite difference schemes, the eigenvalues do indeed have negative real parts.<sup>1</sup> However, the B-spline collocation methods each produce two eigenvalues with positive real parts. As an example, the eigenspectra of the bandwidth 3 and 5 B-spline collocation operator are shown in Fig. 11, along with the corresponding finite difference eigenvalue spectra. The eigenfunction associated with the unstable eigenvalue oscillates with wavelength  $2\Delta x$  and decays rapidly away from the in-flow boundary. Thus, when solving (17) as an initial value problem with the B-spline collocation methods, this growing eigenfunction will be observed, rather than the time-harmonic solutions.

In the applications we have in mind, such as solution of the Navier–Stokes equations, the equations are not strictly hyperbolic. There is a viscous damping term, which if the viscosity is large enough would stabilize this numerical instability. For this to occur, the viscous damping rate, which is approximately  $\nu a/\Delta x^2$  ( $a$  depends on the details of

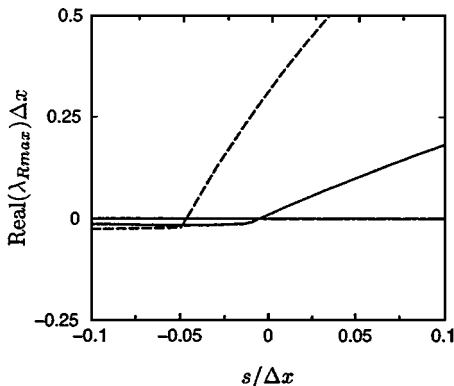
<sup>1</sup> In [5], the eigenvalues for the bandwidth 3 schemes were computed and found to include eigenvalues with positive real parts. However, the boundary conditions are implemented differently. Instead of modifying  $M$  and  $D_1$  as described above, the product  $M^{-1}D_1$  was modified by zeroing the first row. (Private communication with M. H. Carpenter.) These are not equivalent.



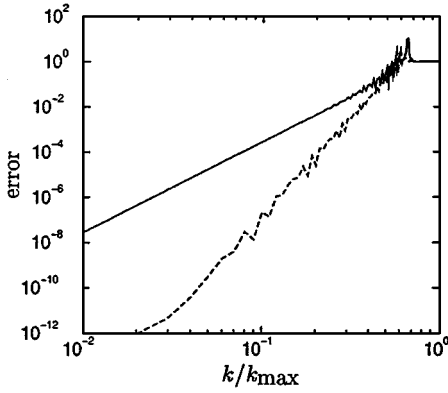
**FIG. 11.** Numerically determined eigenvalue spectrum for B-spline collocation scheme ( $\times$ ) and compact finite scheme ( $\circ$ ) with matrix bandwidth of (a) 3, (b) 5 and  $N = 20$ . For B-spline method, collocation points are chosen at the maxima of the splines. Also, for B-spline with bandwidth 5, there is an eigenvalue  $-42.48$  not shown in the figure.

the second derivative approximation, and for high resolution schemes is a substantial fraction of  $\pi^2$ ; see Section 3.2), must be larger (in magnitude) than the unstable eigenvalue. If  $\lambda_{R\max}\Delta x = 0.3$  as it is for the bandwidth 5 scheme, then the stability requirement would be that the cell Reynolds number is less than  $3.3a$  ( $\Delta x/\nu < 3.3a$  in this case). This is not an arduous cell Reynolds number requirement.

To avoid these stability problems all together, it would be preferable if the schemes could be modified to yield all eigenvalues with a negative real part. In the B-spline collocation scheme, the only aspect that can be modified is the location of the collocation points. Away from the boundaries, the collocation points are fixed by the need to preserve the spatial symmetry of the operators. However, the exceptional collocation points near the boundary (those that are not coincident with a knot; see Section 2.1.3) can be adjusted. As it happens, the maximum real part of the eigenvalues is sensitive to the placement of these collocation points, as is shown in Fig. 12. For the two B-spline schemes shown here, moving these exceptional points toward the boundary a small amount (5% of the knot spacing) stabilizes the unstable eigenvalues. The change in collocation point location has no impact on the



**FIG. 12.** Variation of the maximum real part of the eigenvalues,  $\lambda_{R\max}$ , as near-boundary exceptional collocation points are shifted from the B-spline maximum location by an amount  $s$ : —, bandwidth = 3; - - -, bandwidth = 5. Negative  $s$  is toward the boundary.



**FIG. 13.**  $L_2$  error in the representation of the solution of the wave equation with wavenumber  $k$  on nonuniform grids for matrix bandwidth 5: —, eighth-order interior, sixth-order near-boundary, third-order boundary; - - -, consistent eighth-order.

resolution properties discussed earlier. For bandwidths 7 and 9, it was found that a shift of 6.5% of knot spacing stabilizes the unstable eigenvalues. At this time it is not clear why moving the exceptional collocation points improves the stability of the B-spline schemes at an inflow boundary. It is also not clear if larger shifts will be required to stabilize the higher order (bandwidth) schemes.

If one were numerically simulating Eq. (17), the stability of the time discretization would be an issue. In an explicit scheme, there would be a time step restriction that is fixed by the largest eigenvalues of the homogeneous numerical operators, such as those shown in Fig. 11. The value of  $\lambda_{I \max} \Delta x$  (the maximum imaginary part of the eigenvalues) is insensitive to  $\Delta x$ , and is in good agreement with  $\tilde{k}_{\max} \Delta x$ , where  $\tilde{k}_{\max}$  is the maximum attained value of the effective wavenumber (see Fig. 2).

In compact finite difference methods, boundary schemes of the same order as the interior can be unstable [5, 25]. This stability problem is solved by using boundary schemes of lower order. For high-resolution methods such as the eighth-order pentadiagonal scheme, stable boundary schemes of compatible order can be very difficult to find. Carpenter developed stable sixth-order boundary schemes using a rather involved stability analysis [5]. Mahesh [25] chose sixth- and third-order for near-boundary and boundary points respectively to stabilize the overall scheme. However, this cure comes at a high cost in resolution and order of accuracy. Mahesh [25] pointed out that lower order boundary schemes reduce the formal order of the overall scheme to one greater than that of the boundary [15]. A test was carried out using eighth-order pentadiagonal interior, sixth-order near-boundary (second and  $(n - 1)$ th row), and third-order boundary ( $n$ th row) schemes. This scheme has a much lower resolution than the eighth-order interior and boundary schemes (Fig. 13).

### 5. HEAT EQUATION IN BOUNDED DOMAINS

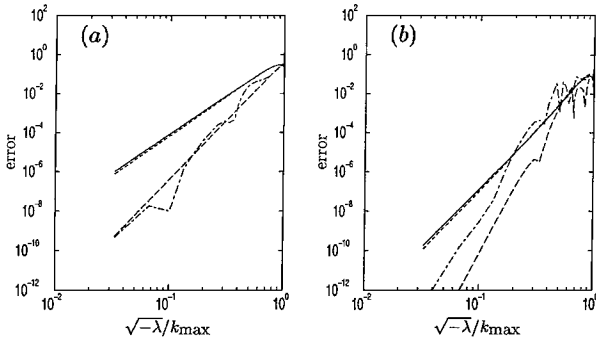
In this section, the eigenvalue problem arising from the heat equation is solved using B-spline collocation and compact finite difference methods. The problem is defined as

$$v'' = \lambda v \quad \text{for } 0 < x < 1, \tag{32}$$

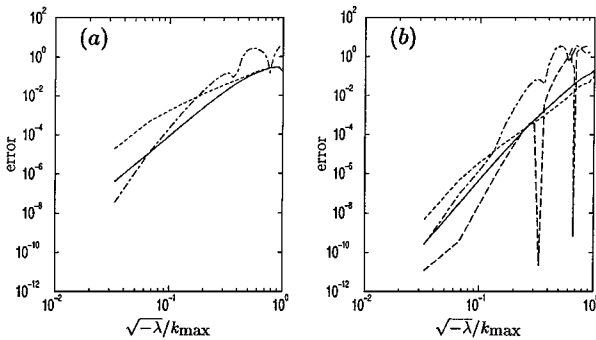
with some boundary conditions, the most common ones being the Dirichlet ( $v(0) = v(1) = 0$ )

**TABLE IX**  
**Coefficients of the Boundary Formulation for the Second Derivative**

Band-width	$i$	$\alpha_{i0}$	$\alpha_{i1}$	$\alpha_{i2}$	$\alpha_{i3}$	
3	0	1	11	0	0	
5	0	1	$\frac{18922}{563}$	$\frac{65943}{563}$	0	
5	1	$\frac{23}{688}$	1	$\frac{2335}{688}$	$\frac{2659}{3096}$	
Band-width	$i$	$a_{i0}$	$a_{i1}$	$a_{i2}$	$a_{i3}$	$a_{i4}$
3	0	13	-27	15	-1	0
5	0	$\frac{2186893}{101340}$	$\frac{526369}{5067}$	$-\frac{3296517}{11260}$	$\frac{1940803}{10134}$	$-\frac{583529}{20268}$
5	1	$\frac{753829}{1114560}$	$\frac{57209}{20640}$	$-\frac{58367}{8256}$	$\frac{172793}{55728}$	$\frac{4453}{8256}$
Band-width	$i$	$a_{i5}$	$a_{i6}$	$a_{i7}$		
3	0	0	0	0	0	0
5	0	0	$\frac{14802}{2815}$	$\frac{14839}{20268}$		$\frac{2659}{50670}$
5	1	1	$-\frac{391}{20640}$	$\frac{529}{1114560}$		0



**FIG. 14.** Error of the eigenvalues of the heat equation with wavenumber  $k$  on uniform grids for matrix bandwidth (a) 3 (b) 5: —, Dirichlet boundary conditions; - - -, Neumann boundary condition, both for B-spline; - · -, Dirichlet boundary condition; · · · ·, Neumann boundary condition, both for compact finite difference.



**FIG. 15.** Error of the eigenfunctions of the heat equation with wavenumber  $k$  on uniform grids for matrix bandwidth (a) 3, (b) 5: —, Dirichlet boundary condition; - - -, Neumann boundary condition, both for B-spline; - · -, Dirichlet boundary condition; · · · ·, Neumann boundary condition, both for compact finite difference. For the compact finite difference, with bandwidth of 3 and Dirichlet boundary conditions, the eigenfunctions are exact to round-off errors.

**TABLE X**  
**Resolved Fraction for the Eigenvalues for Uniform Grid Distribution**

Bandwidth	Boundary condition	B-spline collocation			Compact finite difference		
		10%	1%	0.1%	10%	1%	0.1%
3	Dirichlet	0.33	0.10	0.03	0.66	0.36	0.20
3	Neumann	0.36	0.10	0.03	0.46	0.36	0.16
5	Dirichlet	0.80	0.46	0.26	>0.76	0.50	0.40
5	Neumann	0.76	0.46	0.26	0.43	0.36	0.23

and Neumann ( $v'(0) = v'(1) = 0$ ) boundary conditions. In both cases, the eigenvalues are

$$\lambda_k = -(\pi k)^2, \tag{33}$$

where  $k$  is an integer. The corresponding eigenfunctions are  $\tilde{v}_k(x) = \sin(\pi kx)$  and  $\tilde{v}_k(x) = \cos(\pi kx)$  for Dirichlet and Neumann boundary conditions respectively.

**5.1. B-Spline Collocation Formulation**

Discretizing with B-spline collocation method, we obtain the matrix equation  $\lambda M\alpha = D_2\alpha$ , where  $M$  is the mass and  $D_2$  the second derivative operator matrix, and  $\alpha$  the B-spline coefficient vector for the eigenfunctions. For Dirichlet boundary conditions,  $v_0 = v_N = 0$ . For Neumann boundary conditions,  $v'_0$  and  $v'_N$  are set to zero.

**5.2. Compact Finite Difference Formulation**

Similar to the case of first derivative, the discretized derivative operators are derived from Eq. (9) in the interior. Near the boundary, the symmetry breaks down and the corresponding equation becomes

$$\sum_{j=0}^{m+i} \alpha_{ij} v''_j = \frac{1}{(\Delta x)^2} \sum_{j=0}^{3m+1-i} a_{ij} v_j \quad \text{for } 0 \leq i \leq m-1, \tag{34}$$

$$\sum_{j=0}^{m+N-i} \alpha_{N-i-j} v''_{N-j} = \frac{1}{(\Delta x)^2} \sum_{j=0}^{3m+1-N+i} a_{N-i-j} v_{N-j} \quad \text{for } N \geq i \geq N-m+1,$$

where  $m = \frac{w-1}{2}$ . The coefficients in Eq. (34) are determined by matching the Taylor series coefficients to one order less than the interior for tridiagonal schemes and to the same order of the interior for pentadiagonal schemes (using boundary of the same order for tridiagonal schemes gives rise to poor resolution for unknown reasons). The coefficients are shown in Table IX for the two schemes. After discretization, a generalized eigenvalue problem  $\lambda M_2\alpha = D_2\alpha$  is obtained, where  $M_2$  is the mass,  $D_2$  the second derivative operator matrix, and  $\alpha$  the eigenfunction. The generalized eigenvalue problem can be solved with the appropriate boundary condition.

*5.2.1. Boundary Conditions*

The Dirichlet boundary conditions are implemented by setting  $v_0 = v_N = 0$ . For Neumann boundary conditions, a one-sided explicit (i.e., not compact) finite difference

scheme is used to set  $v'_0 = v'_N = 0$ . Note that this makes the boundary scheme inconsistent with the interior scheme. Also, a very long one-sided finite difference expression is required to maintain the same order as the interior compact finite difference approximations.

### 5.2.2. Nonuniform Grids

To solve the problem on a nonuniform grid, a mesh mapping is used as in the wave equation. Notice from the chain rule,

$$\frac{d^2v}{dx^2} = \left(\frac{d\xi}{dx}\right)^2 \frac{d^2v}{d\xi^2} + \frac{d^2\xi}{dx^2} \frac{dv}{d\xi}. \quad (35)$$

The derivative in the nonuniform  $x$ -coordinate is expressed in terms of those in the transformed uniform  $\xi$ -coordinate. In the  $\xi$ -coordinate, there are finite difference representations of the derivative operators (expressed as  $M_1^{-1}D_1$  and  $M_2^{-1}D_2$ . Note that  $M_1$  and  $M_2$  are different). The finite difference approximation of the second derivative operator can hence be expressed as in Eq. (35).

## 5.3. Comparison

Tests based on the eigenvalue problem are performed using  $N = 30$ . Results based on different  $N$  suggest that  $N$  has no influence on the order of convergence and minor influence on the well-resolved fraction. The results obtained on uniform grids are presented first. The errors in approximating the eigenvalues are shown in Fig. 14. Regardless of the boundary conditions, B-spline collocation schemes have eigenvalue errors which decrease with wavenumber as  $k^{w+1}$ , while compact finite difference has a convergence rate of  $k^{2w}$ . Both of the above are consistent with their corresponding convergence rates in periodic domains (see Table III). For compact finite difference, however, the boundary conditions do have an effect on the magnitude of the error. Neumann boundary conditions give larger errors in the eigenvalues, perhaps due to the boundary approximation of  $v'$ . Also, with the compact finite difference, there are some sharp decreases in error at particular wavenumbers for reasons that are not clear. At high wavenumbers, wiggles appear on the compact finite difference curves irrespective of the boundary conditions. With regard to resolution, we refer to Table X, which gives the resolved fraction for the eigenvalues. In many cases, compact finite difference schemes provide better resolution for the eigenvalues. For pentadiagonal scheme, B-splines have better resolved fractions in many cases. However, due to the high convergence order of the compact finite difference, the more stringent the tolerance for resolved fractions, the better the compact finite difference does.

The  $L_2$  errors of the eigenfunctions of the heat equation are shown in Fig. 15. For the B-spline collocation schemes, the convergence rate for both Dirichlet and Neumann boundary conditions appears to be  $k^{w+1}$ , but in the Neumann case this asymptotic rate is not attained until  $k < 0.06$ , with the resulting impact on resolution. For compact finite difference schemes, the eigenfunctions have errors that converge at a rate approximately equal to  $k^{2w}$ . However, with Neumann boundary conditions, the errors are again larger. In fact, using the Dirichlet boundary condition, the tridiagonal schemes give a solution that is exact to round-off errors. It is also very interesting to note that the pentadiagonal scheme curve shows two sharp decreases at  $\sqrt{-\lambda}/k_{\max} = \frac{1}{3}$  and  $\sqrt{-\lambda}/k_{\max} = \frac{2}{3}$ . At these two particular wavenumbers, the symmetries of the approximate eigenfunctions make the point representations exact. Table XI shows the resolved fraction of eigenfunctions. As for

**TABLE XI**  
**Resolved Fraction for the Eigenfunctions for Uniform Grid Distribution**

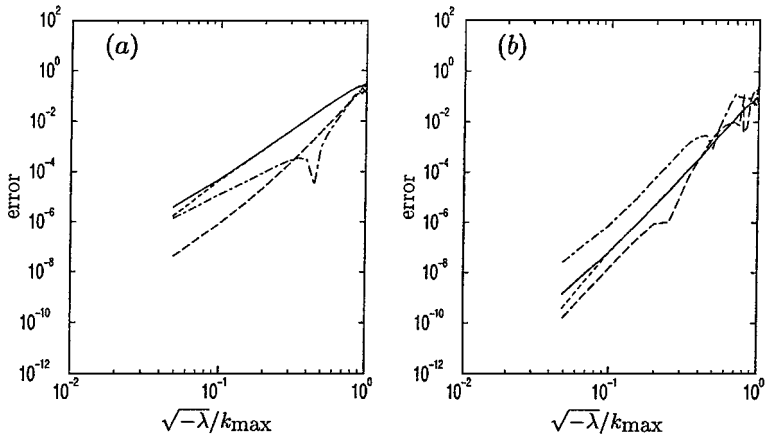
Bandwidth	Boundary condition	B-spline collocation			Compact finite difference		
		10%	1%	0.1%	10%	1%	0.1%
3	Dirichlet	0.43	0.23	0.13	1.00	1.00	1.00
3	Neumann	0.40	0.13	0.06	0.23	0.16	0.10
5	Dirichlet	0.73	0.46	0.30	0.43	0.37	0.33
5	Neumann	0.93	0.56	0.33	0.36	0.20	0.16

**TABLE XII**  
**Resolved Fraction for the Eigenvalues for Nonuniform Grid Distribution**

Bandwidth	Boundary condition	B-spline collocation			Compact finite difference		
		10%	1%	0.1%	10%	1%	0.1%
3	Dirichlet	0.44	0.09	<0.05	0.73	0.44	0.19
3	Neumann	0.44	0.14	0.04	0.74	0.54	0.44
5	Dirichlet	0.92	0.53	0.29	0.73	0.49	0.34
5	Neumann	0.97	0.53	0.29	0.54	0.29	0.14

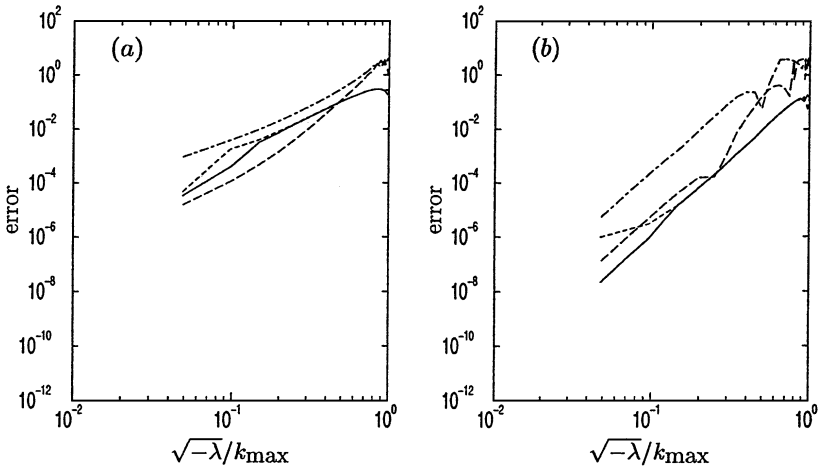
**TABLE XIII**  
**Resolved Fraction for the Eigenfunctions for Nonuniform Grid Distribution**

Bandwidth	Boundary condition	B-spline collocation			Compact finite difference		
		10%	1%	0.1%	10%	1%	0.1%
3	Dirichlet	0.39	0.19	0.09	0.44	0.24	0.14
3	Neumann	0.39	0.14	0.04	0.29	0.09	<0.04
5	Dirichlet	0.72	0.43	0.29	0.39	0.29	0.24
5	Neumann	0.72	0.43	0.29	0.25	0.15	0.09



**FIG. 16.** Error of the eigenvalues of the heat equation with wavenumber  $k$  on nonuniform grids for matrix bandwidth (a) 3, (b) 5: —, Dirichlet boundary condition; - - -, Neumann boundary condition, both for B-spline; ····, Dirichlet boundary condition; - · - ·, Neumann boundary condition, both for compact finite difference.





**FIG. 17.** Error of the eigenfunctions of the heat equation with wavenumber  $k$  on nonuniform grids for matrix bandwidth (a) 3, (b) 5: —, Dirichlet boundary condition; ---, Neumann boundary condition, both for B-spline; ---, Dirichlet boundary condition; ····, Neumann boundary condition, both for compact finite difference.

the eigenvalues, compact finite difference schemes in general provide better resolution for tridiagonal methods while B-splines do better for pentadiagonal schemes.

On nonuniform grids, the behavior of both B-spline collocation and compact finite difference schemes is shown in Fig. 16 and 17 and Tables XII and XIII. The errors in the eigenvalues of the heat equation are shown in Fig. 16. B-spline collocation methods maintain the same convergence rate of  $k^{w+1}$  as in the case of uniform grids irrespective of the boundary conditions. Compact finite difference schemes, however, show a degradation. Convergence rates of the eigenvalues is 2 to 3 orders less than the corresponding rate of  $k^{2w}$  on uniform grids, with Neumann boundary conditions giving worse convergence rates. Regarding resolution, compact finite difference provides better resolution for bandwidth  $w = 3$ , while B-spline collocation schemes provides better resolution for bandwidth  $w = 5$ .

The  $L_2$  errors of the eigenfunctions of the heat equation are shown in Fig. 17. B-spline collocation schemes give convergence rates of  $k^w$  approximately, with Dirichlet boundary conditions giving slightly better solutions at low  $k$ . The degradation of resolution is not serious when nonuniform grids are used instead of uniform ones. Compact finite difference schemes, however, show quite serious degradation of convergence and resolution on nonuniform grids. They have convergence rates of about  $k^w$ , compared to  $k^{2w}$  on uniform grids. A very interesting result is that B-spline and compact finite difference schemes appear to have the same convergence rates on nonuniform grids. From table XIII, it can also be seen that B-spline collocation methods have better resolution properties on nonuniform grids.

## 6. DISCUSSION AND CONCLUSIONS

The results of this paper indicate that in many situations compact finite difference schemes have better resolution and convergence properties than the other numerical methods tested. The comparisons were done for schemes with the same matrix bandwidths, which we use as a surrogate for computational cost. Furthermore, it was shown that finite element and B-spline Galerkin methods had inferior resolution to compact finite difference and B-spline collocation. There are several aspects of these results that deserve further discussion.

Regarding high-order finite element methods, it was noted that a reason for their lower resolution in these tests was their lower order ( $C_0$  and  $C_{(d-1)/2}$ ) continuity at the element boundaries (i.e., knots), whereas the spline basis retains as high a degree of continuity as possible, given the order of the piecewise polynomial representation. In essence, in spline methods, an increase in the degree of the polynomials is used to increase the degree of continuity, while in  $C_0$  and  $C_{(d-1)/2}$  finite elements, it is used to increase the number of degrees of freedom of the representation. The results of the tests here suggest that the added degrees of freedom do not produce much in the way of added accurately represented modes, resulting in poor resolution properties. However, the improved resolution of splines is not without cost; that is, the representation of the polynomials in each interval (element) is not iso-parametric, a very convenient property of finite element representations. Consequently, it is much easier to formulate multidimensional finite elements on complex and/or unstructured grids, than it is to formulate spline methods.

It was also noted that piecewise polynomial Galerkin methods yielded poorer representations of complex exponentials (the derivative eigenfunction) than collocation methods. This is true for both finite element methods and spline methods. This is a curious result because Galerkin approximations minimize  $L_2$  error for a given representation. The reason for the curious result is that we are comparing methods with the same matrix bandwidth. For example, a Galerkin method that yields pentadiagonal matrices has cubic polynomials, whereas a pentadiagonal collocation methods has quintic polynomials. The result is a fourth-order accurate representation for Galerkin and a sixth-order accurate representation for collocation. However, there are other reasons one might choose a Galerkin method, despite its higher cost; for example, a Galerkin method is trivially shown to be conservative.

The two methods discussed here with the best convergence and resolution properties are compact finite difference and spline collocation, and the comparison between them includes four major issues that must be traded off against the improved order of accuracy and in many cases better resolution of the finite difference methods:

1. The generally superior convergence and resolution of compact finite difference compared to B-spline collocation is simply due to the fact that in the finite difference case, the “mass matrix” is not constrained to be the same for all derivatives. There may, however, be costs in code complexity or computational effort in having different mass matrices, depending on the details of the problem being solved.

2. Another difference is in the treatment near a boundary. In the finite difference case, special difference schemes must be formulated near the boundary, and such boundary treatments can be difficult to formulate. For the wave equation, a criterion for a boundary treatment with good resolution was developed in Section 4.2.1, and schemes that satisfy the criterion were found by imposing a formal order of accuracy consistent with the interior scheme. However, consistent order of accuracy is a necessary but not necessarily sufficient condition for the criterion to be satisfied, and directly constructing schemes to satisfy the criterion is prohibitively cumbersome in all but the simplest cases (e.g., the tridiagonal scheme). Thus, we do not have a practical constructive prescription for boundary schemes that satisfies the criterion in Section 4.2.1. Furthermore, the criterion does not guarantee the stability of the resulting combined interior/boundary schemes. Indeed, the bandwidth 5 scheme constructed this way was unstable. Modifying such schemes to be stable as in Carpenter [5] is an arduous task, which has not been done for schemes of order greater than sixth. The lack of stable high-order boundary schemes prompted some

authors to combine high-order interior representations with much lower order boundary schemes (e.g., Mahesh [25]), which spoils the accuracy and resolution of the combined scheme.

In the heat equation problem, the development of a criterion like that used in the advection equation for good boundary schemes is not as obvious, so we have an even less well defined procedure for the boundary treatment in this case. Finally, recall that for the heat equation, with Neumann conditions, an approximation of the first derivative at the boundary had to be devised, so the derivative boundary condition could be imposed. This was essentially ad hoc, and was not inherently consistent with the remainder of the scheme.

These boundary complications are in principle surmountable in finite difference methods; but, they are completely obviated in B-spline methods, since the B-spline representation (like any other functional representation) unambiguously defines the near-boundary scheme. The only complication is that in B-spline collocation the location of the near-boundary collocation points that are not attached to knots must be specified. An algorithm based on the maxima of the B-spline functions was proposed, but at the inlet their location affects the stability of the scalar advection scheme. By slightly adjusting the location of these collocation points ( $0.065\Delta x$ ) toward the boundary, it is possible to obtain stable representations for bandwidth up to 9 (tenth order).

3. On a nonuniform mesh, the spline method can be used directly, without recourse to a mapping to a domain with a uniform mesh, as we did for the finite difference case. Thus, the method can easily be applied to an arbitrary mesh, for which no analytic mapping is known. Besides, the resulting approximations are simpler, with no explicit metric terms, and in the case of approximating the second derivative, no first derivative term appears. Of course, one can construct finite difference methods on arbitrary meshes as well, either by direct construction or by numerically defined mappings. But the process is cumbersome, and for direct construction generally yields schemes that are lower order than the uniform mesh schemes (for the same matrix bandwidth). With the spline methods, the nonuniform mesh formulation is no different from the uniform mesh.

4. In Sections 4 and 5, the error associated with the spline collocation method was well behaved and consistent with the results of the Fourier analysis in Section 3. The same cannot be said for the finite difference schemes. For them, the error spectra were more erratic, with a variety of unexplained features. Furthermore, in at least one case (i.e., nonuniform mesh heat equation with Neumann conditions), the convergence rates appeared to be the same as its spline counterpart, inconsistent with the simple Fourier analysis in Section 3.

Thus, when using a high-order spline collocation scheme instead of compact finite difference with the same bandwidth, one is trading away a potentially higher convergence rate and somewhat better resolution in many cases for a more straightforward and robust formulation. And, as suggested by the results of Section 4 and 5, in complicated situations, there is no guarantee that the finite difference method would actually yield the theoretical higher convergence rate. Finally, Shariff and Moser [34] showed that a B-spline representation could be used on embedded meshes while preserving the high-order convergence of the schemes. They used a Galerkin formulation, but the same general technique is applicable using collocation. The only uncertainty is the location of the collocation points. The B-spline maxima are an obvious choice. However, the stability of the resulting approximations need to be assessed.

**APPENDIX A: THE B-SPLINE AND FINITE-ELEMENT BASIS**

Consider a domain  $[0, L]$  divided into  $N$  intervals with  $N + 1$  grid points (knots)  $t_0, t_1, \dots, t_N$ .  $N + k$  B-splines of order  $k$  can be generated according to recursion relationship [8],

$$B_j^k(x) = \frac{x - t_{j-k-1}}{t_{j-1} - t_{j-k-1}} B_{j-1}^{k-1}(x) + \frac{t_j - x}{t_j - t_{j-k}} B_j^{k-1}(x), \quad j = 1, 2, \dots, N + k \quad (A.1)$$

where  $B_j^k(x)$  is the  $j$ th B-spline of order  $k$ . The B-spline of order 0 is simply the top hat function

$$B_j^0(x) = \begin{cases} 1 & \text{if } t_{j-1} \leq x \leq t_j \\ 0 & \text{otherwise.} \end{cases} \quad (A.2)$$

Close to the boundaries, evaluation of  $B_j^k(x)$  involves “virtual points”  $t_j$ ’s outside the range of knots ( $j < 0$  or  $j > N$ ). In periodic domains, these virtual points are given as

$$t_j = \begin{cases} t_{N+j} - L & \text{if } j < 0 \\ t_{j-N} + L & \text{if } j > N. \end{cases} \quad (A.3)$$

In bounded domains, the virtual points can be placed arbitrarily (either at the boundary or outside the domain). The choice of virtual points determines the near-boundary basis functions, but it does not affect the spline solution space, and thus does not impact the solution. It is most convenient to locate the virtual points at the boundary, thus increasing the multiplicity of knots there, i.e.,

$$t_j = \begin{cases} 0 & \text{if } j < 0 \\ L & \text{if } j > N. \end{cases} \quad (A.4)$$

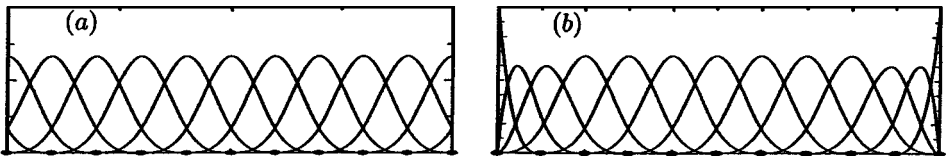
Figure 18 shows B-splines in periodic and bounded domains.

The  $i$ th derivative of the B-spline  $B_j^k$  is written in terms of lower order B-splines as

$$\frac{d^i B_j^k}{dx^i} = \sum_{l=j-1}^j \alpha_{jl}^i B_l^{k-i}, \quad (A.5)$$

where the coefficients  $\alpha$  are found from the recursion [8]

$$\alpha_{jl}^{i+1} = (k - i) \frac{\alpha_{j(l+1)}^i - \alpha_{jl}^i}{t_l - t_{l-k+i}}, \quad (A.6)$$



**FIG. 18.** Cubic B-splines in (a) periodic and, (b) bounded domains. Knots are denoted by ●.

with the starting condition

$$\alpha_{jl}^0 = \delta_{jl}. \quad (\text{A.7})$$

To implement the Galerkin and finite element B-spline methods it is necessary to evaluate definite integrals of the splines and their derivatives, and/or to evaluate the splines and their derivatives at specified points in the domain. The recursions described above allow the B-splines to be evaluated. To compute integrals, Gauss quadrature in each interval between the knots is used, with enough quadrature points for the resulting integrals to be exact. This calculation procedure is stable to roundoff error [8].

The finite element representation in terms of piecewise polynomials is similar to the spline representation, differing only in the degree of continuity at the knots (element boundaries). The same recursion relations can be used to define and evaluate the  $C_0$  and  $C_{(d-1)/2}$  finite element basis by introducing a multiplicity of knots at each knot location [8]. For the  $C_0$  and  $C_{(d-1)/2}$  finite elements each knot point has multiplicity of  $d$  and  $(d+1)/2$ , respectively. Then evaluation of the finite element basis and the Galerkin integrals is accomplished as discussed above. Note that the finite element basis functions defined here are not those most often used in finite element methods, but they result in the same matrix bandwidth and necessarily yield identical results, since the finite element solution space is the same.

## APPENDIX B: DETAILS OF THE FOURIER ANALYSIS

As discussed in Section 2, when the numerical schemes discussed here are applied, the representation of derivative operators is given by

$$M\alpha' = D\alpha, \quad (\text{B.1})$$

where  $M$  is the mass matrix,  $D$  is the derivative matrix (of whatever order),  $\alpha$  is the vector representing the function (either point values or coefficients), and  $\alpha'$  is the vector representing the derivative.

For the B-splines or the compact finite difference methods in a periodic domain with a uniform grid, these matrices are both banded and circulant. The mass matrix is also symmetric, while the derivative matrix is symmetric for even derivatives and anti-symmetric for odd derivatives. The  $i$ th row of these equations is given by

$$\sum_{j=-n}^n m_j \alpha'_{i+j} = \sum_{j=-n}^n d_j \alpha_{i+j}, \quad (\text{B.2})$$

where  $2n+1$  is the bandwidth of the matrices,  $m_j$  are the mass matrix elements (with  $m_0$  on the main diagonal), and  $d_j$  are the derivative matrix elements. The  $m_j$  and  $d_j$  satisfy

$$m_j = m_{-j} \quad d_j = \begin{cases} d_{-j} & \text{even derivatives} \\ -d_{-j} & \text{odd derivatives.} \end{cases} \quad (\text{B.3})$$

The values of the matrix elements are determined as described in Appendix A for the B-splines and are given in Tables I and II for compact finite difference schemes.

Since the matrices are circulant, the elements of the  $j$ th eigenvector are all given by  $e^{ijl2\pi/N}$ , where  $N$  is the size of the matrix and we can identify  $k = 2\pi j/N$  as the wavenumber of the mode. The eigenvalues  $\lambda(k)$  are also easily determined. For odd derivatives,

$$\lambda(k) = \frac{i \sum_{j=1}^n 2d_j \sin(kj)}{m_0 + \sum_{j=1}^n 2m_j \cos(kj)}, \tag{B.4}$$

and for even derivatives

$$\lambda(k) = \frac{d_0 + \sum_{j=1}^n 2d_j \cos(kj)}{m_0 + \sum_{j=1}^n 2m_j \cos(kj)}. \tag{B.5}$$

This is how the effective wavenumbers presented in Section 3 were computed for B-splines and compact finite differences.

As explained in Section 3.1, the accuracy of the eigenfunctions in the B-spline representation is simply the accuracy of representing the complex exponential. This is measured as the  $L_2$  error  $\epsilon$  in the representation, measured per length of the periodic domain

$$\epsilon^2 = \frac{1}{L_x} \int_0^{L_x} \left| e^{ikx} - \sum_j \alpha_j B_j(x) \right|^2 dx, \tag{B.6}$$

where  $\alpha_j$  and  $B_j$  are the B-spline coefficients and functions, and the wavenumber  $k$  is given by  $k = \tilde{k}2\pi/L_x$ , with  $\tilde{k}$  an integer between  $-N/2 + 1$  and  $N/2$ , where  $N$  is the number of intervals in the domain. Expanding the integrand and integrating the  $e^{ikx}e^{-ikx}$  term, one obtains

$$\epsilon^2 = 1 - \frac{2}{L_x} \mathcal{R} \left( \sum_j \alpha_j \int_0^{L_x} e^{-ikx} B_j(x) dx \right) + \frac{1}{L_x} \sum_j \sum_l \alpha_j^* \alpha_l \int_0^{L_x} B_j(x) B_l(x) dx. \tag{B.7}$$

First, note that since the matrices describing the scheme are circulant, the B-spline coefficients are given by  $\alpha_j = a e^{ikj\Delta x}$ , where  $\Delta x = L_x/N$ . Then the second term can be simplified as

$$\sum_j \alpha_j \int_0^{L_x} e^{-ikx} B_j(x) dx = \sum_j a e^{ikj\Delta x} e^{-ikj\Delta x} \beta_k = a N \beta_k, \tag{B.8}$$

where

$$\beta_k = \int_0^{L_x} e^{-ikx} B_0(x) dx. \tag{B.9}$$

The third term in (B.7) is simplified by noting that the integral is the Galerkin mass matrix, and that  $e^{ikj\Delta x}$  is an eigenvector. Then

$$\sum_j \sum_l \alpha_j^* \alpha_l \int_0^{L_x} B_j(x) B_l(x) dx = |a|^2 \sum_j e^{ikj\Delta x} e^{-ikj\Delta x} \lambda_k = |a|^2 N \lambda_k, \tag{B.10}$$

where  $\lambda_k$  is the eigenvalue of the Galerkin mass matrix that is associated with the eigenvector  $e^{ikj\Delta x}$ . The eigenvalue is given by

$$\lambda_k = m_{g0} + \sum_{l=1}^{n_g} 2m_{gl} \cos lk\Delta x, \quad (\text{B.11})$$

where  $m_{gl}$  and  $n_g$  are the elements and half-bandwidth of the Galerkin mass matrix for the order splines considered. Finally, the error is written

$$\epsilon^2 = 1 - \frac{2}{\Delta x} \mathcal{R}(a\beta_k) + \frac{1}{\Delta x} |a|^2 \lambda_k. \quad (\text{B.12})$$

The error is minimized when  $a = \beta_k^*/\lambda_k$ , and in that case the error is given by

$$\epsilon = \sqrt{1 - \frac{|\beta_k|^2}{\Delta x \lambda_k}}. \quad (\text{B.13})$$

The value of  $\beta_k$  can in principle be determined analytically, given the piecewise polynomial description of the splines. However, it is more convenient to evaluate  $\beta_k$  using Gauss quadrature in each interval. The number of quadrature points is selected to give results accurate to machine precision. This is how the B-spline errors given in Section 3 were determined.

The analysis of high-order finite elements is somewhat more complicated, since the basis functions are not all simple shifts of each other. For  $C_j$  elements of order  $o$ , ( $j < o$ ), the number of degrees of freedom per element is  $d = o - j$ , and therefore, there are  $d$  different types of basis functions. The entire basis is formed of shifts of these  $d$  basic types. In a Galerkin method, the derivatives are represented by matrices with several diagonals of  $d \times d$  blocks, in which the blocks in each diagonal are identical. The matrix can be thought of as block circulant or block toeplitz. For finite element representations with continuity up to  $C_{(o-1)/2}$ , the matrices are block tridiagonal, and the  $i$ th block row of the derivative representation  $M\alpha' = D\alpha$  is written

$$\sum_{j=-1}^1 \tilde{M}_j \tilde{\alpha}'_{i+j} = \sum_{j=-1}^1 \tilde{D}_j \tilde{\alpha}_{i+j}, \quad (\text{B.14})$$

where  $\tilde{M}_j$  and  $\tilde{D}_j$  are the  $d \times d$  blocks on the  $j$  diagonal of the mass and derivative matrix, respectively, and  $\tilde{\alpha}_i$  is a vector of length  $d$  representing the coefficients of the  $d$  basis functions associated with element  $i$ . Due to symmetry,  $\tilde{M}_{-j} = \tilde{M}_j^T$  and  $\tilde{D}_{-j} = \tilde{D}_j^T$  for even derivatives while  $\tilde{D}_{-j} = -\tilde{D}_j^T$  for odd derivatives.

The eigenvectors of such a system are of the form  $\tilde{\phi} e^{ikj\Delta x}$ , where  $\Delta x$  is the element size and  $\tilde{\phi}$  is an eigenvector of the  $d \times d$  system

$$\hat{D}\tilde{\phi} = \tilde{\lambda}\hat{M}\tilde{\phi}, \quad (\text{B.15})$$

where

$$\hat{D} = \tilde{D}_0 + D_1 e^{ik\Delta x} + D_{-1} e^{-ik\Delta x} \quad (\text{B.16})$$

$$\hat{M} = \tilde{M}_0 + M_1 e^{ik\Delta x} + M_{-1} e^{-ik\Delta x}. \quad (\text{B.17})$$

We expect  $d$  solutions for each  $-N/2 < k \leq N/2$ ; each of which will be an approximation to the the exact derivative eigenvalue and eigenfunction associated with wavenumber  $\tilde{k}_\ell = k + 2\pi\ell/\Delta x$  for some  $-d/2 > \ell > d/2$ . The approximate eigenvalues and eigenfunctions can thus be determined by solving the matrix eigenvalue problem given in (B.15), but it is then necessary to determine which approximate eigenvalue/eigenvector pair  $(\tilde{\lambda}, \tilde{\phi})$  is associated with which  $\tilde{k}_\ell$ . This is accomplished by determining which eigenvalue best approximates each of the complex exponentials. As was done above, the  $L_2$  error is measured

$$\epsilon^2 = 1 - \frac{2}{\Delta x} \mathcal{R} \left( a \sum_{l=1}^d \phi_l \int_0^{L_x} e^{-i\tilde{k}_\ell x} B_0^l(x) dx \right) + \frac{|a|^2}{\Delta x} \sum_{l=1}^d \sum_{j=1}^d \hat{M}_{lj} \tilde{\phi}_l \tilde{\phi}_j^*, \quad (\text{B.18})$$

where  $B_0^l(x)$  is the basis function of type  $l$  at element 0. The scale factor  $a$  can be determined as before to minimize the error, and the resulting minimum error has the same form as Eq. (B.13) with

$$\beta_k = \sum_{l=1}^d \phi_l \int_0^{L_x} e^{-i\tilde{k}_\ell x} B_0^l(x) dx \quad (\text{B.19})$$

$$\lambda_k = \sum_{l=1}^d \sum_{j=1}^d \hat{M}_{lj} \tilde{\phi}_l \tilde{\phi}_j^*. \quad (\text{B.20})$$

By computing this error for each  $\tilde{k}_\ell - \hat{\phi}$  pair, the best fit to each represented complex exponential is determined, and in this way, each eigenvalue is associated with a wavenumber it represents. Comparing them yields the error in the eigenvalue, and (B.13) is the corresponding error in the eigenfunction.

The error shown in Fig. 6 is not the error in the eigenfunction, but the error in the Galerkin representation of the complex exponential, which, for finite elements, is somewhat less than the eigenfunction error. The coefficients of the finite element approximation to the complex exponential  $e^{ikx}$  are of the form  $\hat{\phi} e^{ikj\Delta x}$ , where  $\hat{\phi}$  is the solution vector to the system

$$\hat{M} \hat{\phi} = R, \quad (\text{B.21})$$

and the right-hand side vector  $R$  is

$$R_l = \int_0^{L_x} B_0^l e^{ikx} dx. \quad (\text{B.22})$$

The error in this approximation is given by

$$\epsilon^2 = 1 - \frac{1}{\Delta x} \hat{\phi} \times R^*, \quad (\text{B.23})$$

where  $\hat{\phi} \times R$  is real and positive because  $\hat{M}$  is conjugate symmetric and positive definite. This is the error that was plotted in Fig. 6.

## ACKNOWLEDGMENT



## REFERENCES

1. A. J. Baker, *Finite Element Computational Fluid Mechanics* (McGraw-Hill, New York, 1983).
2. G. K. Batchelor, *The Theory of Homogeneous Turbulence* (Cambridge Univ. Press, Cambridge, UK, 1982).
3. K. J. Bathe and V. Sonnad, in *Proceedings of AIAA Computers and Fluid Dynamics Conference Palo Alto, California, 1981*.
4. C. Canuto, M. Y. Hussaini, A. Quarteroni, and T. A. Zang, *Special Methods in Fluid Dynamics* (Springer-Verlag, Berlin, 1987).
5. M. H. Carpenter, D. Gottlieb, and S. Abarbanel, The stability of numerical boundary treatments for compact high-order finite-difference schemes, *J. Comput. Phys.* **108**, 272 (1993).
6. A. M. Davies, Applications of the Galerkin method to the solution of burger's equation, *Comput. Methods Appl. Mech. Eng.* **14**, 305 (1978).
7. A. M. Davies, The use of the Galerkin method with a basis of B-splines for the solution of the one-dimensional primitive equations, *J. Comput. Phys.* **27**, 123 (1978).
8. C. de Boor, *A Practical Guide to Splines* (Springer-Verlag, New York/Berlin, 1978).
9. H. Deconinck and C. Hirsch, in *Proceedings of 7th International Conference on Numerical Methods in Fluid Dynamics*, edited by W. C. Reynolds and R. W. MacCormack (Springer-Verlag, New York/Berlin, 1980) p. 138.
10. P. Devloo, J. T. Oden, and P. Pattani, An  $h - p$  adaptive finite element method for the numerical simulation of compressible flow, *Comput. Meths. Mechs. Eng.* **70**, 203 (1988).
11. C. A. J. Fletcher, *Computational Techniques for Fluid Dynamics* (Springer-Verlag, Berlin/New York, 1991, Vol. 1).
12. J. B. Freund, P. Moin, and S. K. Lele, Compressibility effects in a turbulent annular mixing layer, Technical Report TF-72 (Mechanical Engineering, Stanford University, 1997).
13. D. J. Fyfe, The use of cubic spline in the solution of two point boundary value problem, *Comput. J.* **12**, 188 (1969).
14. D. Gottlieb and S. A. Orszag, *Numerical Analysis of Spectral Methods* (Soc. for Industr. & Appl. Math., Philadelphia, 1977).
15. B. Gustafsson, The convergence rate for difference approximations to mixed initial boundary value problems, *Math. Comput.* **29**(130), 396 (1975).
16. M. T. Heath, *Scientific Computation: An Introductory Survey* (WCB, McGraw-Hill, New York, 1997).
17. R. S. Hirsh, Higher order accurate differenece solutions of fluid mechanics problems by a compact differencing technique, *J. Comput. Phys.* **19**, 90 (1975).
18. G. E. Karniadakis, Spectral element simulations of laminar and turbulent flows in complex geometries, *Appl. Numer. Math.* **6**, 85 (1989).
19. G. E. Karniadakis and S. A. Orszag, *Some Novel Aspects of Spectral Methods, Algorithmic Trends in Computational Fluid Dynamics* (Springer-Verlag, Berlin/New York, 1993).
20. J. Kim, P. Moin, and R. D. Moser, Turbulence in channel flow at low reynolds number, *J. Fluid Mech.* **177**, 133 (1987).
21. M. H. Kobayashi, On a class of padé finite volume methods, *J. Comput. Phys.* **156**, 137 (1999).
22. A. G. Kravchenko, P. Moin, and R. D. Moser, Zonal embedded grids for numerical simulation of wall-bounded turbulent flows, *J. Comput. Phys.* **127**, 412 (1996).
23. S. A. Orszag and M. Israeli, Numerical simulation of viscous incompressible flows, *Ann. Rev. Fluid Mech.* **6**, 281 (1974).
24. S. K. Lele, Compact finite difference schemes with spectral-like resolution, *J. Comput. Phys.* **103**, 16 (1992).
25. K. Mahesh, A family of high order finite difference schemes with good spectral resolution, *J. Comput. Phys.* **145**, 332 (1998).
26. S. A. Orszag, Numerical simulation of incompressible flows within simple boundaries: accuracy, *J. Fluid Mech.* **49**, 75 (1971).
27. S. A. Orszag, On the resolution requirements of finite-difference schemes, *Stud. Appl. Math.* **50**, 395 (1971).

28. D. Papamoschou and S. K. Lele, Vortex-induced disturbance field in a compressible shear layer, *Phys. Fluids A* **5**(6), 1412 (1993).
29. A. T. Patera, A spectral element method for fluid dynamics: Laminary flow in a channel expansion, *J. Comput. Phys.* **54**, 468 (1984).
30. R. S. Rogallo and P. Moin, Numerical simulation of turbulent flows, *Ann. Rev. Fluid Mech.* **16**, 99 (1984).
31. S. G. Rubin and R. A. Graves, Viscous flow solutions with a cubic spline approximation, *Comput. Fluids* **3**, 1 (1975).
32. S. G. Rubin and P. K. Khosla, Higher-order numerical solutions using cubic splines, *AIAA J.* **14**, 851 (1976).
33. S. G. Rubin and P. K. Khosla, Polynomial interpolation methods for viscous flow calculations, *J. Comput. Phys.* **24**, 217 (1977).
34. K. Shariff and R. D. Moser, Two-dimensional mesh embedding for b-spline methods, *J. Comput. Phys.* **145**, 471 (1998).
35. P. R. Spalart, Direct simulation of a turbulent boundary layer up to  $re_\theta = 1410$ , *J. Fluid Mech.* **187**, 61 (1988).
36. B. Swartz and B. Wendroff, The relation between the gallerkin and collocation methods using smooth splines, *SIAM J. Numer. Anal.* **11**, 994 (1974).
37. F. Thomasset, *Implementation of Finite Element Methods for Navier–Stokes Equations* (Springer-Verlag, Berlin, 1981).
38. R. Vichnevetsky and J. B. Bowles, *Fourier Analysis of Numerical Approximations of Hyperbolic Equations* (Soc. for Industr. & Appl. Math, Philadelphia, 1982).
39. K. N. S. Kasi Viswanadham and S. R. Koneru, Finite element method for one-dimensional and two-dimensional time dependent problems with b-splines, *Comput. Meth. Appl. Mech. Eng.* **108**, 201 (1993).