

A USER'S VIEW OF SOLVING STIFF ORDINARY DIFFERENTIAL EQUATIONS*

L. F. SHAMPINE[†] AND C. W. GEAR[‡]

Abstract. This paper aims to assist the person who needs to solve stiff ordinary differential equations.

First we identify the problem area and the basic difficulty by responding to some fundamental questions: Why is it worthwhile to distinguish a special class of problems termed "stiff"? What are stiff problems? Where do they arise? How can we recognize them?

Second we describe the characteristics shared by methods for the numerical solution of stiff problems. These characteristics have important implications as to the convenience and efficiency of solution of even routine problems. Understanding them is indispensable to the assembling of codes for the very efficient solution of special problems or for solving exceptionally large problems at all.

Third we shall briefly discuss what is meant by "solving" a differential equation numerically and what might be reasonably expected in the case of stiff problems.

1. Introduction. The numerical solution of ordinary differential equations is an old topic and, perhaps surprisingly, methods discovered around the turn of the century are still the basis of the most effective, widely used codes for this purpose [23]. Great improvements in efficiency have been made, but it is probably fair to say that the most significant achievements have been in reliability, convenience, and diagnostic capabilities. The typical scientific problem can be solved by casual users of these codes both easily and cheaply. Nevertheless, there are several kinds of problems which classical methods do not handle very efficiently. The problems called "stiff" are too important to ignore, and are too expensive to overpower. They are too important to ignore because they occur in many physically important situations. They are too expensive to overpower because of their size and the inherent difficulty they present to classical methods, no matter how great an improvement in computer capacity becomes available. Even if one can bear the expense, classical methods of solution require so many steps that roundoff errors may invalidate the solution. It is all the more frustrating that the solutions of stiff problems look like they should be particularly easy to compute. After a few general remarks about solving differential equations, we shall use some simple examples to show where the trouble originates and what might be done about it. We shall mention a number of contexts in which stiffness was recognized and dealt with by scientists through the use of special features of the problem. Our attention here will be directed towards the phenomenon of stiffness and towards general purpose procedures for the solution of stiff differential equations. There are effective codes available based on these procedures, but it is necessary that the user have some idea how they work in order to take full advantage of them. Lastly we discuss what are realistic goals when solving a stiff differential equation.

2. What are stiff problems? When solving the (vector) system of equations

$$(2.1) \quad y' = f(x, y), \quad y(0) \text{ given},$$

we must consider the behavior of solutions near to the one we seek. This is because as we step along from $y_n \doteq y(x_n)$ to y_{n+1} approximating $y(x_n + h)$ we make inevitable errors causing us to move from the desired integral curve to a nearby one. If we make no further errors, we follow this new curve so that the resulting error depends on the

* Received by the editors October 14, 1976, and in revised form May 12, 1977.

[†] Numerical Mathematics Division 5122, Sandia Laboratories, Albuquerque, New Mexico 87115.

[‡] Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801.

relative behavior of the two solution curves. Let us consider the example of the single equation

$$(2.2) \quad y' = A(y - p(x)) + p'(x), \quad y(0) = v,$$

where A is a constant. The analytical solution is

$$(2.3) \quad y(x) = (v - p(0)) \exp(Ax) + p(x).$$

If A is large and positive, the solution curves for the various v fan out and we say the problem is *unstable*. Such a problem obviously is difficult for any general numerical method which proceeds in a step-by-step fashion. When A is small in magnitude, the curves are more or less parallel and such *neutrally stable* problems are easily handled by conventional means. When A is large and negative, the solution curves converge very quickly. In fact, whatever the value $y(0)$, the solution curve is virtually identical to the particular solution $p(x)$ after a short distance called an *initial transient*. This super-stable situation is ideal for the propagation of error in the differential equation but not, as it turns out, for the propagation of error in a numerical scheme. The last class of problems is called *stiff*.

Equation (2.2) is of more general significance than its special form suggests. The behavior of solutions of (2.1) near a particular solution $g(x)$ can be studied by a Taylor series expansion into

$$(2.4) \quad \begin{aligned} y' &\doteq J(x, g(x))(y - g(x)) + f(x, g(x)) \\ &= J(x, g(x))(y - g(x)) + g'(x) \end{aligned}$$

where the *Jacobian* matrix J has as its (i, j) entry the partial derivative of the i th component of f with respect to the j th component of y . This approximation can be justified in a limiting argument, but we are only going to use it qualitatively and will not attempt rigor. We further suppose that $J(x, g(x))$ is slowly varying in x so that we can approximate it locally by a constant matrix. After a principal axis transformation, these equations are uncoupled into a set of equations each of the form (2.2). In this general situation, A is an eigenvalue of the Jacobian; hence, it may be a complex number. Thus, within the limits of the approximations made, a set of simultaneous equations of the form (2.2) with complex numbers A are representative of the general situation (2.1).

By a stiff problem we mean one for which no solution component is unstable (no eigenvalue has a real part which is at all large and positive) and at least some component is very stable (at least one eigenvalue has a real part which is large and negative). Further, we will not call a problem stiff unless its solution is slowly varying with respect to the most negative real part of the eigenvalues. (Roughly, we mean that the derivatives of the solution are small compared to the corresponding derivatives of e^{Ax} . The meaning will be qualified later in this section.) Consequently, a problem may be stiff for some intervals of the independent variable and not for others.

If A is very negative and $p(x)$ is slowly varying, equation (2.3) represents a stiff problem after the transient e^{Ax} has died out (that is, e^{Ax} is below the error tolerance of interest) but it is not stiff in the transient region. If (2.1) is linear with a constant Jacobian J , it will not be stiff in the initial transient, but will be stiff after the fastest transient has died out.

Although the examples just discussed only exhibit one period of rapid change, the general problem can exhibit several. One reason is that the approximations in equation (2.4) only apply locally. This may, or may not, be obvious from the equations

themselves. A case which is not obvious is the relaxation oscillation of the Van der Pol equation famous in applied mechanics. One should suspect stiffness because the limit curve is rapidly approached by all integral curves, and examination of the Jacobian matrix does show that the integration has large negative eigenvalues in some regions, although not in others. The limit solution has periodic changes which are so sharp as to approach a discontinuity. (The equation is *not* stiff at such places by our definition.) The methods of classical applied mathematics, namely, singular perturbation theory, deal with this particular equation quite well.

A second reason is that the "driving" term— $p(x)$ in equation (2.2)—may suddenly change. An example of this with transients expected on physical grounds is that of chemical kinetic rate equations involving photo-dissociation such as those describing the behavior of atmospheric pollutants. When the sun rises or sets there are reactions which are extremely fast with respect to a basic period of a day. A. Hindmarsh has communicated to us a mockup of the behavior of the oxygen singlet—D(01D) which is a valuable test problem exhibiting the proper physical behavior. It is

$$y'(t) = d - by + aE(t)$$

where

$$a = 10^{-18}, \quad b = 10^8, \quad c = 4, \quad d = 10^{-19}$$

$$E(t) = \left(1 + \frac{c\omega^2}{b} \frac{\cos \omega t}{\sin^2 \omega t}\right) e(t),$$

$$e(t) = \begin{cases} \exp(-c\omega/\sin \omega t) & \text{if } \sin \omega t > 0, \\ 0 & \text{otherwise} \end{cases}$$

and

$$\omega = \pi/43200.$$

Note that the Jacobian matrix is just the number $-b$ which in this case says that the equation is very stable. The time t is in seconds and every 12 hours the solution (which can be obtained analytically) exhibits a change which is almost discontinuous on a time scale of days.

We should also note that the approximations developed for (2.1) lead to only one limit solution, but that in practice there may be others. This originates in the fact that the approximation is local about the solution being studied. The area of nonlinear mechanics is a fertile source of examples with several possible limit curves. Later we shall refer to one in the area of chemical engineering.

To expose the difficulty in solving stiff systems, let us integrate (2.2) by Euler's method when A is a large negative number and $p(x)$ is slowly varying. This scheme advances from y_n to an approximation at $x_n + h_n = x_{n+1}$ by $y_{n+1} = y_n + h_n f(x_n, y_n) = y_n + h_n y'_n$. Since these are the linear terms of a Taylor series expansion at x_n , the local truncation error of the method is $h^2 y''(x_n)/2 + O(h^3)$. A code which selects its step size so that the local truncation error is approximately ε (as most codes do) will choose h so that $\varepsilon \doteq |h^2 y''(x_n)/2|$. If the numerical solution is close to the true solution, we can deduce the behavior of h from

$$y''(x) = (v - p(0))A^2 \exp(Ax) + p''(x).$$

When x is small, it is clear that

$$h_n \doteq \left(\frac{2\varepsilon}{|y''(x_n)|} \right)^{1/2} \doteq \left(\frac{2\varepsilon}{|(v-p(0))A^2|} \right)^{1/2}.$$

When x is large the exponential term disappears so that

$$h_n \doteq \left(\frac{2\varepsilon}{|p''(x)|} \right)^{1/2}.$$

By assumption $|A|$ is large and $|p''(x)|$ is small, so that we have quantified the statement that the step size needed for accuracy must initially be small to resolve the rapid change of the transient but eventually becomes large and independent of A .

This is not the whole story. There are two main factors affecting the size of the step—accuracy and stability. Accuracy refers to smallness of the local error, that is, the error introduced in a single step. Stability refers to errors not growing in subsequent steps. We have seen for this example that accuracy is easily handled. Let us now examine stability. Because this is a linear differential equation and a linear method, it is easy to solve for the global error which is the difference between the numerical solution at any point and the true solution. If we define the global error as

$$\delta_n = y_n - y(x_n),$$

we find that

$$\delta_{n+1} = (1 + h_n A) \delta_n + [y(x_n) + h_n y'(x_n) - y(x_{n+1})].$$

This says that the global error after the n th step consists of the error propagated from the previous point x_n plus the local truncation error in the n th step. This error is amplified unless $-2 \leq h_n A \leq 0$. Clearly this restriction dominates in the selection of the step size once outside of the transient region. Note that a problem is not stiff in the transient region because $|h_n A|$ must be small to control the local truncation error.

The essence of the matter is that for most problems the accuracy requirement dictates the choice of step size, but for some, the stiff problems, the stability requirement does. In general we must discuss the stability of the difference scheme when A is a complex number. Analyzing stability as we did with Euler's method, we now find a region in the left half complex plane, called the region of absolute stability, in which $h_n A$ must lie for the difference scheme to be stable. For Euler's method this region is the disc of radius 1 centered at $(-1, 0)$. Though the approximations are crude, this analysis does furnish a good qualitative understanding of the local behavior of the difference schemes. The stability restriction takes the form that $|h_n A|$ not be too large. As a practical matter this is no restriction unless $|A|$ is "large" and the accuracy requirement is easy to meet.

One worry should be dispelled at once. When implemented properly, the instability on encountering stiffness of classical methods such as Euler's is automatically detected and handled by reducing the step size [24], [25]. Computer programs suitable for nonstiff problems do not "blow up" in the presence of stiffness, they just become inefficient. The reason for this is easy to see for methods like Euler's. Automatic codes estimate the local error by estimating a derivative of the solution. The only way in which this can be done economically involves applying some form of difference operator to the computed solution. With the Euler method, for example, we could form the second difference of the solution to estimate the second derivative. The second difference will consist of two parts, the second difference of the true

solution and the second difference of the global error. A simple calculation for our example (2.2) with constant step size h shows that the latter is $(hA)^2 \delta_{n-1} - h^3(Ay''(x_{n-1}) + y'''(x_{n-1}))$ plus higher order terms. If hA is large, this term dominates, so the step control mechanism will reduce the step size accordingly.

The problem of instability for large step sizes when solving stiff equations is common to all methods that are efficient for nonstiff equations. All methods give rise to a global error equation of the form

$$\delta_{n+1} = S_n \delta_n + \varepsilon_n,$$

where ε_n is the local truncation error and S_n is the error amplification matrix whose size depends on the Jacobian of the differential equations and, in the case of multistep methods, δ_n is a vector containing the global errors in the numerical values of all past values used in the computing of the next step. Methods for nonstiff problems are chosen so as to make the local truncation error term ε_n small for best efficiency. When stiff problems are to be solved, it is necessary to sacrifice some of the accuracy in order to improve the stability. Methods suitable for stiff equations are such that S_n is small for Jacobians with large negative eigenvalues. The example of the backward Euler method

$$y_{n+1} = y_n + hy'_{n+1} = y_n + hf(x_{n+1}, y_{n+1})$$

is informative because it is easy to analyze and yet typical of many methods for stiff equations. When it is applied to equation (2.2), we get a global error equation of the form

$$\delta_{n+1} = \delta_n + h_n A \delta_{n+1} + [y(x_n) - y(x_{n+1}) + h_n y'(x_{n+1})]$$

or

$$\delta_{n+1} = (1 - h_n A)^{-1} \delta_n + (1 - h_n A)^{-1} h^2 y''(x_n)/2 + O(h^3).$$

The propagated error is damped whenever $|1/(1 - h_n A)| \leq 1$ —which includes the whole left half plane. This is a marvelous improvement since an apparently minor change in the scheme has completely done away with a stability limitation. However, the backward Euler scheme is implicit, meaning that a nonlinear equation must be solved at each time step to determine y_{n+1} . This is characteristic of methods with very good stability properties and we shall examine the cost later. For now we just comment that because the stability requirement is so much more stringent for the forward Euler method than the accuracy requirement, the backward Euler method permits orders of magnitude improvement in efficiency for typical stiff problems even though each step is much more expensive.

The essence of stiffness is that one has a slowly varying solution which is such that some perturbations to it are rapidly damped. Most physical systems of interest are going to be stable; those which permit very rapid change are the ones which are potentially stiff. Thus one should be alert to physical components with greatly different time constants. For example, control systems are intended to provide stability. When they very quickly correct a deviation from a desired slowly changing path, the differential equations describing them will be stiff.

It is important to appreciate that stiffness depends on the behavior of all nearby solutions, that is, on the differential equation rather than the behavior of the solution itself. For example, the equation

$$y' = (v - p(0))A \exp(Ax) + p'(x), \quad y(0) \text{ given},$$

has exactly the same solution as (2.2) when both have the initial value v , but is not stiff at all. For this quadrature problem the integral curves are parallel and the problem is of neutral stability. Thus the presence of some components which change at a rate much faster than others is not necessarily an indication of stiffness. Still, physical systems are usually stable so this is a pretty reliable indication of stiffness in the proper context. As examples, chemical reactions with large rate constants and nuclear reactions with species decaying at rates varying widely typically lead to stiff equations. Electrical circuitry involving fast elements such as high speed transistor models ordinarily leads to stiff problems. A survey of such applications with examples can be found in [3]. Rather than take up many examples, we shall look in some detail at the use of semi-discretization to solve a simple partial differential equation. The resulting stiff system of equations will illustrate a number of points we shall need later.

Suppose we want to solve the heat equation

$$\frac{\partial y(x, t)}{\partial t} = \frac{\partial^2 y(x, t)}{\partial x^2},$$

$$y(0, t) = 0, \quad y(1, t) = 0, \quad y(x, 0) \text{ given.}$$

Let $\Delta x = 1/N$ and $x_i = i \Delta x$ for $i = 0, 1, \dots, N$. Suppose that $y_i(t)$ is to approximate $y(x_i, t)$ where $y_i(t)$ arises from replacing the space derivative in the heat equation by a centered difference:

$$\frac{dy_i(t)}{dt} = \frac{1}{(\Delta x)^2} (y_{i+1}(t) - 2y_i(t) + y_{i-1}(t)), \quad i = 1, \dots, N-1,$$

$$y_0(t) = 0, \quad y_N(t) = 0, \quad y_i(0) = y(x_i, 0) \quad \text{for each } i.$$

One easily finds the constant Jacobian and that its eigenvalues are

$$A_{N-n} = - \left[\frac{\sin(n\pi \Delta x/2)}{\Delta x/2} \right]^2,$$

from which we see that

$$A_1 = \frac{-4}{(\Delta x)^2}, \quad A_{N-1} = -\pi^2.$$

If Euler's method is used to solve the set of ordinary differential equations with step size Δt , we conclude that for stability $|h A_1| \leq 2$ or

$$\frac{\Delta t}{(\Delta x)^2} \leq \frac{1}{2}.$$

Here how stiff the problem is depends on how fine a spatial mesh we choose initially. If N is small, one is well advised to use nonstiff methods, but for large N this is not economic. Note that if we proceed in this way, we are using the classical (fully discrete) forward difference scheme for solving the heat equation. The backward Euler method corresponds to the classical backward difference scheme and, as is well known, there is then no limitation on Δt for reasons of stability. Considerable experience on the part of people interested solely in partial differential equations shows that it is much cheaper to use the more stable method with the more expensive steps. Experience applying general purpose codes for ordinary differential equations agrees with this and

adds the observation that enhanced efficiency and reliability can be obtained by variation of step size and formula.

It should be noted that the eigenvalues obtained in this example are not due solely to the spatial discretization used. The original partial differential equation has eigenvalues of $-(k\pi)^2$, $k = 1, 2, \dots$, so the first eigenvalue of the discretized system is approximately the first eigenvalue of the differential operator, and the others are approximations to some of the larger ones. This points out that the stiffness is inherent in the problem, not part of the method of solution. Either the problem (that is, the model of the physical situation) must be changed, or stiffness must be faced in the solution process.

In a number of areas, particularly chemical kinetics, problem solvers have removed stiffness by changing the model. The idea is that physical considerations allow some components to be recognized as changing on time scales much shorter than those of other components. Think, for example, of a chemical reaction taking place in a moving medium or the rolling of a rocket as compared to its motion along its trajectory. Quasi-static or pseudo steady-state approximations hold one set of components fixed in value over suitable time periods either because they change so slowly that changes can be neglected or because they change so rapidly that steady-state values are achieved almost instantly. Such approximations lead to sets of algebraic equations coupled to (hopefully nonstiff) sets of ordinary differential equations. In some cases approaches of this kind have worked very well, but it is hard to relate the solution of the modified model to that of the original model. We shall reconsider this technique in the next section where it will be seen that the methods for stiff problems do much the same thing automatically. Current codes for stiff differential equations are sufficiently efficient that there is no need to consider such model changes for most problems for reasons of cost, and there are excellent reasons of convenience and theoretical support for not changing the model. To be sure, there are exceptions because we are discussing general purpose codes. For specific problems one may be able to break off groups of equations which can be solved easily, perhaps analytically, and so reduce the difficulty or size of the tasks otherwise addressed by the code. On the other hand, it is possible to get into trouble by such manipulations without realizing it. In [20] an example from chemical engineering is discussed for which assumptions of the kind described here were made and the scientists were led to the wrong steady state solution.

The example of the partial differential equation shows how large systems can arise. (Consider the situation with several space variables.) The heat equation, Van der Pol's equation, and the remarks we have made so far about steady-state approximations show that stiffness is not something unfamiliar; it has arisen, been studied, and dealt with, in other contexts. The example also exhibits the limited coupling that is usually present in large systems—here each equation involves only three unknowns or fewer, regardless of the number of equations in the system. Exploitation of this property will be mentioned in the fourth section. It is an important aspect of efficiency.

We have touched on several ways of realizing stiffness is present. Before discussing solution methods let us summarize them. Often one has a considerable understanding of the qualitative behavior of solutions of (2.1) on physical grounds. If the system is known to be very stable, it is likely to be stiff. If some variables are known to change on time scales very different from others and the physical problem is well posed, the governing equations are likely to be stiff. Analysis or experience with similar or model problems is often very useful. A common sign of stiffness is that a code aimed at nonstiff problems proves conspicuously inefficient for no obvious

reason (such as a severe lack of smoothness in the equation or the presence of singularities). There are codes for nonstiff problems [25], [27] which rather reliably diagnose stiffness automatically.

The integration during the transient has the step size limited by accuracy rather than considerations of stability so this part of a differential equation problem is not stiff. The distinction might be made more vivid by a little anecdote. The authors recently participated in a conference on the solution of stiff differential equations arising in models of lasers. One speaker discussed a refined physical model involving some 250 energy levels which was consuming hours of computer time. Though simpler models with these states aggregated into a smaller number of states were clearly stiff (they could be solved with current stiff codes far more efficiently than with nonstiff codes), his problem was being solved very inefficiently by codes aimed at stiff problems. Indeed, he found codes aimed at nonstiff problems to be conspicuously more efficient. Each solution component had the same qualitative behavior. After a while the level would become populated and the population would rapidly grow to a number about which it varied slowly. The difficulty originates in the fact that the various levels are populated successively and there are a great many of them. As far as the codes are concerned, they are always on the transient for some energy level. Eventually, of course, the whole system would get into a steady state and stiff methods would show their worth, but the cost of reaching this was prohibitive and the scientist was not particularly interested in the long-time behavior. The nonstiff methods perform better in the transient because this is what they are designed to do. Besides illustrating the role of transients, this example also points out that we do not have codes, or even methods, capable of adequately solving all the problems of scientific interest in the area of differential equations.

3. Characteristics of solution methods. All of the methods used in general purpose codes for the solution of stiff differential equations are implicit of necessity. This means that an equation must be solved at each step to obtain the numerical approximation. For example, in the backward Euler method,

$$y_{n+1} = y_n + h_n f(x_{n+1}, y_{n+1})$$

must be solved for y_{n+1} . If the problem is linear (that is, if f is linear), then a linear equation must be solved, but if the problem is nonlinear, a nonlinear equation must be solved. Simple functional iteration is used in codes for nonstiff problems to solve such equations:

$$y_{n+1}^{(m+1)} = y_n + h_n f(x_{n+1}, y_{n+1}^{(m)}).$$

For the model problem (2.2) we easily find that the iteration error satisfies

$$y_{n+1}^{(m+1)} - y_{n+1} = h_n A (y_{n+1}^{(m)} - y_{n+1}).$$

Once again we encounter the effects of stiffness—this simple iteration will not converge if we have a stiff problem (for which $|h_n A| > 1$). The usual procedure for stiff problems is some variant of Newton's method. This uses an approximate Jacobian J and one solves repeatedly the linearized system

$$y_{n+1}^{(m+1)} = y_n + h_n f(x_{n+1}, y_{n+1}^{(m)}) + h_n J (y_{n+1}^{(m+1)} - y_{n+1}^{(m)}).$$

For the model problem (2.2) the iteration error satisfies

$$(3.1) \quad y_{n+1}^{(m+1)} - y_{n+1} = (1 - h_n J)^{-1} (h_n A - h_n J) (y_{n+1}^{(m)} - y_{n+1}).$$

Any reasonable approximation J to A will cause this iteration to converge. If the problem is not stiff so that $h_n A$ is small, all that matters is that $h_n J$ be small. If the problem is stiff so that $h_n A$ is large and negative, all that matters is that $h_n J$ be within about 50% of $h_n A$. For this example, and for any linear problem, the iteration converges in one step if J is exactly A .

This analysis extends to systems, in which case the iteration (3.1) involves solving a linear equation at each iterate. If the problem is linear, only one iteration is necessary if we take $J = A$. In fact, most codes do *not* take $J = A$ so that iteration is used even in the linear case. When iteration is used, as it must for nonlinear problems, the starting approximation is very important because it must be good enough that the process will converge and moreover good enough that only a few iterations are necessary. Fortunately, the situation is such that a good starting approximation is almost always available by use of an explicit integration formula, usually called a predictor. Although a predictor is an integration formula, it does not have any of the main costs associated with a formula suitable for stiff problems because it is really a polynomial extrapolation process using information already known about the solution. For example, if the backward Euler formula is being used to solve stiff equations, the forward Euler formula can be used as a predictor. The forward Euler formula uses the function value y_n and the derivative y'_n computed in the last step to estimate the value of y_{n+1} to be used as the first iterate $y_{n+1}^{(0)}$. The accuracy of the predictor is as good, and usually better than that of the actual integrator (called the corrector) for stiff problems. The purpose of the corrector is to provide stability. When this technique is used along with a number of other techniques to detect convergence quickly, an average of less than 1.5 iterations of equation (3.1) are needed at each step in typical problems. (Using a predictor has several other advantages. In particular, the difference between the predictor and corrector provides a reasonable error estimator for local truncation error—an important part of any code. See Gear [12].)

It is worth noting that convergence of a quasi-Newton method for the corrector iteration is guaranteed for small enough step sizes because as the step size is reduced, an iteration such as (3.1) becomes a contraction. In addition, the accuracy of the predictor improves as h is reduced so that the initial approximation is more likely to be in the region of convergence.

The linearity of a problem does little to reduce the solution time in current codes for stiff problems. This could indicate a need for the development of better methods for linear problems, but it seems more likely that the reason is that the effects of nonlinearity are being handled very efficiently.

Returning now to the idea of pseudo steady-state approximations, we suppose that the equations can be written in the form

$$(3.2) \quad y' = f(x, y, z), \quad \varepsilon z' = g(x, y, z),$$

where the solution components are split into two groups y and z . Except in the transient, or boundary layer as it is commonly termed in this context, the small parameter ε suggests one neglect the term $\varepsilon z'$ and solve instead the algebraic-differential system

$$(3.3) \quad y' = f(x, y, z), \quad 0 = g(x, y, z).$$

In a number of significant physical applications this kind of approximation results in an easy (that is, nonstiff) set of ordinary differential equations and a set of algebraic equations. Even in this favorable situation it is not clear how to assess the errors which arise. Most often one has to solve the algebraic system by numerical means. If a code

for stiff equations is applied to equations (3.2) directly, it effectively solves equations (3.3) in the stiff region. This can be seen by considering the application of the backward Euler method to the second of equations (3.2) when $\varepsilon z'$ is small. We get

$$\frac{\varepsilon}{h_n}(z_{n+1} - z_n) = g(x_{n+1}, y_{n+1}, z_{n+1}) \doteq 0.$$

The fact that a code for stiff equations is automatically doing much the same thing as an analyst might, indicates the power of the kind of schemes we are discussing and shows how well-conceived and well-executed software can provide the casual problem solver with great assistance.

4. Codes for stiff problems. The backward Euler method, and similar methods such as the trapezoidal rule, were the first discovered which could handle stiff problems reasonably well so they became widely used. Because the large codes written to facilitate specific application areas such as simulation, chemical kinetics, circuit analysis, and the like are very slow to change, these methods are still seen in practice. More efficient methods are now available as a result of researches into higher order methods. The initial transients alone represent difficult nonstiff problems and the great success of high order methods for nonstiff problems has encouraged such investigations. There has been no lack of ideas for high order procedures with excellent stability but improvements have not been won easily because each seems to raise some new difficulty. Relatively few ideas have been implemented as software of sufficiently high quality to merit general consideration. We shall mention a few such ideas in order to allude to some of the difficulties later. Also, we shall mention specific codes because there are so few general purpose codes for stiff problems which are widely available. Because in some cases we do not have experience with the codes we are not necessarily endorsing them; we do have reason to believe that all are serious attempts at providing mathematical software for this problem so we hope that learning of them will prove useful to the reader.

Using the general idea of extrapolation [19] there is a code of Schryer [22] which does repeated extrapolation of the backward Euler formula. In doing this one adapts the order to the requirements of the problem. Lindberg [21] has written a code which extrapolates the midpoint rule a single time to generate a method of order four. The popular Runge-Kutta methods furnish methods suitable for stiff equations if one considers implicit methods. Hulme [18] has written a code which provides an arsenal of implicit Runge-Kutta methods; there are two families of methods each with a large range of orders. Though a fixed formula is used for the integration, the code can select an appropriate one automatically. The most popular formulas being used in general purpose codes are the backward differentiation formulas (see Gear [9]). Like their relatives the Adams formulas, these formulas are usually implemented so as to automatically choose the step size and the order. An early code is that of Gear [11]. It has been often reprogrammed and various improvements made; a generally available and widely used version is that of Hindmarsh [14]. Extensions which are reported to improve the efficiency of the methods are given by Skeel and Kong [28]. These extensions "blend" the backward differentiation formulas with Adams formulas in a ratio determined by the size of the Jacobian, so that they look like Adams formulas for nonstiff problems and backward differentiation formulas for stiff problems. A different form of extension appears in Hindmarsh and Byrne [16], where a different step changing formula believed to lead to better behavior in problems with multiple transient regions is used. Bickart, Tendler, and Picel have developed codes for two

subsets of composite multistep methods—methods invoking a set of multistep formulas solved simultaneously for a set of solution points. The cyclic methods implemented in [29] may be thought of as generalizations of the backward differentiation formulas and the one-step methods in [2] as generalizations of the Runge–Kutta methods.

The backward differentiation formulas illustrate an important point about codes for stiff problems. As the stability plots in [10] show, if the Jacobian has eigenvalues near the imaginary axis, the higher order formulas will be unstable. The unstable region increases as the order does to the point that formulas of order greater than six are not stable at all. Most codes do not use the sixth order formula because of this limitation though it is stable in most of the left half plane. The formulas of order five and lower do not have much of a limitation of this kind but occasionally one notices it with a real problem. Just because one has a good code implementing a good method does not mean that he will not encounter difficulties with stiffness. One needs some understanding of the characteristics of his code and it is best to have a repertoire of methods on which to draw.

The implicit Runge–Kutta methods avoid the problems associated with the high-order backward differentiation formulas because they can be stable in the whole left half plane. However, there is a price associated with this additional stability. The system of equations which must be solved is two or more times larger, implying a large increase in storage for large systems. When these methods are effective, it is because they can use a much longer step than methods based on backward differentiation. Sometimes, however, the solution will change character as we move from a smooth stiff region back into a rapidly varying transient within one long step which is then wasted. To some degree the composite multistep formulas share this difficulty though they have some compensating factors. The difficulty is even more pronounced with extrapolation methods which gain their speed from a very long basic step.

A good code selects its step size, and hence x_n , automatically. If the user is prepared to accept output at only the points x_n , there is no impact on the efficiency of the integration. If, however, the user must compute the solution at many other points, the cost can become prohibitive. If there are a sufficient number of integration points that the desired accuracy can be achieved by interpolation, the cost is small and there is no effect on the integrator. In fact, most codes based on the backward differentiation formulas have scaled derivatives or divided differences available so that the interpolation cost is minimal; the better codes provide interpolation subroutines for the user. If, on the other hand, the code must be asked to compute intermediate points by integration, the cost can be high, particularly in codes which use long steps such as implicit Runge–Kutta and especially in extrapolation.

We see that the solution of stiff problems is practical provided we are prepared to solve at each step a linear system based on the Jacobian. This can have serious consequences. The function f has N components but in general its Jacobian has N^2 components. For a general problem of even moderate size, providing the Jacobian analytically can be burdensome for the programmer and it can be very difficult to be sure that the partial derivatives have been obtained correctly. For these reasons, automatic generation of a Jacobian by numerical differencing is an indispensable part of a general purpose code.

The evaluation of Jacobians is a relatively expensive part of the solution process even though good codes try to do this as infrequently as possible. The key to reducing this cost is to take advantage of special structure of the problem. Ordinarily, medium to large systems have weak coupling so that most partial derivatives are zero; linear

terms in the equations giving rise to immediately available partial derivatives are fairly common; and often the functions are rather simple in form, especially for very large systems. It is easier to take advantage of these possibilities with analytic differentiation than numerical differentiation with the consequence that although the former may be more trouble to set up, it will prove very much cheaper. Those concerned with application packages should be especially alert to this possibility because of the restriction to a special class of equations in their area. An example is the package [5] in which the user describes his chemical kinetics problem in a manner natural to him as a chemist. The package sets up the differential equations and forms the (easy) Jacobian analytically on its own. At some computing installations symbol manipulation languages furnish a quite practical tool for the generation of Jacobians. Depending on the problem and the efficiency of the code generated by the symbolic processor, the analytical derivatives may, or may not, be evaluated more cheaply than by differencing.

One of the reasons for preferring analytical Jacobians is that scaling difficulties in differencing can lead to poor numerical Jacobians. Good differencing algorithms try to cope with this sort of trouble and, of course, the codes will work with poor Jacobians anyway, but the net effect is that scaling troubles increase the cost and decrease the reliability of differencing as compared to analytical Jacobians.

An intermediate procedure is for the user to provide structure information which the code then uses to construct Jacobians by differencing. A trick discovered by Curtis, Powell and Reid [4] may achieve substantial reductions in the cost of differencing because of known structure. A particularly simple and quite common case is that of a banded system. We say the system is banded, with half bandwidth m , if for each i , the j th equation does not involve the variables y_j for $|j-i| > m$. For such a system the Jacobian has nonzero elements only within a band about the main diagonal, whence the name. As it turns out, a Jacobian of half bandwidth m can be formed by differencing in only $2m+1$ extra evaluations of f regardless of the size of the system of equations. Clearly this leads to large reductions in machine time even when m is not particularly small compared to N .

Storage can be a limiting factor in the solution of large stiff problems because one must store the Jacobian or a closely related matrix for which one must solve a linear system of equations. A major disadvantage of most implicit Runge-Kutta schemes is that the size of the nonlinear system to be solved at each step is a multiple of the number of equations, the higher the order, the bigger the multiple. Because the Jacobian grows like the square of the number of unknowns, for even moderate sized systems the order in the code COLODE [18] has to be restricted severely to hold the program in rapid memory.

When solving nonstiff systems of equations, one is accustomed to ignoring all costs except that of evaluating the equation. With stiff systems this cannot be done because the cost of repeatedly solving linear systems often represents a substantial fraction of the total cost. We shall examine the efficient solution of these systems, but first we note that methods behave differently with respect to the size and number of systems which must be solved. The origin of the difficulty is in the fact that the iteration matrix depends explicitly on the step size and so must be treated afresh if the step size is significantly altered. Because of this certain algorithms must spend an unusually large portion of their effort on solving linear systems. For example, in extrapolation procedures, assuming that the approximate Jacobian will serve for a whole step, one repeatedly advances a method such as the backward Euler method through this interval with successively smaller fractions of the basic step size, and combines these results at the end of the basic step. Each sweep involves a different

linear system which must be solved repeatedly at each fractional step of the sweep. Other examples for which the cost of the linear algebra is prominent for similar reasons include the composite multistep methods, implicit Runge-Kutta methods, methods based on averaging various implicit solutions at each step, and the use of step halving to estimate local errors. This situation has been the object of intensive research and the recent papers of Enright [8] and Bickart [1] show that progress is being made.

For small to moderate sized systems one can use elimination to factor the matrix and then use these factors to do the necessary iterations efficiently. To solve medium to large problems we must again resort to structural information. Band structure is the easiest to accommodate. Because the factors of a band matrix are also of band form it is possible to compute and store only those elements known to be potentially nonzero. If the band width is relatively small, $m \ll N$, one greatly reduces the storage and greatly increases the efficiency of the solution of the linear system by taking advantage of the band form. More generally a matrix is said to be *sparse* if most of its elements are zero. This is typical of many very large stiff problems. There are schemes for storing only the nonzero elements of sparse matrices and for doing elimination in such a way as to introduce relatively few nonzero elements. Some storage is required for the schemes themselves and some extra effort in the elimination, so they do not begin to pay off in either storage or speed until only a few percent of the entries in the matrices are nonzeros. Fortunately this is often the case. Problems have been solved which are so large that iterative methods had to be used for the linear systems because such methods require less storage than a factorization. This is an active field of research at the present and we may well see iterative techniques being used for less special problems.

Exploitation of structural information has been incorporated in a number of codes, but they are not as widely available or as easy to use as those codes cited at the beginning of this section. An exception is the code of Hindmarsh [17] which is a version of the code [14] modified to use a band structured Jacobian. It uses the differencing scheme and the special solution of banded systems which we have described. One of the earliest such codes is the sparse tableau approach. This is described in Hachtel, et al. [13], and is apparently part of an IBM proprietary software package for electronic circuit design.

Since it is often true that evaluation of the function and even the Jacobian are not particularly expensive for large problems, the extensive linear algebra and data transfers are often a large fraction of the total cost. FORTRAN does not take full advantage of the overlapping of data transfers and computation and of hardware possibilities of pipeline and vector machines. FORTRAN callable assembly language modules are becoming available which allow more efficient handling of basic tasks than are possible in the FORTRAN codes being widely disseminated. The report [15] gives some codes and timing comparisons on a CDC 7600 of routines for the LU decomposition of a matrix and the forward and backward substitution process for solving a linear system. For only three equations, the FORTRAN programs are slightly faster because of some additional overhead in the linear algebra modules, but the more specialized routines quickly pull ahead and show an improvement of a *factor* of about three for 200 equations. (The solving routine benefits a little more than the decomposition routine.) For large problems there are clearly important cost savings which can be achieved on some machines in this manner.

Many problems lead to implicit sets of differential equations such as

$$(4.1) \quad My' = Ky + p(x),$$

where the matrices M and K may depend on y , x , and even y' . It is not necessary, nor even desirable, to invert these to the explicit form

$$(4.2) \quad y' = M^{-1}[Ky + p(x)],$$

when a quasi-Newton iteration is used for the corrector. If, for example, the backward Euler method is used as an integrator, it can be substituted in

$$(4.3) \quad F(y, y', x) = 0$$

at $x = x_{n+1}$ to eliminate y'_{n+1} and get

$$F\left(y_{n+1}, \frac{y_{n+1} - y_n}{h}, x_{n+1}\right) = 0$$

which is to be solved for y_{n+1} . Dealing with the implicit equation (4.3) directly can save arithmetic and storage in large sparse problems. This is clear from equations (4.1) and (4.2) when M and K are constant. Whereas the Jacobian of (4.1) will be $K - M/h$, which will be sparse, the Jacobian of (4.2) is $M^{-1}K - I/h$, which will normally be dense.

A very important part of any code for solving differential equations, whether stiff or nonstiff, is automatic control of the step size and order of method used. While the basic strategy is straightforward—the two are chosen to try and minimize the amount of work done to integrate over the interval—the implementation is not. The problem lies in the fact that there is not yet an adequate theory to tell us how to choose these parameters, so the choice is based on the extension of existing theory to situations in which it probably does not apply, coupled with a lot of testing and tuning of codes to make them as reliable as possible. Because of this, it is possible for a person to implement a set of formulas into a code and use the same basic step and order control strategy as another code, and yet have a code that is orders of magnitude slower on some difficult problems. (This is illustrated in [23] which compares a number of codes for nonstiff problems, including four based on Adams methods and rather similar step and order control strategies.) The lesson to be learned from this is that, whenever possible, an existing piece of well tested and well documented software should be used, and if possible, it should be used without change.

Variable order codes adapt the formula used to the observed characteristics of the solution and have proved very efficient in general use. One should appreciate that such codes are relatively inefficient during the initial stages of the computation while they find an appropriate order. Ordinarily this is an unimportant part of the overall expense, but there are exceptions. We have seen examples in flame chemistry where one has partial differential equations describing the gas flow coupled to ordinary differential equations describing reactions taking place in the flow. A pseudo steady-state approximation is made in which one holds the gas flow parameters constant for a time period during which one integrates the ordinary differential equations. Using these values one re-solves the partial differential equations for the time zone or advances into the next time zone—the details do not matter here. The point is that at the end of each zone, the ordinary differential equations change. Codes based on backward differentiation formulas prove very inefficient because they must be restarted at each zone. Things do not change much from zone to zone (else they would be shortened) so that the last step size used in one zone would be reasonable in the next if one retained the same formula. Fixed order methods appear to offer considerable advantage in this situation and the limited experience agrees. To some degree, extrapolation methods which permit high orders must also be at a disadvantage

because the order and the length of the step will be held back by the length of the zone. A low order and small step size can be obtained more cheaply by a code with fixed order.

There is one situation in which automatic step and order control may not be desirable; that is, when a problem is to be repeatedly integrated to study, for example, effects of parameters. Here it is desirable to use the same sequence of steps and formulas in each integration. If that is done, the numerical solution depends more smoothly on the parameter values than if not. However, even in this case, it is desirable to choose a step and formula sequence automatically for the first integration and either perform the other integrations in parallel (and possibly use the same Jacobian and its decomposition) or save the sequence for subsequent integrations.

5. Solving a differential equation. There are a number of meanings of "solving" a differential equation numerically. Sometimes one wants an approximate solution at a single point. More often a solution is desired on a set of points so as to get a table of the solution. Other times one wants either a continuous approximate solution curve or a table so dense as to be equivalent. It is important to realize that where and how frequently one wants output points can have a serious effect on the cost. The most efficient action depends on the method implemented, the principal factor being how far, relatively speaking, the code steps before producing an approximate solution. For example, high order implicit Runge-Kutta schemes and extrapolation schemes advance very much further in a step than does a code based on the backward differentiation formulas. The former produce output at a given point by stepping so as to actually hit the point. Requesting output more frequently than the natural step size chosen by the code will severely degrade the efficiency of such methods for a variety of reasons. If one seeks an answer at only a few points, especially if he is willing to accept the natural output points, output is not a problem. Some methods, like the backward differentiation formulas, provide continuous polynomial solutions which can be evaluated at negligible cost. Not all codes implement this equally well. In this matter one needs to consider what is required in the way of output and how it impacts the repertoire of codes.

The user's meaning of accuracy can affect the results considerably. One common scheme is to measure the error relative to the maximum (absolute) value of the solution component seen so far in the integration. Another common scheme is a mixture of absolute error and error relative to the solution magnitude. It is important to be able to specify error tolerances for each component of the solution because scaling of components often differs radically for stiff problems. Stiff problems almost always involve transients during which the solution changes sharply. In the survey [26] about half the users required the accurate solution of these transients as well as that of the slowly varying portions. To do this one must use a suitable error control and the net effect is that in the transient, accuracy dominates the choice of step size rather than stability. As a result the transients become a relatively expensive part of the integration for most codes. Even if the user is interested only in the long term behavior of the solution, a resolution of the transient may be necessary to assure that the proper equilibrium solution is picked up and to assure reliability of the basic algorithms. If one knows that all errors will be heavily damped and is interested only in equilibrium behavior, one can economize by computing the transient crudely.

An extremely common misuse of codes is to seek a solution accurate in a relative sense when the solution vanishes initially or tends to zero rapidly later in the integration. The attempt can prove exceedingly expensive and is rarely meaningful. By far

the most common situation is that, when a quantity drops below a certain level, the user of the code is no longer interested in this quantity and so should not waste time trying to compute it accurately. A proper choice of error criterion will accomplish this. Of course, it may happen that nonphysical values get computed in this way, for example, a negative concentration. Ordinarily this is unimportant, but it can happen that the differential equations become unstable if such a value should be generated. Discussion of this difficulty in the context of mass action kinetics and an example can be found in [6].

Stiff problems are relatively expensive to solve and the expense depends much more strongly on the tolerance than is true of the best codes for nonstiff problems. The physical origin of stiff problems rarely makes high accuracy meaningful because fundamental quantities are known inaccurately. The case of semi-discretization of partial differential equations is an easily understood and important case in point. If the spatial discretization is crude, it makes no sense to solve the ordinary differential equations very accurately at all. In the survey [26] accuracies of one or two digits were by far the most common requests. An accuracy of five digits was considered stringent. Apparently experience says that accuracies in this general area represent a bearable expense with our currently available codes and machines. At this point we should remind the reader that the codes control their local truncation errors, not the global errors which interest the user of the code. The present state of the art is such that for some problems one can get reasonable looking numbers which are not close to the desired solution. This difficulty can arise in a number of ways, including a step control routine that makes the step so large that an active region of the solution is missed entirely, and a formula that is stable for unstable problems and completely ignores an increasing component of the solution. A conservative choice of tolerance, alertness to scaling considerations in the error criterion and differencing of Jacobians, experimentation, and a thoughtful examination of the numerical results are indispensable for solving stiff differential equations.

Acknowledgment. We would like to acknowledge the valuable comments and suggestions of the referees.

REFERENCES

- [1] T. A. BICKART, *An efficient solution process for implicit Runge-Kutta methods*, SIAM J. Numer. Anal., to appear.
- [2] T. A. BICKART AND Z. PICAL, *High order stiffly stable composite multistep methods for numerical integration of stiff differential equations*, BIT, 13 (1973), pp. 272-286.
- [3] G. BJUREL, G. DAHLQUIST, B. LINDBERG, S. LINDE AND L. ODEN, *Survey of stiff ordinary differential equations*, Rep. NA 70.11, Dept. of Information Processing, Royal Institute of Technology, Stockholm, 1970.
- [4] A. R. CURTIS, M. J. D. POWELL AND J. K. REID, *On the estimation of sparse Jacobian matrices*, J. Inst. Math. Appl., 13 (1974), pp. 117-119.
- [5] L. EDSBERG, *Integration package for chemical kinetics*, Stiff Differential Systems, R. A. Willoughby, ed., Plenum Press, New York, 1974, pp. 81-94.
- [6] ———, *Numerical methods for mass action kinetics*, Numerical Methods for Differential Systems, L. Lapidus and W. E. Schiesser, eds., Academic Press, New York, 1976, pp. 181-195.
- [7] W. H. ENRIGHT, T. E. HULL AND B. LINDBERG, *Comparing numerical methods for stiff systems of ODEs*, BIT, 15 (1975), pp. 10-48.
- [8] W. H. ENRIGHT, *Improving the efficiency of matrix operations in the numerical solution of stiff ODEs*, Rep. no. 98, Dept. of Computer Science, University of Toronto, Toronto, Canada, 1976.
- [9] C. W. GEAR, *The automatic integration of stiff ordinary differential equations*, Information Processing, 1969, pp. 187-193.

- [10] ———, *Numerical Initial Value Problems in Ordinary Differential Equations*, Prentice-Hall Inc., Englewood Cliffs, NJ, 1969.
- [11] ———, *Algorithm 407: DIFSUB for solution of ordinary differential equations [D2]*, *Comm. ACM*, 14 (1971), pp. 185–190.
- [12] ———, *Estimation of errors and derivatives in ordinary differential equations*, *Information Processing*, 3 (1974), pp. 447–451, 1974.
- [13] G. D. HACHTEL, R. K. BRAYTON AND F. G. GUSTAVSON, *The sparse tableau approach to network analysis and design*, *IEEE Trans. Circuit Theory*, CT-18 (1971), pp. 101–113.
- [14] A. C. HINDMARSH, *GEAR: Ordinary differential equation solver*, Rep. no. UCID-30001, rev. 3, Lawrence Livermore Laboratory, Livermore, CA, 1974.
- [15] A. C. HINDMARSH, L. J. SLOAN, K. W. FONG AND G. H. RODRIQUE, *DEC/SOL: Solution of dense systems of linear algebraic equations*, Rep. UCID-30137, Lawrence Livermore Laboratory, Livermore, CA, 1976.
- [16] A. C. HINDMARSH AND G. D. BYRNE, *EPISODE: An experimental package for the integration of ordinary differential equations*, Rep. no. UCID-30112, Lawrence Livermore Laboratory, Livermore, CA, 1975.
- [17] A. C. HINDMARSH, *GEARB: Solution of ordinary differential equations having banded Jacobian*, Rep. no. UCID-30059, rev. 1, Lawrence Livermore Laboratory, Livermore, CA, 1975.
- [18] B. L. HULME, *COLODE: A collocation subroutine for ordinary differential equations*, Rep. SAND74-0380, Sandia Laboratories, Albuquerque, NM, 1974.
- [19] D. C. JOYCE, *Survey of extrapolation processes in numerical analysis*, this Review, 13 (1971), pp. 435–488.
- [20] L. LAPIDUS, R. C. AIKEN AND Y. A. LIU, *The occurrence and numerical solution of physical and chemical systems having widely varying time constants*, *Stiff Differential Systems*, R. A. Willoughby, ed., Plenum Press, NY, pp. 187–200, 1974.
- [21] B. LINDBERG, *A stiff system package based on the implicit midpoint method with smoothing and extrapolation*, *Ibid.*, pp. 201–215.
- [22] N. L. SCHRYER, *An extrapolation step size monitor for solving ordinary differential equations*, *Proc. ACM*, 1974, pp. 140–148.
- [23] L. F. SHAMPINE, H. A. WATTS AND S. M. DAVENPORT, *Solving nonstiff ordinary differential equations—The state of the art*, this Review, 18 (1976), pp. 376–411.
- [24] L. F. SHAMPINE, *Stiffness and non-stiff differential equation solvers*, *Numerische Behandlung von Differentialgleichungen*, L. Collatz, ed., ISNM 27, Birkhäuser Verlag, Basel, Switzerland, 1975, pp. 287–301.
- [25] L. F. SHAMPINE AND M. K. GORDON, *Computer Solution of Ordinary Differential Equations; The Initial Value Problem*, W. H. Freeman, San Francisco, CA, 1975.
- [26] ———, *Typical problems for stiff differential equations*, *SIGNUM Newsletter*, 10 (1975), p. 41.
- [27] L. F. SHAMPINE, *Stiffness and non-stiff differential equation solvers II: detecting stiffness with Runge-Kutta methods*, *ACM Trans. Math. Software*, 3 (1977), pp. 44–53.
- [28] R. D. SKEEL AND A. K. KONG, *Blended linear multistep methods*, UIUCDCS-R-76-800, Dept. of Computer Science, University of Illinois, Urbana, IL, 1976.
- [29] J. M. TENDLER, T. A. BICKART AND Z. PICEL, *A stiffly stable integration process using cyclic composite methods*, to appear.