

1. Introduction

1.1 Historical Background

Spectral methods are a class of spatial discretizations for differential equations. The key components for their formulation are the trial functions (also called the expansion or approximating functions) and the test functions (also known as weight functions). The trial functions, which are linear combinations of suitable trial basis functions, are used to provide the approximate representation of the solution. The test functions are used to ensure that the differential equation and perhaps some boundary conditions are satisfied as closely as possible by the truncated series expansion. This is achieved by minimizing, with respect to a suitable norm, the residual produced by using the truncated expansion instead of the exact solution. The residual accounts for the differential equation and sometimes the boundary conditions, either explicitly or implicitly. For this reason they may be viewed as a special case of the method of weighted residuals (Finlayson and Scriven (1966)). An equivalent requirement is that the residual satisfy a suitable orthogonality condition with respect to each of the test functions. From this perspective, spectral methods may be viewed as a special case of Petrov-Galerkin methods (Zienkiewicz and Cheung (1967), Babuška and Aziz (1972)).

The choice of the trial functions is one of the features that distinguishes the early versions of spectral methods from finite-element and finite-difference methods. The trial basis functions for what can now be called *classical spectral methods* – spectral methods on a single tensor-product domain – are global, infinitely differentiable and nearly orthogonal, i.e. the matrix consisting of their inner products has very small bandwidth; in many cases this matrix is diagonal. (Typically the trial basis functions for classical spectral methods are tensor products of the eigenfunctions of singular Sturm-Liouville problems). In contrast, for the h version of finite-element methods, the domain is divided into small elements, and low-order trial functions are specified in each element. The trial basis functions for finite-element methods are thus local in character and still nearly orthogonal, but not infinitely differentiable. They are thus well suited for handling complex geometries. Finite-difference methods are typically viewed from a pointwise approximation perspective rather than from a trial function/test function perspective. However, when

appropriately translated into a trial function/test function formulation, the finite-difference trial basis functions are likewise local.

The choice of test functions distinguishes between the three earliest types of spectral schemes, namely, the Galerkin, collocation, and tau versions. In the Galerkin (1915) approach, the test functions are the same as the trial functions. They are, therefore, infinitely smooth functions that individually satisfy some or all of the boundary conditions. The differential equation is enforced by requiring that the integral of the residual times each test function be zero, after some integration-by-parts, accounting in the process for any remaining boundary conditions. In the collocation approach the test functions are translated Dirac delta-functions centered at special, so-called collocation points. This approach requires the differential equation to be satisfied exactly at the collocation points. Spectral tau methods are similar to Galerkin methods in the way the differential equation is enforced. However, none of the test functions need satisfy the boundary conditions. Hence, a supplementary set of equations is used to apply the boundary conditions.

The collocation approach appears to have been first used by Slater (1934) and by Kantorovic (1934) in specific applications. Frazer, Jones and Skan (1937) developed it as a general method for solving ordinary differential equations. They used a variety of trial functions and an arbitrary distribution of collocation points. The work of Lanczos (1938) established for the first time that a proper choice of trial functions and distribution of collocation points is crucial to the accuracy of the solution. Perhaps he should be credited with laying down the foundation of the orthogonal collocation method. This method was revived by Clenshaw (1957), Clenshaw and Norton (1963) and Wright (1964). These studies involved the application of Chebyshev polynomial expansions to initial-value problems. Villadsen and Stewart (1967) developed this method for boundary-value problems.

The earliest applications of the spectral collocation method to partial differential equations were made for spatially periodic problems by Kreiss and Olinger (1972) (who called it the Fourier method) and Orszag (1972) (who termed it pseudospectral). This approach is especially attractive because of the ease with which it can be applied to variable-coefficient and even nonlinear problems. The essential details will be furnished below.

The Galerkin approach enjoys the esthetically pleasing feature that the trial functions and the test functions are the same, and the discretization is derived from a weak form of the mathematical problem. Finite-element methods customarily use this approach. Moreover, the first serious application of spectral methods to PDE's – that of Silberman (1954) for meteorological modeling – was a Galerkin method. However, spectral Galerkin methods only became practical for high resolution calculations of such nonlinear problems after Orszag (1969, 1970) and Eliassen, Machenhauer and Rasmussen (1970) developed transform methods for evaluating the convolution sums arising from quadratic nonlinearities. (Nonlinear terms also increase the

cost of finite-element methods, but not nearly as much as they do for spectral Galerkin methods.) For problems containing more complicated nonlinear terms, high-resolution spectral Galerkin methods remain impractical.

The tau approach is a modification of the Galerkin method that is applicable to problems with nonperiodic boundary conditions. It may be viewed as a special case of the so-called Petrov-Galerkin method. Lanczos (1938) developed the spectral tau method, and Orszag's (1971b) application of the Chebyshev tau method to produce highly accurate solutions to fluid dynamics linear stability problems inspired considerable use of this technique, not just for computing eigenvalues but also for solving constant-coefficient problems or subproblems, e.g., for semi-implicit time-stepping algorithms.

In the middle 1980's newer spectral methods, which combined the Galerkin approach with Gaussian quadrature formulas, came into common use. These methods share with the Galerkin approach the weak enforcement of the differential equation and of certain boundary conditions. In their original version the unknowns are the values of the solution at the quadrature points, as in a collocation method. We shall refer to such approaches as Galerkin with numerical integration, or G-NI, methods.

The first unifying mathematical assessment of the theory of spectral methods was provided in the monograph by Gottlieb and Orszag (1977). The theory was extended to cover a large variety of problems, such as variable-coefficient and nonlinear equations. A sound approximation theory for the polynomial families used in spectral methods was developed. In his monograph Mercier (1981) advanced the understanding of the role of Gaussian quadrature points for orthogonal polynomials as collocation points for spectral methods, as had originally been observed in 1979 by Gottlieb. Stability and convergence analyses for spectral methods were produced for a variety of approaches. The theoretical analysis of spectral methods in terms of weak formulations proved very successful. As a matter of fact, this opened the door to the use of functional analysis techniques to handle complex problems and to obtain the sharpest results. Application developments were equally extensive, and by the late 1980's spectral methods had become the predominant numerical tool for basic flow physics investigations of transition and turbulence. All in all, the 10 years that followed were extremely fruitful for the theoretical development and the application deployment of spectral methods.

Developments of the first five years that followed Gottlieb and Orszag (1977) were reviewed in the symposium proceedings edited by Voigt, Gottlieb and Hussaini (1984). Indeed, that very symposium in 1982 inspired the youthful incarnations of the present authors to produce their first text on this subject (Canuto, Hussaini, Quarteroni and Zang (1988)). Subsequently, numerous other texts and review articles on various aspects of spectral methods appeared. Boyd (1989, and especially the 2001 second edition) contains a wealth of detail and advice on spectral algorithms and is an especially good reference for problems on unbounded domains and in cylindrical and spherical

coordinate systems. A sound reference for the theoretical aspects of spectral methods for elliptic equations was provided by Bernardi and Maday (1992b, 1997). Funaro (1992) and Guo (1998) discussed the approximation of differential equations by polynomial expansions. Fornberg (1996) is a guide for the practical application of spectral collocation methods, and it contains illustrative examples, heuristic explanations, basic Fortran code segments, and a succinct chapter on applications to turbulent flows and weather prediction. Trefethen (2000) is a lively introduction to spectral collocation methods and includes copious examples in Matlab. Focused applications of spectral methods on particular classes of problems were provided by Tadmor (1998) and Gottlieb and Hesthaven (2001) for first-order hyperbolic problems, by Cohen (2002) for wave equations, and by Bernardi, Dauge and Maday (1999) for problems in axisymmetric domains. Peyret (2002) provided a rather comprehensive discussion of Fourier and Chebyshev spectral methods for the solution of the incompressible Navier-Stokes equations, specifically in the primitive equations and vorticity-streamfunction formulations.

By the late 1980's classical spectral methods were reasonably mature, and the research focus had clearly shifted to the use of high-order methods for problems on complex domains. We shall refer to this class of spectral methods generically as *multidomain spectral methods* or as *spectral methods in arbitrary geometries*. The 1988 book by the present authors closed with an overview of this then nascent subject. Funaro (1997) treats spectral-element methods in the context of elliptic boundary-value problems, especially convection-dominated flows, and includes a multidomain treatment for complex geometry. The first comprehensive texts on spectral methods in complex domains appeared around the year 2000. Karniadakis and Sherwin (1999) provides a unified framework for spectral-element methods (as introduced by Patera (1984)) and *hp* finite-element methods (see, for example, Babuška, Szabó and Katz (1981)). It includes structured and unstructured domains, and applications to both incompressible and compressible flows. The Deville, Fischer and Mund (2002) text focuses on high-order methods in physical space (collocation and spectral-element methods) with applications to incompressible flows. Its coverage of the implementation details of such methods on vector and parallel computers distinguishes it from other books on the subject. Although specifically devoted to the *hp*-version of finite-element methods, the book by Schwab (1998) provides many useful theoretical results about the approximation properties of high-order polynomials in complex domains.

The present book is focused on the fundamentals of spectral methods on simple domains. A companion book (Canuto, Hussaini, Quarteroni and Zang (2007)) discusses specific spectral algorithms for fluid dynamics applications and describes the evolution of spectral methods to complex domains. We shall refer to the companion book as CHQZ3. Citations in the present text that refer to specific material in the companion book will have the format

CHQZ3, Chap. x or CHQZ3, Sect. x.y. For example, a reference such as CHQZ3, Chap. 1 refers to Chapter 1 of Canuto, Hussaini, Quarteroni and Zang (2007).

1.2 Some Examples of Spectral Methods

Spectral methods are distinguished not only by the fundamental type of the method (Galerkin, collocation, Galerkin with numerical integration, or tau), but also by the particular choice of the trial functions. The most frequently used trial functions are trigonometric polynomials, Chebyshev polynomials, and Legendre polynomials. In this section we shall illustrate the basic principles of each method and the basic properties of each set of polynomials by examining in detail one particular spectral method on each of several different types of differential equations. Each of these examples will be reconsidered in Chap. 6 from a rigorous theoretical point of view.

1.2.1 A Fourier Galerkin Method for the Wave Equation

Many evolution equations can be written as

$$\frac{\partial u}{\partial t} = \mathcal{M}(u) , \quad (1.2.1)$$

where $u(\mathbf{x}, t)$ is the solution, and $\mathcal{M}(u)$ is an operator (linear or nonlinear) that contains all the spatial derivatives of u . Equation (1.2.1) must be coupled with an initial condition $u(\mathbf{x}, 0)$ and suitable boundary conditions.

For simplicity suppose that there is only one spatial dimension, that the spatial domain is $(0, 2\pi)$, and that the boundary conditions are periodic. Most often spectral methods are used only for the spatial discretization. The approximate solution is represented as

$$u^N(x, t) = \sum_{k=-N/2}^{N/2} a_k(t) \phi_k(x) . \quad (1.2.2)$$

The ϕ_k are the trial functions, whereas the a_k are the expansion coefficients. In general, u^N will not satisfy (1.2.1), i.e., the residual

$$\frac{\partial u^N}{\partial t} - \mathcal{M}(u^N)$$

will not vanish everywhere. The approximation is obtained by selecting a set of test functions ψ_k and by requiring that

$$\int_0^{2\pi} \left[\frac{\partial u^N}{\partial t} - \mathcal{M}(u^N) \right] \psi_k(x) dx = 0 , \quad (1.2.3)$$

for $k = -N/2, \dots, N/2$, where the test functions determine the weights of the residual. In this sense the approximation is obtained by a method of weighted residuals. Most often the numerical analysis community describes discretizations of differential equations formulated by integral expressions such as (1.2.3) (possibly after applying integration-by-parts) as discrete *weak formulations*. This more common terminology is the one that we follow in this text. The alternative, discrete *strong formulation* is characterized by enforcing that the approximate representation of the solution, e.g., (1.2.2), satisfy the differential equation exactly at a discrete set of points. Finite-difference methods use a strong formulation, as do spectral collocation methods – see the example in Sect. 1.2.2. A more comprehensive discussion of alternative formulations of differential problems is provided in Sect. 3.2.

The most straightforward spectral method for a problem with periodic boundary conditions is based on trigonometric polynomials:

$$\phi_k(x) = e^{ikx}, \quad (1.2.4)$$

$$\psi_k(x) = \frac{1}{2\pi} e^{-ikx}. \quad (1.2.5)$$

Note that the trial functions and the test functions are essentially the same, and that they satisfy the (bi-)orthonormality condition

$$\int_0^{2\pi} \phi_k(x) \psi_l(x) dx = \delta_{kl}. \quad (1.2.6)$$

If this were merely an approximation problem, then (1.2.2) would be the truncated Fourier series of the known function $u(x, t)$ with

$$a_k(t) = \int_0^{2\pi} u(x, t) \psi_k(x) dx \quad (1.2.7)$$

being simply the familiar Fourier coefficients. For the partial differential equation (PDE), however, $u(x, t)$ is not known; the approximation (1.2.2) is determined by (1.2.3).

For the linear hyperbolic problem

$$\frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} = 0, \quad (1.2.8)$$

i.e., for

$$\mathcal{M}(u) = \frac{\partial u}{\partial x}, \quad (1.2.9)$$

condition (1.2.3) becomes

$$\frac{1}{2\pi} \int_0^{2\pi} \left[\left(\frac{\partial}{\partial t} - \frac{\partial}{\partial x} \right) \sum_{l=-N/2}^{N/2} a_l(t) e^{ilx} \right] e^{-ikx} dx = 0.$$

The next two steps are the analytical (spatial) differentiation of the trial functions:

$$\frac{1}{2\pi} \int_0^{2\pi} \left[\sum_{l=-N/2}^{N/2} \left(\frac{da_l}{dt} - ila_l \right) e^{ilx} \right] e^{-ikx} dx = 0 ,$$

and the analytical integration of this expression, which produces the dynamical equations

$$\frac{da_k}{dt} - ik a_k = 0 , \quad k = -N/2, \dots, N/2 . \quad (1.2.10)$$

The initial conditions for this system of ordinary differential equations (ODEs) are the coefficients for the expansion of the initial condition. For this Galerkin approximation,

$$a_k(0) = \int_0^{2\pi} u(x, 0) \psi_k(x) dx . \quad (1.2.11)$$

For the strict Galerkin method, integrals such as those that appear in (1.2.11) should be computed analytically. For the simple example problem of this subsection this integration can indeed be performed analytically. For more complicated problems, however, numerical quadratures are performed. This is discussed further in Sect. 1.2.3.

We shall use the initial condition

$$u(x, 0) = \sin(\pi \cos x) \quad (1.2.12)$$

to illustrate the accuracy of the Fourier Galerkin method for (1.2.8). The exact solution,

$$u(x, t) = \sin[\pi \cos(x + t)] , \quad (1.2.13)$$

has the Fourier expansion

$$u(x, t) = \sum_{k=-\infty}^{\infty} a_k(t) e^{ikx} , \quad (1.2.14)$$

where the Fourier coefficients are

$$a_k(t) = \sin\left(\frac{k\pi}{2}\right) J_k(\pi) e^{ikt} , \quad (1.2.15)$$

and $J_k(t)$ is the Bessel function of order k .

The asymptotic properties of the Bessel functions imply that

$$k^p a_k(t) \rightarrow 0 \quad \text{as} \quad k \rightarrow \infty \quad (1.2.16)$$

for all positive integers p . As a result, the truncated Fourier series,

$$u^N(x, t) = \sum_{k=-N/2}^{N/2} a_k(t) e^{ikx}, \quad (1.2.17)$$

converges faster than any finite power of $1/N$. This property is often referred to as spectral convergence.

An illustration of the superior accuracy available from the spectral method for this problem is provided in Fig. 1.1. Shown in the figure are the maximum errors after one period at $t = 2\pi$ for the spectral Galerkin method, a second-order finite-difference method, an (explicit) fourth-order finite-difference method, a fourth-order compact method, and a sixth-order compact method. The integer N denotes the degree of the expansion (1.2.17) for the Fourier Galerkin method and the number of grid points for the finite-difference and compact methods. The time discretization was the classical fourth-order Runge-Kutta method and the exact initial Fourier coefficients were used for the spectral method. In all cases the time-step was chosen so small that the temporal discretization error was negligible. (Appendix D furnishes the formulas (and stability regions) for commonly used time discretizations. The familiar formula for the classical fourth-order Runge-Kutta methods is given in (D.2.17).)

The second-order and fourth-order finite-difference methods used here and elsewhere in this book for examples are the standard central-difference methods with 3-point and 5-point explicit stencils, respectively. The fourth-order and sixth-order compact methods used in our examples are the classical 3-point Padé approximations (see, for example, Collatz (1966) and Lele (1992))

$$u'_{j-1} + 4u'_j + u'_{j+1} = \frac{3}{\Delta x}(u_{j+1} - u_{j-1}) \quad (1.2.18)$$

and

$$u'_{j-1} + 3u'_j + u'_{j+1} = \frac{7}{3\Delta x}(u_{j+1} - u_{j-1}) + \frac{1}{12\Delta x}(u_{j+2} - u_{j-2}), \quad (1.2.19)$$

respectively, where Δx is the grid spacing and u'_j denotes the approximation to the first derivative at $x_j = j\Delta x$. Of course, when nonperiodic boundary conditions are present, special stencils are needed for points at, and sometimes also adjacent to, the boundary.

Figure 1.2 compares these various numerical solutions for $N = 16$ with the exact answer. Note that the major errors in the finite-difference solutions are ones of *phase* rather than *amplitude*. In many problems the very low phase error of spectral methods is a significant advantage.

Because the solution is infinitely smooth, the convergence of the spectral method on this problem is more rapid than any finite power of $1/N$. Actually, since the solution is analytic, convergence is exponentially fast. (The errors for the $N \geq 64$ spectral results are so small that they are swamped by the round-off error of these calculations. Unless otherwise noted, all numerical examples presented in this book were performed in 64-bit arithmetic.)

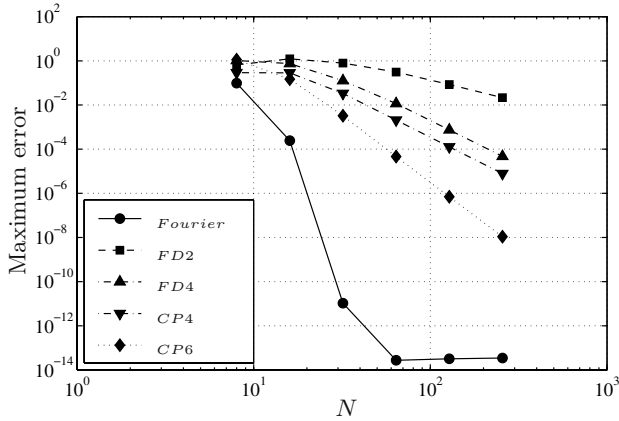


Fig. 1.1. Maximum errors for the linear hyperbolic problem at $t = 2\pi$ for Fourier Galerkin and several finite-difference schemes

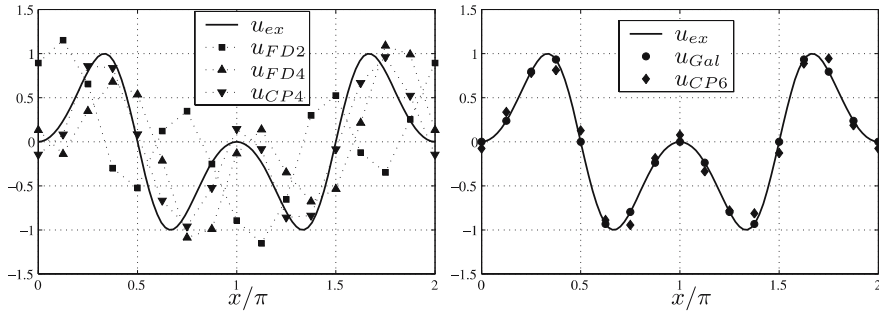


Fig. 1.2. Numerical solutions for the linear hyperbolic problem at $t = 2\pi$ for $N = 16$ for Fourier Galerkin and several finite-difference schemes

In most practical applications the benefit of the spectral method is not the extraordinary accuracy available for large N but rather the small size of N necessary for a moderately accurate solution.

1.2.2 A Chebyshev Collocation Method for the Heat Equation

Fourier series, despite their simplicity and familiarity, are not always a good choice for the trial functions. In fact, for reasons that will be explored in the next chapter, Fourier series are only advisable for problems with periodic boundary conditions. A more versatile set of trial functions is composed of the Chebyshev polynomials. These are defined on $[-1, 1]$ by

$$T_k(x) = \cos(k \cos^{-1} x), \quad (1.2.20)$$

for $k = 0, 1, \dots$

Let us focus on the linear heat equation

$$\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = 0, \quad (1.2.21)$$

i.e.,

$$\mathcal{M}(u) = \frac{\partial^2 u}{\partial x^2}, \quad (1.2.22)$$

on $(-1, 1)$ with homogeneous Dirichlet boundary conditions,

$$u(-1, t) = 0, \quad u(1, t) = 0. \quad (1.2.23)$$

Choosing the trial functions

$$\phi_k(x) = T_k(x), \quad k = 0, 1, \dots, N, \quad (1.2.24)$$

the approximate solution has the representation

$$u^N(x, t) = \sum_{k=0}^N a_k(t) \phi_k(x). \quad (1.2.25)$$

In the collocation approach the requirement is that (1.2.21) be satisfied exactly by (1.2.25) at a set of collocation points x_j in $(-1, 1)$:

$$\left. \frac{\partial u^N}{\partial t} - \mathcal{M}(u^N) \right|_{x=x_j} = 0, \quad j = 1, \dots, N-1. \quad (1.2.26)$$

The boundary conditions

$$u^N(-1, t) = 0, \quad u^N(1, t) = 0 \quad (1.2.27)$$

and the initial condition

$$u^N(x_k, 0) = u(x_k, 0), \quad k = 0, \dots, N, \quad (1.2.28)$$

accompany (1.2.26).

Equations (1.2.26) are based on the strong formulation of the differential equation, since the approximate solution is required to satisfy the differential equation exactly at a set of discrete points, in this case called the collocation points. One can formally obtain the same equations starting from a weak formulation of the problem by taking as test functions the (shifted) Dirac delta-functions (distributions)

$$\psi_j(x) = \delta(x - x_j), \quad j = 1, \dots, N-1, \quad (1.2.29)$$

and enforcing the conditions

$$\int_{-1}^1 \left[\frac{\partial u^N}{\partial t} - \mathcal{M}(u^N) \right] \psi_j(x) dx = 0, \quad j = 1, \dots, N-1 \quad (1.2.30)$$

(where the integral should really be interpreted as a duality; see (A.10)).

A particularly convenient choice for the collocation points x_j is

$$x_j = \cos \frac{\pi j}{N} . \quad (1.2.31)$$

Not only does this choice produce highly accurate approximations, but it also is economical. Note that

$$\phi_k(x_j) = \cos \frac{\pi j k}{N} . \quad (1.2.32)$$

This enables the Fast Fourier Transform (FFT) to be employed in the evaluation of $\mathcal{M}(u^N)|_{x=x_j}$, as is discussed in Sect. 2.4.

For the particular initial condition

$$u(x, 0) = \sin \pi x , \quad (1.2.33)$$

the exact solution is

$$u(x, t) = e^{-\pi^2 t} \sin \pi x . \quad (1.2.34)$$

It has the infinite Chebyshev expansion

$$u(x, t) = \sum_{k=0}^{\infty} b_k(t) T_k(x) , \quad (1.2.35)$$

where

$$b_k(t) = \frac{2}{c_k} \sin \left(\frac{k\pi}{2} \right) J_k(\pi) e^{-\pi^2 t} , \quad (1.2.36)$$

with

$$c_k = \begin{cases} 2 , & k = 0 , \\ 1 , & k \geq 1 . \end{cases} \quad (1.2.37)$$

Because of the rapidly decaying $J_k(\pi)$ factor, the truncated series converges at an exponential rate. A well-designed collocation method will do the same. (Since the finite series (1.2.25) is not simply the truncation of the infinite series (1.2.35) at order N , the expansion coefficients $a_k(t)$ and $b_k(t)$ are not identical.)

Unlike a Galerkin method, which in its conventional version is usually implemented in terms of the expansion coefficients $a_k(t)$, a collocation method is implemented in terms of the nodal values $u_j(t) = u^N(x_j, t)$. Indeed, in addition to (1.2.25), we have the expansion

$$u^N(x, t) = \sum_{j=0}^N u_j(t) \phi_j(x) ,$$

where now ϕ_j denote the discrete (shifted) delta-functions, i.e., the unique N -th degree polynomials satisfying $\phi_j(x_i) = \delta_{ij}$ for $0 \leq i, j \leq N$.

(These particular functions will be more commonly denoted by the symbol ψ_j in the sequel and referred to as characteristic Lagrange polynomials; see, e.g., (1.2.55)). The expansion coefficients are used only in an intermediate step, namely, in the analytic differentiation (with respect to x) of (1.2.25). The details of this step, which will be derived in Sect. 2.4, follow.

The expansion coefficients are given by

$$a_k(t) = \frac{2}{N\bar{c}_k} \sum_{l=0}^N \bar{c}_l^{-1} u_l(t) \cos \frac{\pi l k}{N}, \quad k = 0, 1, \dots, N, \quad (1.2.38)$$

where

$$\bar{c}_k = \begin{cases} 2, & k = 0 \text{ or } N, \\ 1, & 1 \leq k \leq N-1 \end{cases}. \quad (1.2.39)$$

The exact derivative of (1.2.25) is

$$\frac{\partial^2 u^N}{\partial x^2}(t) = \sum_{k=0}^N a_k^{(2)}(t) T_k(x), \quad (1.2.40)$$

where

$$\begin{aligned} a_{N+1}^{(1)}(t) &= 0, & a_N^{(1)}(t) &= 0, \\ \bar{c}_k a_k^{(1)}(t) &= a_{k+2}^{(1)}(t) + 2(k+1)a_{k+1}^{(1)}(t), & k &= N-1, N-2, \dots, 0, \end{aligned} \quad (1.2.41)$$

and

$$\begin{aligned} a_{N+1}^{(2)}(t) &= 0, & a_N^{(2)}(t) &= 0, \\ \bar{c}_k a_k^{(2)}(t) &= a_{k+2}^{(2)}(t) + 2(k+1)a_{k+1}^{(1)}(t), & k &= N-1, N-2, \dots, 0. \end{aligned} \quad (1.2.42)$$

The coefficients $a_k^{(2)}$ obviously depend linearly on the nodal values u_l ; hence, there exists a matrix D_N^2 such that

$$\left. \frac{\partial^2 u^N}{\partial x^2}(t) \right|_{x=x_j} = \sum_{k=0}^N a_k^{(2)}(t) \cos \frac{\pi j k}{N} = \sum_{l=0}^N (D_N^2)_{jl} u_l(t) \quad (1.2.43)$$

(see Sect. 2.4.2 for more details). By (1.2.27), we actually have $u_0(t) = u_N(t) = 0$. Substituting the above expression into (1.2.26), we end up with a system of ordinary differential equations for the nodal unknowns:

$$\frac{du_j}{dt}(t) = \sum_{l=0}^N (D_N^2)_{jl} u_l(t), \quad j = 1, \dots, N-1. \quad (1.2.44)$$

Supplemented by the initial conditions (1.2.28), the preceding system of ordinary differential equations for the nodal values of the solution is readily integrated in time.

The maximum errors at $t = 1$ in the numerical solutions for a Chebyshev collocation method, a second-order finite-difference method and a fourth-order compact method are given in Fig. 1.3, along with the maximum errors for the truncated Chebyshev series of the exact solution at $t = 1$. The Chebyshev method used the $N + 1$ non-uniformly distributed collocation points (1.2.31), whereas the finite-difference methods used $N + 1$ uniformly distributed points. The maximum errors have been normalized with respect to the maximum value of the exact solution at $t = 1$. The fourth-order scheme is the classical 3-point Padé approximation,

$$u''_{i-1} + 10u''_i + u''_{i+1} = \frac{12}{(\Delta x)^2}(u_{i-1} - 2u_i + u_{i+1}), \quad i = 1, \dots, N-1, \quad (1.2.45)$$

supplemented with a compact, third-order approximation at the boundary points (see Lele (1992)), e.g.,

$$u''_0 + 11u''_1 = \frac{1}{(\Delta x)^2}(13u_0 - 27u_1 + 15u_2 - u_3), \quad i = 0. \quad (1.2.46)$$

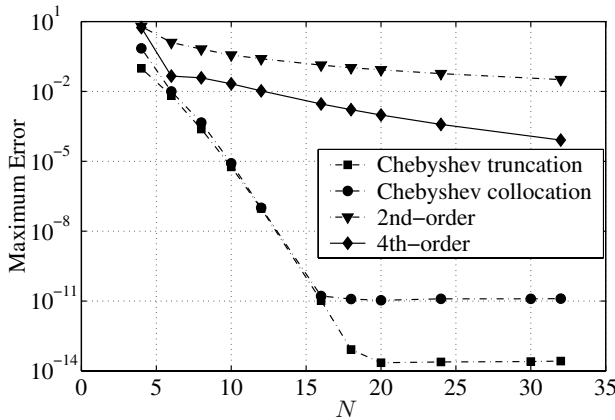


Fig. 1.3. Maximum errors for the heat equation problem at $t = 1$ for Chebyshev collocation and several finite-difference schemes. The Chebyshev truncation result is shown for comparison

Before leaving this example, we consider a more general equation than (1.2.21), namely,

$$\frac{\partial u}{\partial t} - \frac{\partial}{\partial x} \left(\kappa \frac{\partial u}{\partial x} \right) = 0, \quad (1.2.47)$$

where the conductivity coefficient κ varies in $(-1, 1)$ and may even depend on the solution u . In this case, it is not convenient to apply the collocation scheme (1.2.26) to equation (1.2.47) directly, as this would require the exact differentiation of the heat flux $\mathcal{F}(u^N) = \kappa \frac{\partial u^N}{\partial x}$. Instead, one first computes the nodal values $F_l(t) = \mathcal{F}(u^N)(x_l)$, $l = 0, \dots, N$, of this flux, then applies a transformation similar to (1.2.38), and follows that with a differentiation of the flux as in (1.2.41); the resulting expansion of the derivative is then evaluated at the collocation points. This process amounts to differentiating exactly the numerical flux $\mathcal{F}^N(u^N) = I_N(\mathcal{F}(u^N))$, which is obtained by interpolating the flux $\mathcal{F}(u^N)$ at the collocation points by a global N -degree algebraic polynomial. (Here and in the rest of the book, I_N is a general symbol that denotes an interpolation operator.) The resulting collocation scheme reads as follows:

$$\left. \frac{\partial u^N}{\partial t} - \frac{\partial}{\partial x} I_N \left(\kappa \frac{\partial u^N}{\partial x} \right) \right|_{x=x_j} = 0, \quad j = 1, \dots, N-1. \quad (1.2.48)$$

Equivalently, we have

$$\frac{du_j}{dt}(t) = \sum_{l=0}^N (D_N)_{jl} F_l(t), \quad j = 1, \dots, N-1, \quad (1.2.49)$$

where D_N is the Chebyshev collocation derivative matrix, which is discussed in detail in Sect. 2.4.2.

The approach used for the discretization of (1.2.47) highlights a general strategy that is adopted for collocation methods: differentiation is applied to a function only after the argument of the function is interpolated by a global polynomial at a suitable set of collocation points. Obviously, when the argument is itself a polynomial of degree $\leq N$, as in the constant-coefficient heat equation (1.2.21), the interpolation returns the value of the argument.

1.2.3 A Legendre Galerkin with Numerical Integration (G-NI) Method for the Advection-Diffusion-Reaction Equation

Spectral methods are also applicable to time-independent equations. The general boundary-value problem is given by the equation

$$\mathcal{M}(u) = f \quad (1.2.50)$$

to be solved in a specified domain, along with the boundary conditions

$$\mathcal{B}(u) = 0. \quad (1.2.51)$$

As a first example, we consider the one-dimensional advection-diffusion-reaction equation

$$\mathcal{M}(u) = \frac{d\mathcal{F}(u)}{dx} + \gamma u = f, \quad (1.2.52)$$

where the advection-diffusion flux is defined as

$$\mathcal{F}(u) = -\nu \frac{du}{dx} + \beta u.$$

The domain for the equation is $(-1, 1)$, and the boundary conditions are

$$\mathcal{B}_1(u) = u(-1) = 0, \quad (1.2.53a)$$

$$\mathcal{B}_2(u) = \mathcal{F}(u)(1) + g = 0. \quad (1.2.53b)$$

We assume that the coefficients ν , β and γ as well as the data f may vary in the domain, and that the diffusion coefficient satisfies $\nu \geq \bar{\nu}$ for some constant $\bar{\nu} > 0$.

Trial and test functions are defined as follows. Consider the N -th degree Legendre orthogonal polynomial $L_N(x)$. (A detailed discussion of the properties of Legendre polynomials is furnished in Sect. 2.3.) The polynomial L_N has $N - 1$ extrema x_j , i.e., $L'_N(x_j) = 0$, for $j = 1, \dots, N - 1$; they belong to the interval $(-1, 1)$. Adding the boundary points $x_0 = -1$ and $x_N = 1$, we obtain $N + 1$ points, which are high-precision quadrature nodes (they are termed the *Legendre Gauss-Lobatto nodes*); indeed, there exist weights w_j such that the quadrature formula

$$\int_{-1}^1 p(x) dx \sim \sum_{j=0}^N p(x_j) w_j \quad (1.2.54)$$

is exact for all polynomials p of degree $\leq 2N - 1$. Based on these nodes, we now introduce the *characteristic Lagrange polynomials*

$$\psi_j(x) = \frac{1}{N(N+1)} \frac{(1-x^2)}{(x_j-x)} \frac{L'_N(x)}{L'_N(x_j)}, \quad j = 0, \dots, N, \quad (1.2.55)$$

which are discrete (shifted) delta-functions, i.e., they are N -th degree polynomials which approximate the (shifted) Dirac delta-functions $\delta(x - x_j)$, as they satisfy

$$\psi_j(x_k) = \delta_{jk}, \quad j, k = 0, \dots, N. \quad (1.2.56)$$

In view of the boundary condition (1.2.53a), we drop ψ_0 . The remaining functions ψ_j , $j = 1, \dots, N$, will be our trial and test functions. The approximate solution is sought in the form

$$u^N(x) = \sum_{l=1}^N u_l \psi_l(x). \quad (1.2.57)$$

Note that the coefficients in the expansion are precisely the values of u^N at the nodes $u_l = u^N(x_l)$, $l = 1, \dots, N$.

In order to arrive at the equations which uniquely define u^N , we have to go back to the exact solution u of our boundary-value problem. We shall derive a set of integral conditions satisfied by the exact solution (which constitute the weak formulation of the problem). The same integral conditions are enforced on the discrete solution. To this end, consider (1.2.52), multiply both sides by any test function ψ_j and integrate over the interval $(-1, 1)$; we obtain the equations

$$\int_{-1}^1 \frac{d\mathcal{F}(u)}{dx} \psi_j dx + \int_{-1}^1 \gamma u \psi_j dx = \int_{-1}^1 f \psi_j dx, \quad j = 1, \dots, N. \quad (1.2.58)$$

Integrating the first term by parts, we get

$$\begin{aligned} \int_{-1}^1 \frac{d\mathcal{F}(u)}{dx} \psi_j dx &= - \int_{-1}^1 \mathcal{F}(u) \frac{d\psi_j}{dx} dx + [\mathcal{F}(u) \psi_j]_{-1}^1 \\ &= - \int_{-1}^1 \mathcal{F}(u) \frac{d\psi_j}{dx} dx - g \delta_{jN}, \end{aligned}$$

where we have used the boundary condition (1.2.53b), as well as the relations (1.2.56). Thus, recalling the definition of the flux $\mathcal{F}(u)$, we see that u satisfies

$$\begin{aligned} \int_{-1}^1 \nu \frac{du}{dx} \frac{d\psi_j}{dx} dx - \int_{-1}^1 \beta u \frac{d\psi_j}{dx} dx + \int_{-1}^1 \gamma u \psi_j dx & \quad (1.2.59) \\ = \int_{-1}^1 f \psi_j dx + g \delta_{jN}, & \quad j = 1, \dots, N. \end{aligned}$$

This is precisely the set of equations which we ask to be satisfied by u^N as well. If we replace u by u^N in (1.2.59), we obtain the numerical scheme

$$\begin{aligned} \int_{-1}^1 \nu \frac{du^N}{dx} \frac{d\psi_j}{dx} dx - \int_{-1}^1 \beta u^N \frac{d\psi_j}{dx} dx + \int_{-1}^1 \gamma u^N \psi_j dx & \quad (1.2.60) \\ = \int_{-1}^1 f \psi_j dx + g \delta_{jN}, & \quad j = 1, \dots, N. \end{aligned}$$

Note that u^N satisfies (1.2.53a) exactly; conversely, (1.2.53b) is not enforced directly on u^N , yet it has been incorporated into (1.2.59). We say that we enforce this boundary condition in a *weak*, or *natural*, manner.

Since the integrals in (1.2.59) are evaluated exactly, we have obtained a *pure Galerkin* scheme. However, only in special situations (e.g., constant coefficients and data) can the integrals above be computed analytically. Otherwise, we have to resort to numerical integration, in which case the natural choice is the quadrature formula (1.2.54). In this way, we obtain the following modified scheme, which we term the *Galerkin with numerical integration scheme*, or in short, the G-NI scheme:

$$\begin{aligned}
& \sum_{k=0}^N \left(\nu \frac{du^N}{dx} \frac{d\psi_j}{dx} \right) (x_k) w_k - \sum_{k=0}^N \left(\beta u^N \frac{d\psi_j}{dx} \right) (x_k) w_k + \sum_{k=0}^N (\gamma u^N \psi_j) (x_k) w_k \\
&= \sum_{k=0}^N (f \psi_j) (x_k) w_k + g \delta_{jN}, \quad j = 1, \dots, N.
\end{aligned} \tag{1.2.61}$$

Inserting the expansion (1.2.57) for u^N , we can rephrase this scheme as a system $K\mathbf{u} = \mathbf{b}$ of N algebraic equations in the unknowns u_l ; in particular, they are

$$\sum_{l=1}^N K_{jl} u_l = b_j, \quad j = 1, \dots, N, \tag{1.2.62}$$

where the matrix entries are

$$K_{jl} = \sum_{k=0}^N \left(\nu \frac{d\psi_l}{dx} \frac{d\psi_j}{dx} \right) (x_k) w_k - \left(\beta \frac{d\psi_j}{dx} \right) (x_l) w_l + \gamma(x_j) w_j \delta_{lj},$$

and the right-hand side components are

$$b_j = f(x_j) w_j + g \delta_{jN}.$$

Efficient solution techniques for such a system are described in Sect. 4.2.

The G-NI scheme can be given a pointwise, or collocation-like, interpretation, which serves to highlight the effect of the weak enforcement of the boundary condition (1.2.53b). To this end, we denote by $I_N \varphi$ the N -th degree algebraic polynomial that interpolates a function φ at the Gauss-Lobatto nodes x_j , $j = 0, \dots, N$; this allows us to introduce the numerical flux

$$\mathcal{F}^N(u^N) = I_N(\mathcal{F}(u^N)).$$

The two first sums in (1.2.61) can be written as

$$\begin{aligned}
& \sum_{k=0}^N \left(\nu \frac{du^N}{dx} \frac{d\psi_j}{dx} \right) (x_k) w_k - \sum_{k=0}^N \left(\beta u^N \frac{d\psi_j}{dx} \right) (x_k) w_k = \\
&= - \sum_{k=0}^N \left(\mathcal{F}(u^N) \frac{d\psi_j}{dx} \right) (x_k) w_k = - \sum_{k=0}^N \left(\mathcal{F}^N(u^N) \frac{d\psi_j}{dx} \right) (x_k) w_k.
\end{aligned}$$

Now it is crucial to observe that both the terms $\mathcal{F}^N(u^N) \frac{d\psi_j}{dx}$ and $\frac{d\mathcal{F}^N(u^N)}{dx} \psi_j$ are polynomials of degree $\leq 2N-1$; hence, they can be integrated exactly by the quadrature formula (1.2.54). Thus, we are allowed to counter-integrate by parts in the last sum appearing above, obtaining

$$\begin{aligned}
-\sum_{k=0}^N \left(\mathcal{F}^N(u^N) \frac{d\psi_j}{dx} \right) (x_k) w_k &= -\int_{-1}^1 \mathcal{F}^N(u^N) \frac{d\psi_j}{dx} dx \\
&= \int_{-1}^1 \frac{d\mathcal{F}^N(u^N)}{dx} \psi_j dx - [\mathcal{F}^N(u^N) \psi_j]_{-1}^1 \\
&= \sum_{k=0}^N \left(\frac{d\mathcal{F}^N(u^N)}{dx} \psi_j \right) (x_k) w_k - \mathcal{F}(u^N)(1) \psi_j(1).
\end{aligned}$$

If we insert this expression into (1.2.61) and use the relations (1.2.56), we obtain the following equivalent formulation of the G-NI scheme:

$$\left(\frac{d\mathcal{F}^N(u^N)}{dx} + \gamma u^N \right) (x_j) w_j - \mathcal{F}(u^N)(1) \delta_{jN} = f(x_j) w_j + g \delta_{jN}, \quad j = 1, \dots, N. \quad (1.2.63)$$

For $j = 1, \dots, N-1$, this is simply

$$\frac{d\mathcal{F}^N(u^N)}{dx} + \gamma u^N - f \Big|_{x=x_j} = 0, \quad (1.2.64)$$

i.e., at the internal quadrature points we are collocating the differential equation after replacing the exact flux $\mathcal{F}(u^N)$ by the numerical one $\mathcal{F}^N(u^N)$. For $j = N$ we get

$$\frac{d\mathcal{F}^N(u^N)}{dx} + \gamma u^N - f \Big|_{x=1} - \frac{1}{w_N} (\mathcal{F}(u^N) + g) \Big|_{x=1} = 0, \quad (1.2.65)$$

i.e., at $x = 1$ we are collocating a particular linear combination of the discrete form of the differential equation and the boundary condition. Since $1/w_N$ grows like N^2 as $N \rightarrow \infty$ (see Sect. 2.3.1), (1.2.65) shows that the boundary condition is approximately fulfilled in a more and more accurate way as the equation residual $\mathcal{M}^N(u^N) - f|_{x=1}$ gets smaller and smaller for $N \rightarrow \infty$ (recall that the residual vanishes at all internal nodes, see (1.2.64)).

The example addressed above is indeed a paradigm for a general class of second-order steady problems. The G-NI discretization consists of collocating the differential equation (with numerical flux) at the internal Gauss Lobatto nodes; Dirichlet boundary conditions (i.e., conditions involving only pointwise values of the unknown function) are fulfilled exactly at the boundary points, whereas Neumann or Neumann-like boundary conditions (i.e., conditions involving also the first derivative(s) of the unknown function) are enforced via an intrinsically (and unambiguously) defined penalty method.

The accuracy of the G-NI method is illustrated by the following example. We consider the problem (1.2.50)–(1.2.53) in the interval $(-1, 1)$ with $\nu = 1$, $\beta(x) = \cos(\pi/4 \cdot (1+x))$ and $\gamma = 1$. The right-hand side $f(x)$ and the datum g are computed so that the exact solution is

$$u(x) = \cos(3\pi(1+x)) \sin(\pi/5 \cdot (x+0.5)) + \sin(\pi/10). \quad (1.2.66)$$

For several values of N , we denote by u^N the G-NI solution (N is the polynomial degree) and by u^p ($p = 1, 2, 3$) the (piecewise-polynomial) finite-element solution corresponding to a subdivision in subintervals of equal size. In all cases, $N + 1$ denotes the total number of nodal values. In Fig. 1.4 (left) we plot the maximum error of the solution, while on the right we plot the absolute error of the boundary flux $|(\nu \frac{du^p}{dx}(1) + \beta u^p(1)) - g|$ for $p = 1, 2, 3, N$. The two errors exhibit a similar decay with respect to N . In particular, the boundary condition at $x = 1$ is fulfilled with spectral accuracy.

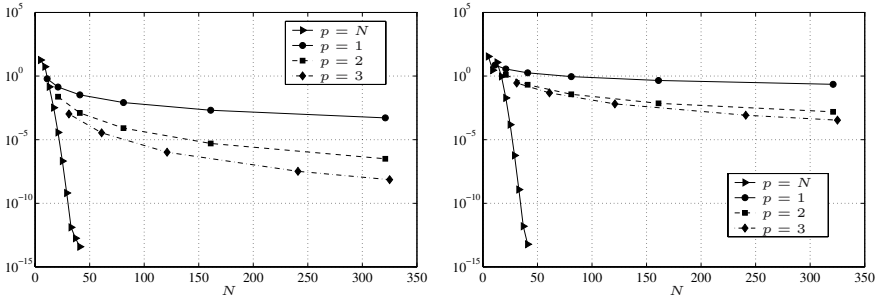


Fig. 1.4. Comparison between the accuracy of the G-NI solution (corresponding to the curve $p = N$) and the finite-element solutions of order $p = 1, 2$ and 3 versus N which represents the total number of nodal values. The maximum error between the numerical solution and the exact one $u(x) = \cos(3\pi(1+x)) \cdot \sin(\pi/5 \cdot (x+0.5)) + \sin(\pi/10)$ (left) and the absolute value of the error on the flux at $x = 1$ (right)

1.2.4 A Legendre Tau Method for the Poisson Equation

Our second example of a steady boundary-value problem is the Poisson equation on $(-1, 1) \times (-1, 1)$, with homogeneous Dirichlet boundary conditions. The choice of \mathcal{M} and \mathcal{B} in (1.2.50) and (1.2.51) is as follows:

$$\mathcal{M}(u) = - \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right), \quad (1.2.67)$$

$$\mathcal{B}_1(u) = u(x, -1), \quad (1.2.68a)$$

$$\mathcal{B}_2(u) = u(x, +1), \quad (1.2.68b)$$

$$\mathcal{B}_3(u) = u(-1, y), \quad (1.2.68c)$$

$$\mathcal{B}_4(u) = u(+1, y). \quad (1.2.68d)$$

(We prefer to use the negative sign in second-derivative operators such as (1.2.67) so that $\mathcal{M}(u)$ is a positive, rather than a negative, operator. Although this might be disconcerting to some, it does simplify the discussion of the mathematical properties of the operator and its numerical approximations. For example, some spectral approximations to (1.2.67)–(1.2.68) yield

symmetric and positive-definite matrices, albeit not the particular approximation discussed in the present subsection. This will become clearer in due course, particularly in Chaps. 4, 6 and 7.)

Both Legendre and Chebyshev polynomials are suitable trial functions. A two-dimensional Legendre expansion is produced by the tensor-product choice

$$\phi_{kl}(x, y) = L_k(x)L_l(y) , \quad k, l = 0, 1, \dots, N , \quad (1.2.69)$$

where L_k is the Legendre polynomial of degree k . The approximate solution is

$$u^N(x, y) = \sum_{k=0}^N \sum_{l=0}^N a_{kl} L_k(x) L_l(y) . \quad (1.2.70)$$

Note that the trial functions do not satisfy the boundary conditions individually. (In most Galerkin methods the trial functions do satisfy the boundary conditions.) In this case two separate sets of test functions are used to enforce the PDE and the boundary conditions. For the PDE the test functions are

$$\psi_{kl}(x, y) = Q_k(x)Q_l(y) , \quad k = 0, 1, \dots, N-2 , \quad (1.2.71)$$

where

$$Q_k(x) = \frac{2k+1}{2} L_k(x) ; \quad (1.2.72)$$

for the boundary conditions they are

$$\chi_k^i(x) = Q_k(x) , \quad \begin{array}{l} i = 1, 2 , \\ k = 0, 1, \dots, N , \end{array} \quad (1.2.73a)$$

$$(1.2.73b)$$

$$\chi_l^i(y) = Q_l(y) , \quad \begin{array}{l} i = 3, 4 , \\ l = 0, 1, \dots, N . \end{array} \quad (1.2.73c)$$

The integral conditions for the differential equations are

$$\int_{-1}^1 dy \int_{-1}^1 \mathcal{M}(u^N) \psi_{kl}(x, y) dx = 0 , \quad k, l = 0, 1, \dots, N-2 , \quad (1.2.74)$$

while the equations for the boundary conditions are

$$\int_{-1}^1 \mathcal{B}_i(u^N) \chi_k^i(x) dx = 0 , \quad \begin{array}{l} i = 1, 2 , \\ k = 0, 1, \dots, N , \end{array} \quad (1.2.75a)$$

$$\int_{-1}^1 \mathcal{B}_i(u^N) \chi_l^i(y) dy = 0 , \quad \begin{array}{l} i = 3, 4 , \\ l = 0, 1, \dots, N . \end{array} \quad (1.2.75b)$$

Four of the conditions in (1.2.75) are linearly dependent upon the others; in effect the boundary conditions at each of the four corner points have been

applied twice. For the Poisson equation the above integrals may be performed analytically. The result is

$$-(a_{kl}^{(2,0)} + a_{kl}^{(0,2)}) = f_{kl}, \quad k, l = 0, 1, \dots, N-2, \quad (1.2.76)$$

$$\sum_{k=0}^N a_{kl} = 0, \quad \sum_{k=0}^N (-1)^k a_{kl} = 0, \quad l = 0, 1, \dots, N, \quad (1.2.77a)$$

$$\sum_{l=0}^N a_{kl} = 0, \quad \sum_{l=0}^N (-1)^l a_{kl} = 0, \quad k = 0, 1, \dots, N, \quad (1.2.77b)$$

where

$$f_{kl} = \int_{-1}^1 dy \int_{-1}^1 f(x, y) \psi_{kl}(x, y) dx, \quad (1.2.78)$$

$$a_{kl}^{(2,0)} = \left(k + \frac{1}{2}\right) \sum_{\substack{p=k+2 \\ p+k \text{ even}}}^N [p(p+1) - k(k+1)] a_{pl}, \quad (1.2.79a)$$

$$a_{kl}^{(0,2)} = \left(l + \frac{1}{2}\right) \sum_{\substack{q=l+2 \\ q+l \text{ even}}}^N [q(q+1) - l(l+1)] a_{kq}. \quad (1.2.79b)$$

These last two expressions represent the expansions of $\partial^2 u^N / \partial x^2$ and $\partial^2 u^N / \partial y^2$, respectively, in terms of the trial functions.

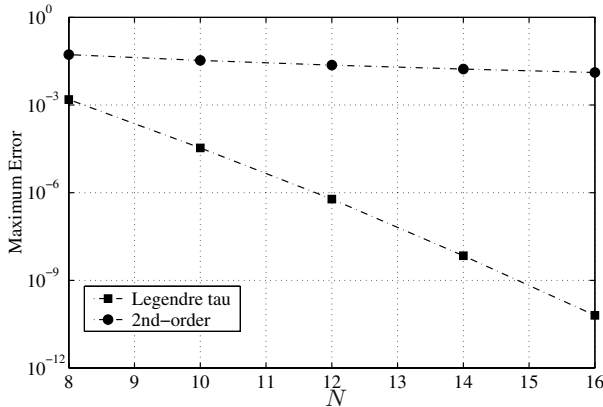


Fig. 1.5. Maximum errors for the Poisson problem for Legendre tau and second-order finite-difference schemes

The Legendre tau approximation to the Poisson equation consists of (1.2.76) and (1.2.77). An efficient scheme for the solution of these equations is provided in Sect. 4.1.

The specific example that will be used to illustrate the accuracy of this method is

$$f(x, y) = 2\pi^2 \sin \pi x \sin \pi y, \quad (1.2.80)$$

which corresponds to the analytic solution

$$u(x, y) = \sin \pi x \sin \pi y. \quad (1.2.81)$$

The results are given in Fig. 1.5 along with results for a second-order finite-difference scheme. The integer N denotes the degree of the expansion (1.2.70) in each dimension for the Legendre tau method and the number of uniform intervals in each dimension for the finite-difference method.

1.2.5 Basic Aspects of Galerkin, Collocation, G-NI and Tau Methods

The Galerkin, collocation, G-NI and tau methods are more general than suggested by any of the above examples. In a broad sense, pure Galerkin and tau methods are implemented in terms of the expansion coefficients, whereas collocation methods and G-NI (Galerkin with numerical integration) methods are implemented in terms of the physical space values of the unknown function. The first example illustrated only one of the key aspects of Galerkin methods – the test functions are the same as the trial functions. The other important aspect is that the trial functions must individually satisfy all or part of the boundary conditions (the remaining ones are enforced weakly within the integral conditions). In the case of periodic boundary conditions the trigonometric polynomials automatically satisfy these requirements. Otherwise, simple linear combinations of the orthogonal polynomials will usually suffice. For example, an obvious choice of trial functions for a Chebyshev Galerkin approximation to the fourth example is

$$\phi_k(x) = \begin{cases} T_0(x) - T_k(x), & k \text{ even } \geq 2, \\ T_1(x) - T_k(x), & k \text{ odd } \geq 3; \end{cases}$$

a computationally more efficient choice (see Sect. 2.3.3) is provided by

$$\phi_k(x) = T_{k-2}(x) - T_k(x), \quad k \geq 2.$$

On the other hand, for the tau method the trial functions do not individually satisfy the boundary conditions. Thus, some equations are needed to ensure that the global expansion satisfies the boundary conditions. Some of

the integral equations corresponding to the highest order test functions are dropped in favor of these boundary condition equations.

The collocation method uses the values of the function at certain physical points as the fundamental representation; the expansion functions are employed solely for evaluating derivatives (and only when a fast transform is available and convenient). The collocation points for both the differential equations and the boundary conditions are usually the same as the physical grid points. The most effective choice for the grid points are those that correspond to quadrature formulas of maximum precision.

The Galerkin with numerical integration (G-NI) method aims at preserving the advantages of both Galerkin and collocation methods. Integrals appearing in the weak formulation of the problem are efficiently approximated by the quadrature formulas mentioned above. Usually, the solution is again represented in physical space through its values at a selected set of nodes. In most cases, as in the example in Sect. 1.2.3, the nodes that serve to represent the solution coincide with the nodes that are used for quadrature. Some exceptions are discussed in later chapters. Certain boundary conditions (for instance, those involving derivatives for second-order operators) are imposed weakly, through a penalty approach that naturally stems from the weak formulation of the problem.

1.3 Three-Dimensional Applications in Fluids: A Look Ahead

Chapters 2–4 of CHQZ3 are devoted to the details of spectral algorithms for investigations of instability, transition and turbulence in fluid flows. The simplest class of flows, termed *laminar flow*, comprises those flows in which the motion is quite regular and predictable, even though possibly unsteady. (Plane Poiseuille flow, discussed in CHQZ3, Sects. 1.3, 2.3 and 3.4, is one example of a laminar flow.) Laminar flows are either stable or unstable. In somewhat oversimplified terms, linearly stable flows are those in which all sufficiently small perturbations to the mean flow decay, whereas unstable flows are those in which some small perturbations grow. Many flows start out as laminar, become unstable (in space or time), and eventually undergo a transition to turbulent flow. The complex category of *turbulent flow* is described by Hinze (1975) as

“Turbulent fluid motion is an irregular condition of flow in which the various quantities show a random variation with time and space coordinates, so that statistically distinct average values can be discerned.”

In this section we illustrate some representative flow physics results from many of the principal fully spectral algorithms that we discuss in Chaps. 2–4.

Turbulent flows contain a wide range of length scales, bounded above by the geometric dimension of the flow field and bounded below by the dissipative action of the molecular viscosity (see, for instance, Tennekes and Lumley (1972, Chap. 3)). The ratio of the macroscopic (largest) *integral length scale* L to the microscopic (smallest) length η (usually known as the *Kolmogorov length scale*) is

$$\frac{L}{\eta} = \text{Re}^{3/4} ,$$

where the Reynolds number Re is

$$\text{Re} = \frac{uL}{\nu} , \quad (1.3.1)$$

with ν denoting the kinematic viscosity and $u = \left(\overline{u'^2}/3\right)^{1/2}$, where u' is the fluctuating velocity, and the bar denotes time averaging. To resolve these scales, N mesh points would be needed in each direction, where

$$N = c_1 \frac{L}{\eta} .$$

(A summary of nondimensionalization in general and Reynolds numbers in particular is provided in CHQZ3, Sect. 1.1.4.)

Two simple classes of turbulent flows are homogeneous turbulence, for which the flow properties are invariant with respect to translations, and isotropic turbulence, for which the flow properties everywhere are invariant with respect to rotations. (Isotropic turbulence is necessarily homogeneous.) For the simulation of homogeneous turbulence with a spectral method, it is appropriate to take $c_1 = 2$; for a fourth-order scheme c_1 would be about 6 and for a second-order scheme about 24. (These estimates are based on the typical requirement of 0.1% or better accuracy per period, using estimates such as those by Kreiss and Oliger (1972), and conclusions from the channel flow computations presented in CHQZ3, Sect. 1.3.) The ratio of the time scales of the macroscopic and microscopic motions is $T/t = \sqrt{\text{Re}}$. Consequently, the number of time-steps required to describe the flow during the characteristic period (or *temporal scale*) of the physically significant events is

$$N_{Ts} = c_0 \sqrt{\text{Re}} , \quad (1.3.2)$$

where the multiplicative factor, c_0 , is between 100 and 1000 depending on the time-stepping algorithm and the time interval needed to obtain reasonable statistics for the flow. Now, the number of operations required to update the solution per time-step of a multistep scheme such as Adams-Bashforth or per stage of a multistage scheme such as Runge-Kutta is

$$c_2 N^3 \log_2 N + c_3 N^3 ,$$

where, for the spectral method, $c_2 = 45$, $c_3 = 35$, for the fourth-order spatial method, $c_2 = 17$, $c_3 = 120$, and for the second-order spatial method, $c_2 = 17$, $c_3 = 60$. (For the finite-difference methods, this assumes that the convection term is treated explicitly, the diffusion term is treated implicitly, a Poisson equation is solved for the pressure, and that the implicit equations for the finite-difference method are solved exactly using FFTs. See CHQZ3, Sect. 3.3 for the details of the spectral algorithm.) Thus, for homogeneous turbulence simulations, the storage requirement is roughly proportional to

$$4c_1^3 \text{Re}^{9/4}, \quad (1.3.3)$$

and the total number of operations is approximately

$$c_0 c_1^3 \text{Re}^{11/4} \left[c_2 \log_2(c_1 \text{Re}^{3/4}) + c_3 \right]. \quad (1.3.4)$$

The estimates above provide the resolution requirements for computations in which all the scales of the flow are resolved numerically. Such a computation is known as a *direct numerical simulation* (DNS). Many of the examples that follow are from DNS computations.

The original Orszag and Patterson (1972) computations were performed in an era in which the fastest supercomputer had a speed of roughly 1 MFlop (10^6 floating point operations per second). Using a typical value of $c_0 = 500$, the computer time required then for one realization of homogeneous turbulence by a spectral method was, according to (1.3.4), about 10 hours for their $\text{Re} = 45$ cases. (Their computations used $N = 32$ modes in each direction.) For sustained performances typical of the fastest supercomputers circa 1980 (100 MFlop), the computer time required for one realization of homogeneous turbulence by a spectral method is 6 minutes for $\text{Re} = 45$ and 2 years for $\text{Re} = 3000$ (for the Brachet et al. (1983) case mentioned below, although they were able to save a factor of 64 by exploiting symmetries). Assuming a sustained performance of 1 TFlop (10^{12} floating point operations per second, typical of the very fastest supercomputers circa 2000), the computer time required for one realization of homogeneous turbulence by a spectral method is about 10 hours for $\text{Re} = 3000$, and about 4 months for $\text{Re} = 40,000$ (for the Kaneda and Ishihara (2006) results mentioned below).

Spectral methods have been singularly successful for this problem since the corresponding requirements for a fourth-order finite-difference method are typically a factor of 10 longer in time and a factor of 20 larger in storage. Second-order finite-difference methods require more than 3 orders of magnitude more resources than spectral methods on this problem. Moreover, Fourier functions arise naturally in the theoretical analysis of homogeneous turbulence, and they are the natural choice of trial functions for spectral methods. Thus, the spectral methods, apart from their computational efficiency, have the added advantage of readily permitting one to monitor and diagnose nonlinear interactions which contribute to resonance effects, energy

transfer, dissipation and other dynamic features. Furthermore, if there are any symmetries underlying a problem, and symmetry-breaking phenomena are precluded, spectral methods permit unique exploitation of these symmetries. (Since the finite-difference methods cannot benefit from the symmetries exploited by Brachet et al. (1983), even the fourth-order method is nearly a thousand times less efficient than the spectral method in this case.) These advantages in computational efficiency are so compelling that they have motivated many flow physics research groups to adopt spectral methods despite their additional complexity rather than simply waiting for increased computational power to make their desired computations feasible. These advantages have also inspired many numerical analysts to develop more efficient spectral methods and to provide their firm theoretical foundation.

Much theoretical work on homogeneous turbulence has focused on the details of the inertial range, which is the range of scales of motion (well observed experimentally) that are not directly affected by the energy maintenance and dissipation mechanisms (Mestayer et al. (1970)) and that possess an energy spectrum exhibiting a scaling behavior (Grant, Stewart, and Moilliet (1962)):

$$E(k, t) = k^{-m}$$

where k is the magnitude of the wavenumber vector and m is close to $5/3$. The spectrum with $m = 5/3$ is the famous Kolmogorov spectrum. The huge Reynolds numbers required to produce an extended inertial range are experimentally accessible only in geophysical flows such as planetary boundary layers and tidal channels.

The pioneering simulations of isotropic turbulence by Orszag and Patterson (1972) evolved over the subsequent decade-and-a-half to the first numerically computed three-dimensional inertial range by Brachet et al. (1983). (See CHQZ3, Sects. 3.3.1 and 3.3.2 for details on this Fourier Galerkin algorithm.) The Reynolds number was 3000 and, of course, crude by experimental standards. This calculation of the Taylor-Green vortex was feasible only because the symmetries of the problem were fully exploitable with the spectral method to obtain an effective resolution of 256^3 , i.e., the equivalent of $N = 256$ modes in each spatial direction. Among the salient results of this study is the physical insight gained into the behavior of turbulence at high Reynolds number, including the formation of an inertial range and the geometry of the regions of high vorticity.

Two decades later Kaneda and Ishihara (2006) (see also Yokokawa et al. (2002)) exploited 512 nodes of the Earth Simulator (then the world's fastest computer) to perform isotropic turbulence simulations using a very similar, Fourier spectral algorithm on grids as large as 4096^3 . (The sustained speed was as fast as 16 TFlop.) Figure 1.6 illustrates the regions of intense vorticity in $1/64$ of the volume of their 2048^3 simulation for $Re = 16,135$. The macroscopic scale L is approximately 80% the size of one edge of the figure, and the microscopic scale η is 0.06% of the edge length. Among the

many results obtained from their high-resolution simulations was convincing evidence that the scaled energy spectrum (where the wavenumber is scaled by the inverse of the Kolmogorov length scale $\eta = (\nu^3/\bar{\epsilon})^{1/4}$, with ν the viscosity and $\bar{\epsilon}$ the average dissipation rate) is not the classical Kolmogorov result of $k^{-5/3}$, but rather k^{-m} with $m \simeq 5/3 - 0.10$.

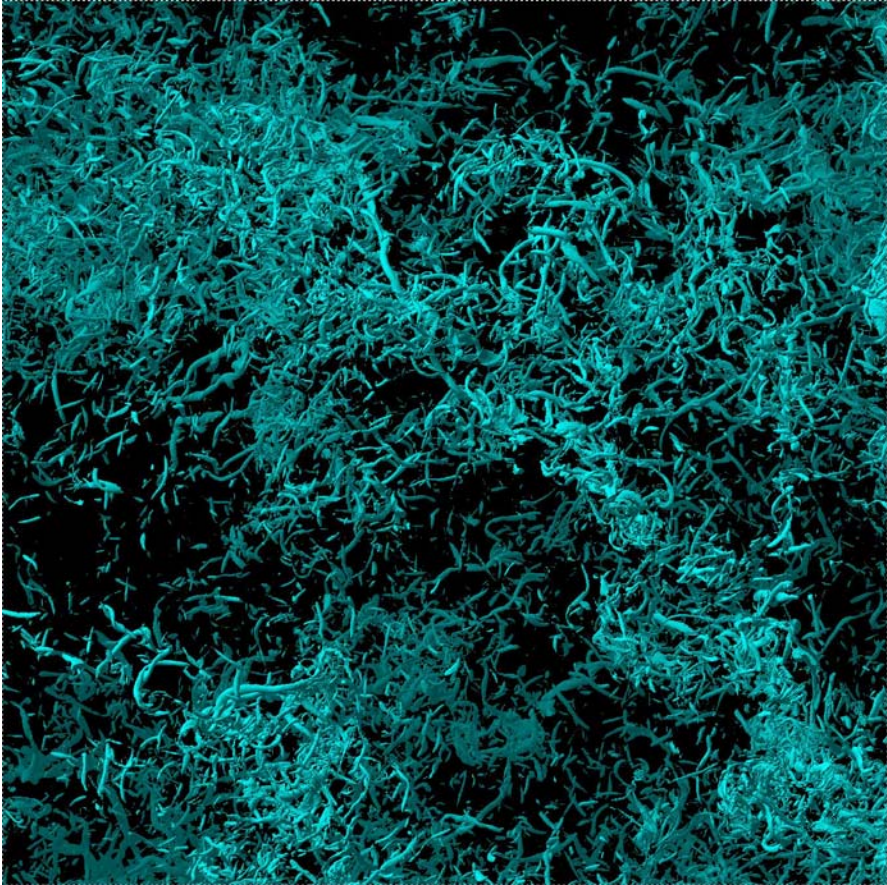


Fig. 1.6. Direct numerical simulation of incompressible isotropic turbulence by Kaneda and Ishihara (2006) on a 2048^3 grid. The figure shows the regions of intense vorticity in a subdomain with $1/4$ the length in each coordinate direction of the full domain [Reprinted with kind permission by the authors]

Rogallo (1977) developed a transformation that permits Fourier spectral methods to be used for homogeneous turbulence flows, such as flows with uniform shear. Blaisdell, Mansour and Reynolds (1993) used the extension of this transformation to the compressible case to simulate compressible, homoge-

neous turbulence in uniform shear on 192^3 grids ($N = 192$ grid points in each spatial direction) using a Fourier collocation method. (In this example, as in all the examples cited in this section for inhomogeneous flows, the y direction is the direction of inhomogeneity.) Figure 1.7 illustrates the coalescence of sound waves that is responsible for enhanced turbulence production in compressible flows. The Rogallo transformation is described in CHQZ3, Sect. 3.3.3 for incompressible flow and in CHQZ3, Sect. 4.3 for compressible flow.

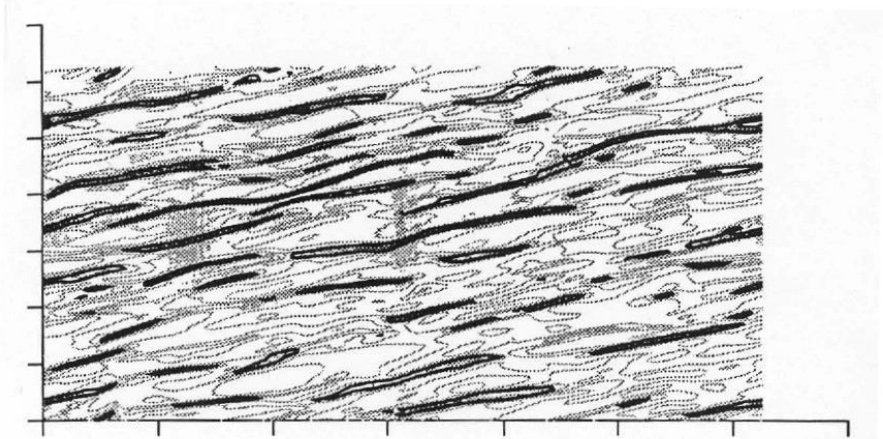


Fig. 1.7. Two-dimensional slice illustrating contours of the pressure field from a compressible homogeneous turbulence DNS by Blaisdell and Zeman (1992) [Reprinted with permission from G.A. Blaisdell, O. Zeman (1992); Center for Turbulence Research, Stanford University/NASA Ames Research Center]

The applications cited above were all for problems with no physical boundaries. Spectral algorithms for problems with solid boundaries are more subtle, largely because a pure Fourier method is no longer appropriate. It was not until the late 1970's that reliable Fourier-Chebyshev algorithms were applied to the simplest wall-bounded flows (Orszag and Kells (1980), Kleiser and Schumann (1980)). The principal advantage of such spectral methods over finite-difference methods is their minimal phase errors (Sect. 1.2.1). This is especially important in numerical simulations of instability and transition to turbulence, because such simulations must follow the evolution and nonlinear interaction of waves through several characteristic periods. Since phase errors are cumulative, a method that admits phase errors of even a few percent per period is unacceptable.

Kleiser and Schumann (1984) devised an influential algorithm for plane channel flow using two Fourier directions and one Chebyshev direction. This algorithm was later used by Gilbert and Kleiser (1990) for the first simulation of the complete transition to turbulence process in a wall-bounded flow using a 128^3 grid. Figure 1.8 illustrates the evolution of one of the principal

diagnostics of a transitional flow – the wall-normal shear of the streamwise velocity $\partial u/\partial y$. The ordinate in the top part of the figure is the Reynolds number based on the wall shear velocity; it is given by $Re_\tau = \sqrt{\frac{1}{\nu} \frac{\partial \bar{u}}{\partial y}} h$, where h is the channel half-width and $\bar{u}(y, t)$ is the average over x and z of the streamwise velocity. The bottom part of the figure illustrates the evolution of the vertical shear at the spanwise station containing the peak shear. These detailed results compared very favorably with the vibrating ribbon experiments of Nishioka, Asai and Iida (1980). The $t = 136$ frame was already computed by Kleiser and Schumann (1984) at lower resolution. (The Kleiser-Schumann algorithm is given in detail in CHQZ3, Sect. 3.4.1.)

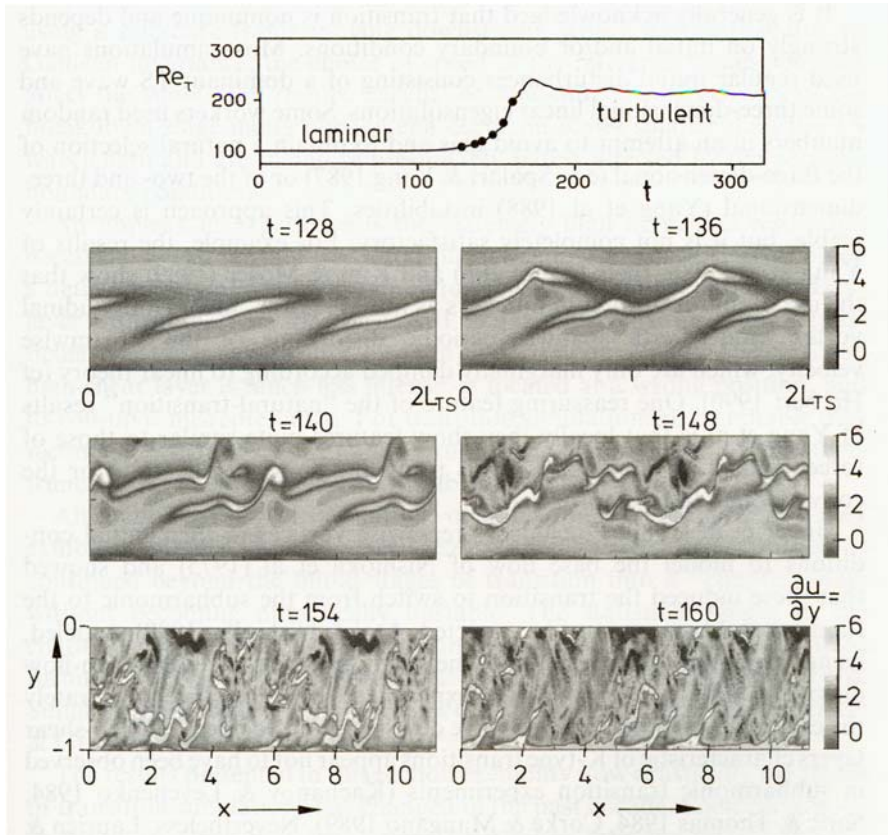


Fig. 1.8. DNS of transition to turbulence in plane channel flow by Gilbert and Kleiser (1990). The top figure illustrates the evolution in time of the Reynolds number based on wall friction velocity. The remaining frames illustrate the shear, $\partial u/\partial y$, in the bottom half of the channel in a two-dimensional slice at the spanwise (z) location containing the maximum shear [Reprinted with permission from N. Gilbert, L. Kleiser (1990); © 1990, Taylor and Francis Group]

Another widely-used algorithm, this one based on the vorticity-velocity equations, was originally developed by Kim, Moin and Moser (1987) for plane channel flow (see CHQZ3, Sect. 3.4.1). Figure 1.9 shows results from Rogers and Moser (1992) using the adaptation of this algorithm to incompressible, free shear layers; Fourier series are employed in the two homogeneous directions (x and z) and Jacobi polynomials (see Sect. 2.5) in the y direction. This figure, based on computations on a $64 \times 128 \times 64$ grid, illustrates several aspects of the vorticity from a simulation that is most representative of experiments on vortex roll-up in mixing layers. The thin, shaded surfaces correspond to the rib vortices (large component of vorticity normal to the spanwise direction), the cross-hatched surfaces denote the “cups” (regions of strong spanwise vorticity) that are critical to free shear layer transition, and the lines are vortex lines that comprise the rib vortices.

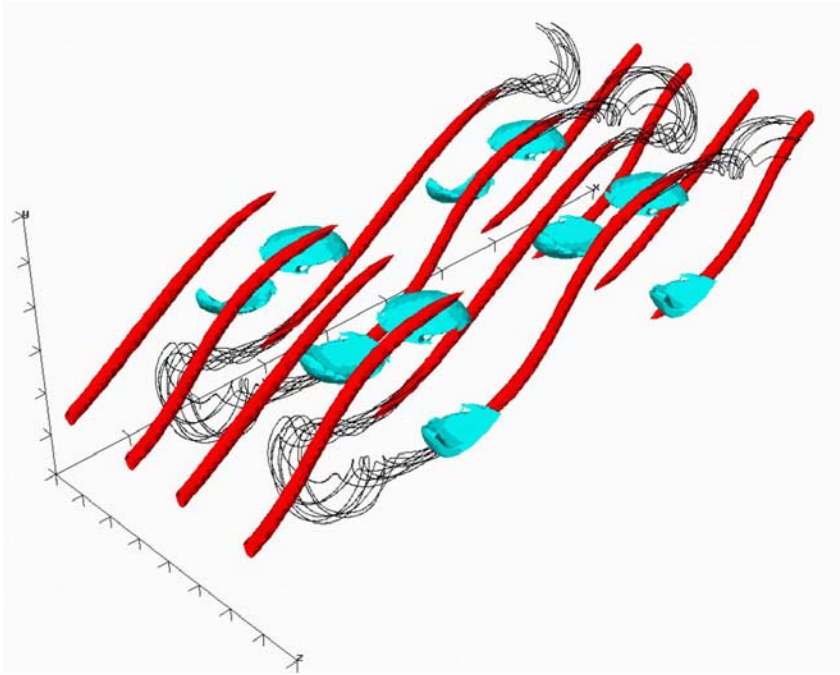


Fig. 1.9. DNS of vortex rollup in an incompressible free shear layer by Rogers and Moser (1992). The surfaces denote two types of regions of strong vorticity and the lines are vortex lines [Reprinted with permission from M.M. Rogers, R.D. Moser (1992); © 1992, Cambridge University Press]

Orszag and Kells (1980) and Orszag and Patera (1983) pioneered the use of splitting methods for wall-bounded flows. Figure 1.10 illustrates results from a later version of a splitting method, due to Zang and Hussaini (1986),

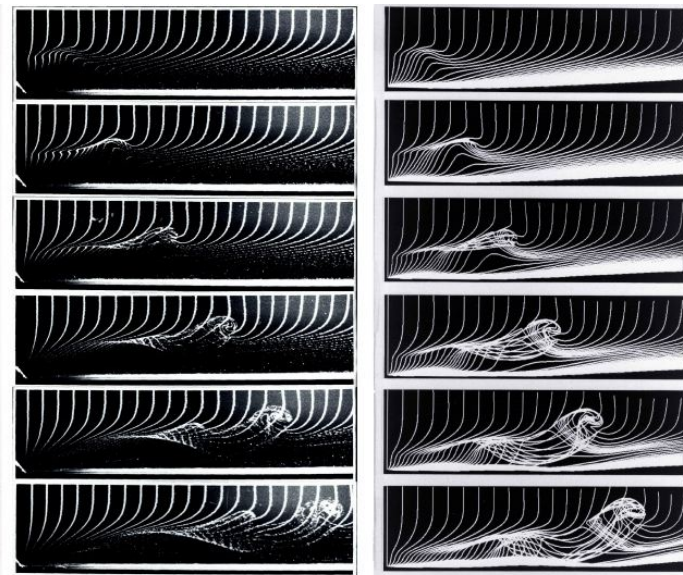


Fig. 1.10. Comparison of hydrogen bubble flow visualizations (*left*) of incompressible flat plate boundary-layer transition with DNS results of Zang, Hussaini and Erlebacher (*right*) [Reprinted with permission from T.A. Zang, M.Y. Hussaini (1987); © 1987 ASME]

applied to transition in a simplified version of flow past a flat plate. (The simplification invokes the parallel flow approximation that is discussed in CHQZ3, Sects. 2.3.2 and 3.4.5.) The left half of the figure is taken from the experiments of Hama and Nutant (1963) who used a hydrogen bubble flow visualization technique to illustrate the strongly nonlinear stage of transition. The right half of the figure, from Zang, Hussaini and Erlebacher (see Zang, Krist, Erlebacher and Hussaini (1987) and Zang and Hussaini (1987)), shows how well this phenomena was reproduced in the numerical computations using a $128 \times 144 \times 288$ grid. These authors demonstrated that the fine details of the vortex roll-ups were not present in the streamwise symmetry plane but only appeared in a streamwise plane displaced by a small fraction of the spanwise wavelength from the symmetric plane. (Details of the splitting algorithms are provided in CHQZ3, Sect. 3.4.2.)

This same splitting algorithm – the Zang-Hussaini version – was used by Scotti and Piomelli (2001) in their 64^3 large-eddy simulations of pulsating channel flow. *Large-eddy simulation* (LES) is one method of accounting for the effects of turbulence by solving an augmented set of equations on a grid much coarser than for a DNS. (See CHQZ3, Sect. 1.1.3 for a summary of LES and Sagaut (2005) for a thorough discussion of the subject.) Figure 1.11 illustrates

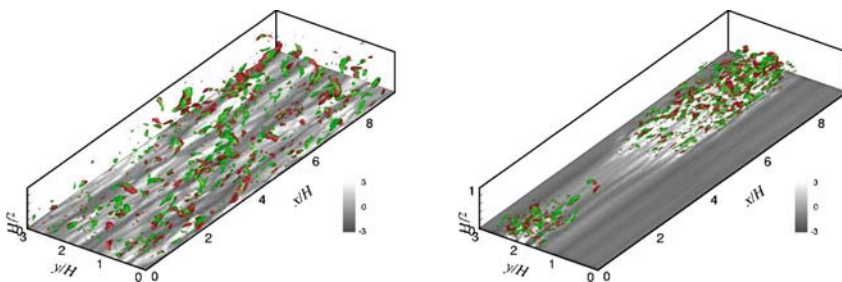


Fig. 1.11. Turbulent fluctuations near the bottom wall in incompressible pulsating channel flow from the LES computations of Scotti and Piomelli (2001). The left frame is near the end of the acceleration phase and the right frame is at the middle of the deceleration phase of the cycle [Reprinted with permission from A. Scotti, U. Piomelli (2001); © 2001, American Institute of Physics]

the flow structures at a fully turbulent phase of the oscillation (left half of the figure) and at a relaminarization phase (right half). The solid surface is a contour of the fluctuating streamwise velocity. The small-scale surfaces are contours of a measure of the coherent vorticity due to rotational motions. Note that the grid used for this large-eddy simulation was significantly coarser than that used in many of the examples above for transitional and turbulent flows. This illustrates a major attraction of the LES approach. The smaller grid permits wide parameter studies to be performed as opposed to the one-of-a-kind simulations typical of direct numerical simulations for such flows. Scotti and Piomelli did parametric studies using LES to characterize the detailed physics of such pulsating flows.

Figure 1.12 illustrates results from three additional classes of spectral algorithms. The physical problem is the study of the instability of flow past a flat plate. Unlike the computation of Zang, Hussaini and Erlebacher, shown above in Fig. 1.10, where the parallel flow approximation was used to study the temporal instability of this important physical problem, the results in Fig. 1.12 were for the unadulterated, spatial instability of the nonparallel flow past a flat plate. This problem requires the resolution of 10's or 100's of wavelengths in the streamwise direction (and has challenging outflow boundary conditions) rather than the mere 1 or 2 wavelengths in x that are needed in the parallel flow approximation. The direct numerical simulation results used Spalart's (1988) ingenious *fringe method*, which permits a highly accurate approximation to be obtained with a Fourier approximation in x . (See CHQZ3, Sect. 3.6.1 for the details.) These two-dimensional DNS computations required approximately 4 points per wavelength in x and no more than 40 Jacobi polynomials in y . The *parabolized stability equations* (PSE) *method* solve a much more economical set of equations using a marching method in x , a low-order Fourier expansion in z and a Chebyshev collocation method in y with $N \leq 40$. (See CHQZ3, Sects. 2.4.1 and 2.5.2 for PSE algorithms.)

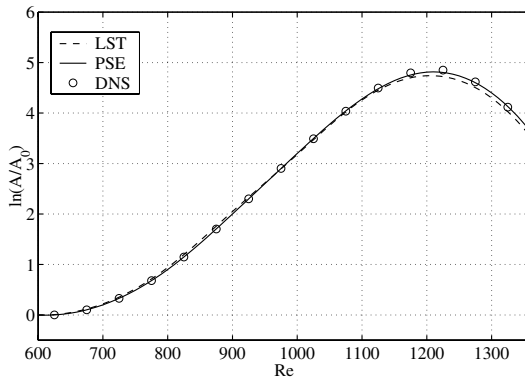


Fig. 1.12. Evolution of the spatial instability of an incompressible flat-plate boundary layer by Bertolotti, Herbert and Spalart (1992). Results are shown for direct numerical simulation (DNS), parabolized stability equations (PSE) and linear stability theory (LST) using the parallel flow approximation [Adapted with permission from F.P. Bertolotti, Th. Herbert, P.R. Spalart (1992); © 1992, Cambridge University Press]

The figure compares the spatial development of the maximum streamwise velocity perturbation as computed by the DNS and by the PSE; also shown for comparison are results of linear stability theory (LST) using the parallel flow approximation. (Spectral algorithms for linear stability are discussed in CHQZ3, Sect. 2.3.) The results of the PSE method agree well with the DNS results and are far cheaper. Hence, the PSE is far better suited to parametric studies.

Simulations of much later stages of transition in spatially developing flows have also been performed with both PSE and DNS techniques utilizing spectral methods. The spatial simulation of oblique transition in a boundary layer on a $1200 \times 64 \times 96$ grid by Berlin, Wiegel and Henningson (1999) is a prime example of a high-resolution DNS using the fringe method with a Fourier-Chebyshev algorithm. Figure 1.13 illustrates a comparison of their numerical results with flow visualizations of their experiment on transition in a boundary layer. (The algorithm uses components discussed in CHQZ3, Sects. 3.4.1, 3.4.4 and 3.6.1.)

In addition to the DNS, LES and PSE computations emphasized in the examples so far, spectral methods have also excelled in computations of eigenvalue problems. Indeed, Orszag's (1971b) demonstration of the power of Chebyshev spectral methods for discretizing the eigenvalue problems arising in linear stability analyses inspired many subsequent workers to adopt spectral methods for such problems in both incompressible and compressible flows. Eventually, in the 1990's computer resources were adequate for solving such problems with two or even three directions treated as inhomogeneous. An example of a large-scale eigenvalue problem solved by Theofilis (2000),

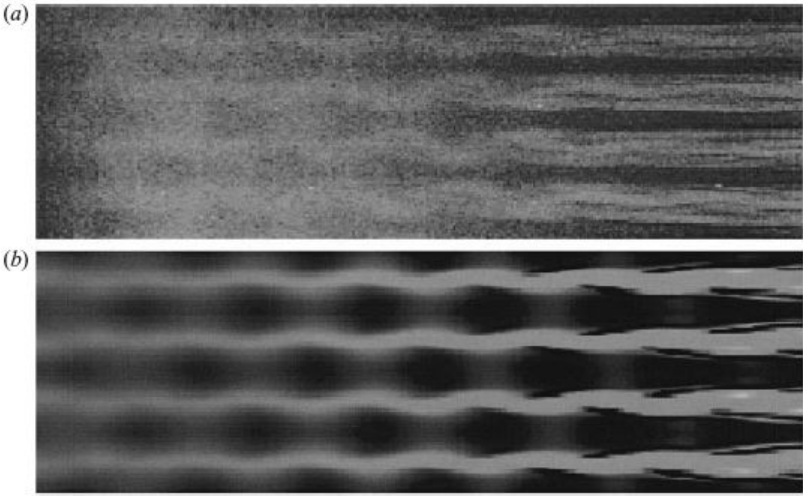


Fig. 1.13. Streamwise velocity flow visualizations of incompressible boundary-layer transition by Berlin, Wiegel and Henningson (1999): experiment (a) and spatial computation (b) [Reprinted with permission from S. Berlin, M. Wiegel, D.S. Henningson (1999); © 1999, Cambridge University Press]

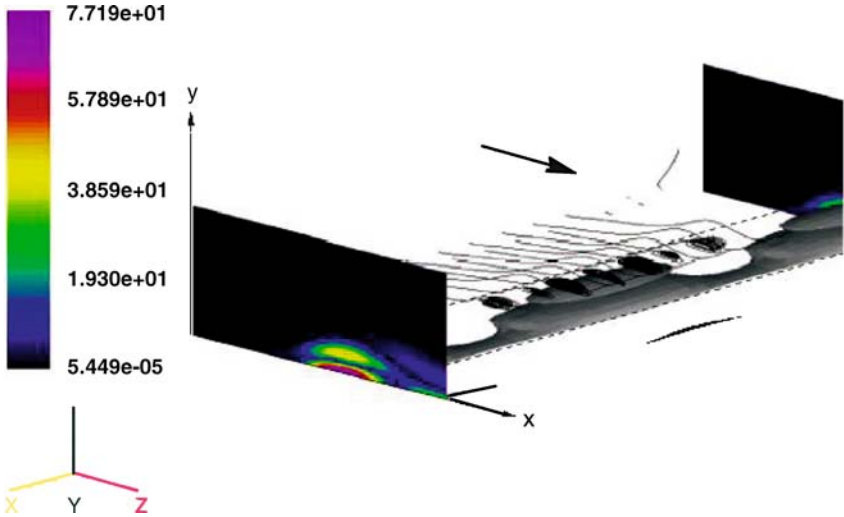


Fig. 1.14. Isosurface of disturbance vorticity of the primary instability of an incompressible separation bubble by Theofilis (2000)

who used two Chebyshev directions and one Fourier direction, is given in Fig. 1.14. Spectral algorithms for discretizing the eigenvalue problems of fluid dynamical linear stability are described in much of CHQZ3, Chap. 2.

This list is by no means exhaustive and certainly neglects applications in related disciplines such as meteorology, oceanography, plasma physics and general relativity. Many of the components of algorithms mentioned above have been analyzed theoretically. The essential elements of the numerical analysis are provided in Chap. 7. Rigorous error estimates for some incompressible Navier-Stokes algorithms are reviewed in CHQZ3, Chap. 3.

The examples in this section have been confined to those using classical spectral methods. We noted earlier in this section that fourth-order methods require a factor of 10 more computational resources than spectral methods. The desire to handle problems in complex domains with greater than fourth-order accuracy has motivated the development of higher order methods using domain decomposition. Chapters 5 and 6 of the companion book (CHQZ3) survey spectral methods in complex domains. Chapters 2–7 of this book and Chaps. 1–4 of CHQZ3 are devoted to classical spectral methods.

Spectral Methods

Fundamentals in Single Domains

Canuto, C.; Hussaini, M.Y.; Quarteroni, A.; Zang, Th.A.

2006, XXII, 581 p. 106 illus., 10 illus. in color., Hardcover

ISBN: 978-3-540-30725-9