# ANALYSIS AND IMPROVEMENT OF UPWIND AND CENTERED SCHEMES ON QUADRILATERAL AND TRIANGULAR MESHES

H. T. Huynh

NASA Glenn Research Center, MS 5-11

Cleveland, Ohio 44135

Phone: (216) 433-5852; E-mail: huynh@grc.nasa.gov

**Abstract.** Second-order accurate upwind and centered schemes are presented in a framework that facilitates their analysis and comparison. The upwind scheme employed consists of a reconstruction step (Van Leer 1977) followed by an upwind step (Roe 1981). The two centered schemes are of Lax-Friedrichs (L-F) type. They are the non-staggered versions of the N-T scheme (called ORD in Nessyahu-Tadmor 1990) and the CE/SE method with $\epsilon = 1/2$ (Chang 1995). The upwind scheme is extended to the case of two spatial dimensions (2D) in a straightforward manner. The N-T and CE/SE schemes are extended in a manner similar to the 2D extensions of the CE/SE schemes by Wang and Chang (1999) and Zhang et al. (2002); the slope estimates, however, are simplified. Fourier stability and accuracy analyses are carried out for these schemes for the standard 1D and the 2D quadrilateral mesh cases. In the nonstandard case of a triangular mesh, the triangles must be paired up when analyzing the upwind and N-T schemes. An observation resulting in an extended N-T scheme which is faster and uses only one third of the storage for flow data compared with the CE/SE method is presented. Numerical results are shown. Other improvements to the schemes are discussed.

**Introduction.** When solving a fluid flow problem, a researcher has the option of choosing between upwind and centered schemes using a quadrilateral or a triangular mesh. In this paper, trade-offs among these choices are discussed. Relations between schemes and their strengths and weaknesses are shown, and improvements are suggested.

The schemes employed are among the simplest second-order accurate schemes that can capture shocks and deal with unsteady problems. Both the upwind and centered schemes here use piecewise-linear reconstructions, i.e., MUSCL interpolants (monotone upwind schemes for conservation laws, Van Leer, 1977), which extend Godunov's piecewise constant method (1959). The key difference is that for the upwind scheme, numerical dissipation is added by the upwind step (Roe 1981, 1986), while for the centered schemes, dissipation is obtained by, loosely put, averaging the neighboring data (scheme ORD of Nessyahu and Tadmor 1990).

The upwind step has a few drawbacks. Roe's flux-difference splitting, which is mathematically rigorous and among the most popular, may cause oscillations as in the case of a slow-moving shock, or instability as in the carbuncle problem. The AUSM scheme (Liou and Steffen Jr. 1992, Wada and Liou 1997) does not have these problems, but it is not clear to this author why the scheme works. The upwind step is also sometimes perceived to be costly and difficult to grasp. In spite of these problems, upwind schemes are popular because they work well for a large class of flows. The upwind step employed here is Roe's splitting with an entropy correction described in Huynh (1995a). It can be derived by diagonalization and coded by stepping across one acoustic wave. The resulting scheme is concise and economical; the presentation below is also simpler than most presentations in the literature. Numerical solutions obtained with this upwind scheme for the 1D Euler equations can be found in Huynh (1995a,b). The 2D extension of this scheme is conceptually straightforward. For other versions of upwind schemes, see, e.g., Barth and Jespersen (1989), Roe (1989), Hirsch (1990), and Venkatakrishnan (1995).

To avoid upwinding, a second-order accurate scheme, which extends the first-order scheme of Lax-Friedrichs (L-F), was introduced by Nessyahu and Tadmor (1990). There, the reconstruction step is the same as that of the upwind scheme, but

the upwind step is avoided by the use of a staggered mesh. The scheme employed here, however, is the nonstaggered version obtained by overlaying two staggered meshes to form a regular mesh. The drawback of the nonstaggered version is that the computing time doubles. The main reason this version is chosen is that the 2D extension retains its key advantage of simplicity especially when the mesh is triangular. In addition, the mesh and the boundary conditions can be chosen to be the same as those of the upwind scheme. The staggered version, on the other hand, requires two sets of meshes and two sets of boundary conditions—compromising the advantage of simplicity. Moreover, for a triangular mesh, the extension of the staggered version is quite involved (Arminjon et al. 1997) and, as will be discussed, the advantage in computing time may no longer hold. Therefore, unless otherwise stated, *we deal only with the the nonstaggered version below.* Note that since there is no one-sided bias, the N-T scheme is centered. Also note that these L-F type centered schemes are different from the semidiscrete centered schemes made popular by Jameson et al. (1984).

While the numerical dissipation of the N-T scheme is a lot less than that of the L-F method, it is still considerable. The CE/SE or conservation element and solution element method (Chang 1995) provides a way to adjust dissipation for these centered schemes. Compared with the N-T scheme, the mesh, the balancing of fluxes, and the updates of the cell average quantities are essentially identical. The difference is in the calculation of the slopes (of the linear interpolant). For CE/SE, the slopes must be stored, and due to the way the slopes are updated, numerical dissipation can be adjusted via a parameter called $\epsilon$. For the general CE/SE scheme ($\epsilon \neq 1/2$), the slope calculation is quite different from that in a typical MUSCL approach. When $\epsilon = 0$, the scheme has no numerical dissipation, i.e., it is reversible in time. Currently, the CE/SE member employed in essentially all practical calculations corresponds to $\epsilon = 1/2$. For this reason, we restrict our attention to this member and, from this point on, unless otherwise stated, the term *CE/SE refers to the member with $\epsilon = 1/2$.* Note that there are numerous differences in terminology between (Nessyahu and Tadmor 1990) and (Chang 1995); here, the terminology in the former is often employed.

The CE/SE schemes were extended to 2D for unstructured triangular meshes by Wang and Chang (1999) using the nonstaggered version. What is novel about this extension is that the spatial domain where each reconstruction is valid at the beginning of the time level is a hexagon, which is roughly twice as big as the triangular cell. This CE/SE approach to extension is also applied here to the N-T scheme. (Such an extension was mentioned as a coupled version for the CE/SE scheme in Chang et al. (1999) and has recently been incorporated—independently from this work—as an option in the CE/SE code; private communication with Drs. Ananda Himansu, Ching Y. Loh, and Xiao-Yen Wang. Note, however, that the extended N-T scheme presented here has numerous differences resulting in a scheme which is faster and requires considerably less storage.) The quadrilateral-mesh extension for the CE/SE method can be found in Zhang et al. (2002); see also Cook (1999). The CE/SE schemes have been applied to solve numerous practical problems in two and three dimensions with a lot of success, especially in aeroacoustics.

For a structured quadrilateral mesh, in a manner similar to the 1D case, the nonstaggered mesh in Zhang et al. (2002) can be obtained by overlaying two staggered meshes in Arminjon et al. (1995) and Jiang and Tadmor (1998) (see also Jiang et al. 1998). For a triangular mesh, however, a similar statement does not hold; in fact, the nonstaggered extension in Wang and Chang (1999) appears to have numerous advantages over a staggered-mesh extension (remark (c) in §6 below).

In this paper, the schemes involved are first presented for the 1D advection equation where key ideas and trade-offs can already be seen. It is shown that the N-T and the CE/SE ($\epsilon = 1/2$) schemes are respectively the centered counterparts of Van Leer's first and second upwind schemes. Next, the extensions of the upwind, N-T, and CE/SE schemes to the 1D Euler equations are explained. Then, extensions to the 2D Euler equations on a quadrilateral and a triangular mesh as well as the simplified slope estimates are described. Comparison of schemes via Fourier stability and accuracy analyses are carried out. Here, for a triangular mesh, we must pair up the downward and upward pointing triangles when analyzing the upwind and N-T schemes.

Concerning the two L-F type methods, the N-T and CE/SE schemes are shown to produce essentially the same numerical solutions. The former has the advantage of better coupling; consequently, it converges better for a steady state problem. In addition, the following observation yields an extended N-T scheme which is faster and requires consider-

ably less storage than the CE/SE method (slopes need not be stored): instead of gathering fluxes to update the solution, fluxes are distributed to the neighboring cells. This observation—made possible in part by the simplified slope estimates—results in each cell needing to be visited only once during each time step instead of three or four times as in the case of the CE/SE scheme (three for a triangular mesh, and four for a quadrilateral mesh).

This paper is essentially self-contained. Readers who are interested only in the Euler equations can start with §4. The paper is organized as follows. In §1, the first-order accurate upwind and centered schemes are presented as two different ways of stabilizing the Euler forward scheme. In §2, upwind schemes are developed for the 1D advection equation using piecewise linear reconstructions. The centered counterparts of these upwind schemes are described in §3; here, a comparison via Fourier analysis is carried out. Section 4 deals with extensions to the case of the 1D Euler equations; §5, the 2D Euler equations on quadrilateral meshes; and §6, on triangular meshes. In §7, Fourier analysis for the 2D case is carried out. Numerical examples are shown in §8, and conclusions are presented in §9.

The author wishes to thank Prof. B. van Leer, Mr. C. Steffen, Jr., Drs. S.-C. Chang, R. Chima, A. Himansu, D. Jacqmin, P. Jorgenson, M.-S. Liou, C. Y. Loh, A. Suresh, and X.-Y. Wang for several illuminating discussions. This work was supported by the Engine System Noise Reduction project, which is part of the Quiet Aircraft Technology Program at the NASA Glenn Research Center.

## 1. First-order accurate schemes for advection.

The first-order case is straightforward, but it conveys the ideas and the trade-offs. In addition, improvements and nonstandard observations will be made. The simplest discretization—forward-time centered-space (FTCS)—leads to the Euler forward scheme, which is unstable. To stabilize this scheme, we need numerical dissipation. Dissipation can be added by using an upwind-biased difference, which results in an upwind scheme, or by using a dissipative start-off value, which results in a centered scheme.

For simplicity, consider the advection equation with constant speed $a$,

$$\frac{\partial u}{\partial t} + a \frac{\partial u}{\partial x} = 0, \qquad (1.1a)$$

$$u(x, 0) = u_0(x), \qquad (1.1b)$$

where $t$ is time, $x$ distance, and $u_0(x)$ the initial condition. By assuming that $u_0(x)$ is periodic, boundary conditions are straightforward and are omitted. The exact solution is

$$u(x, t) = u_0(x - at). \qquad (1.2)$$

To discretize the above problem, let $h$ be the mesh spacing, and $x_j = jh$, $j = 0, 1, 2, \ldots$ be a uniform mesh. Let $\Delta t$ be the time step and $t^n = n\Delta t$ be the time level. At time $t^n$, let $u_j^n$ be an approximation to the the solution $u$ at $x_j$. Assume that we know $u_j^n$ for all $j$; we wish to calculate $u_j^{n+1}$. To simplify the notation, the superscript $n$ in $u_j^n$ is omitted and the data is denoted by $u_j$; all other superscripts, however, are retained.

Next, set

$$\sigma = a\Delta t/h. \qquad (1.3a)$$

Then, the quantity $|\sigma|$ is the Courant number. Assume that the time step satisfies the CFL (Courant, Friedrichs, and Lewy) condition

$$|\sigma| \leq 1; \qquad (1.3b)$$

loosely put, information propagates no more than one mesh size per time step.

**1.1. Euler forward scheme.** This scheme is given by the FTCS differencing,

$$u_j^{n+1} = u_j - \sigma \tfrac{1}{2}(u_{j+1} - u_{j-1}). \qquad (1.4)$$

Again, $u_j^n$ is abbreviated to $u_j$.

For the Fourier (or Von Neumann) stability analysis, set $x_j = j$ and

$$u_j = e^{Ijw},$$

where $I = \sqrt{-1}$ and $w$ is the wave number, $-\pi \leq w < \pi$. Equation (1.4) implies

$$u_j^{n+1} = e^{Ijw}\left[1 - \sigma \tfrac{1}{2}(e^{Iw} - e^{-Iw})\right].$$

The quantity in the square brackets is the amplification factor of the Euler forward scheme:

$$\mathcal{A} = 1 - \sigma\left[\tfrac{1}{2}(e^{Iw} - e^{-Iw})\right]. \qquad (1.5)$$

If $a \geq 0$, then, the CFL condition (1.3b) implies $0 \leq \sigma \leq 1$. Next, the slope $u_x$ by central difference at $j = 0$ is given by $\frac{1}{2}\overrightarrow{CE} = \frac{1}{2}(e^{Iw} - e^{-Iw})$ shown in Fig. 1.1, where the point $C$ corresponds to $j = -1$, $A$ to $j = 0$, and $E$ to $j = 1$. Thus, the amplification factor lies on the line segment AB, where A corresponds to $\sigma = 0$ and B, $\sigma = 1$. Since this range is outside the unit circle, the scheme is unstable. Also
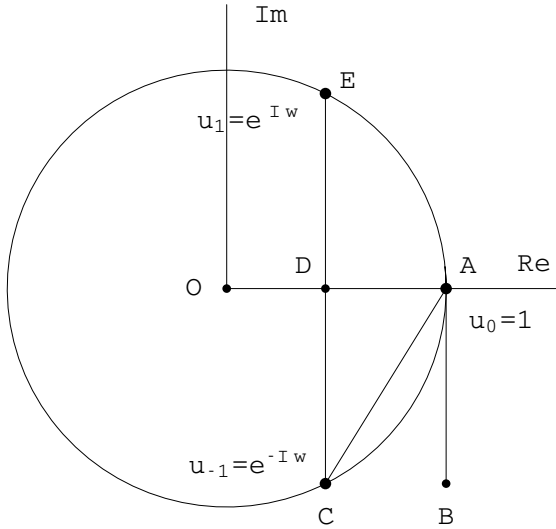
Fig. 1.1. *Fourier Analysis. The circle is of radius 1 in the complex plane. The amplification factor for the Euler forward scheme lies on the line segment AB; that for the first-order upwind scheme, AC; and that for the L-F scheme, DC. To reduce numerical dissipation for the L-F scheme, instead of D, we can employ a start-off value which blends D and A, e.g., $(1 - |\sigma|)u_j + |\sigma|\frac{1}{2}(u_{j-1} + u_{j+1})$.*

note that after a time step corresponding to $\sigma$, the exact solution at $j = 0$ is $e^{-I\sigma w}$. The amplification factor is an approximation of this exact solution.

**1.2. First-order accurate upwind scheme.** The Euler forward scheme can be stabilized by an upwind-biased difference for $(u_x)_j$. The result is (Courant, Isaacson, and Rees, 1952),

$$u_j^{n+1} = \begin{cases} u_j - \sigma(u_j - u_{j-1}) & \text{if } a \geq 0, \\ u_j - \sigma(u_{j+1} - u_j) & \text{otherwise.} \end{cases} \quad (1.6)$$

The amplification factor for the first-order upwind scheme is

$$\mathcal{A} = \begin{cases} 1 - \sigma(1 - e^{-Iw}) & \text{if } a \geq 0, \\ 1 - \sigma(e^{Iw} - 1) & \text{otherwise.} \end{cases} \quad (1.7)$$

If $a \geq 0$, then $0 \leq \sigma \leq 1$, and the amplification factor lies on the line segment AC shown in Fig. 1.1. Since this range is inside the unit circle, the scheme is stable.

Notice that if the wind direction is incorrectly determined—for systems of nonlinear equations, we may potentially run into this problem due to various approximations—the solution would follow the direction EA but would lie outside the circle and, as

a result, the scheme may encounter stability problems.

**1.3. First-order centered scheme (L-F).** Another way to stabilize the Euler forward scheme is by employing a dissipative start-off value: replacing $u_j$ in (1.4) by the average of the two neighboring values. The result is the Lax-Friedrichs (L-F) scheme (Lax 1954):

$$u_j^{n+1} = \frac{1}{2}(u_{j-1} + u_{j+1}) - \sigma\frac{1}{2}(u_{j+1} - u_{j-1}). \quad (1.8)$$

Note that when $\Delta t = 0$, the solution, instead of being $u_j$, is $\frac{1}{2}(u_{j-1} + u_{j+1})$. Such a quantity is called a *start-off* value here.

The amplification factor for the L-F scheme is

$$\mathcal{A} = \frac{1}{2}(e^{Iw} + e^{-Iw}) - \sigma\frac{1}{2}(e^{Iw} - e^{-Iw}). \quad (1.9)$$

If $a \geq 0$, then $0 \leq \sigma \leq 1$, and the amplification factor lies on the line segment DC shown in Fig. 1.1. As a result, the scheme is stable. Observe that the two key advantages of the L-F scheme are simplicity and stability (the solution lies well inside the unit circle).

Concerning the disadvantages, by (1.8), the solution at $j$ does not depend on the data $u_j$, but depends only on the data at the two neighboring indices. Consequently, the scheme has the odd-even decoupling problem. This problem can also be seen via the staggered-mesh formulation. Indeed, suppose at time level $n$, we have data at even indices only. At time level $n + 1$, (1.8) yields the solution at all odd indices. At time level $n + 2$, we can obtain the solution at all even indices again, etc. Such a scheme employs a staggered mesh. The other staggered-mesh solution has data at odd indices at time level $n$. Thus, the regular-mesh version (1.8) above consists of two completely independent staggered-mesh solutions. As a consequence, the L-F scheme needs twice the number of mesh points to have the same resolution as the first-order upwind scheme. (This observation can also be seen from Fig. 1.1.) In other words, the convenience of a regular mesh is achieved at the cost of a) twice the amount of calculations and b) the odd-even decoupling problem.

The next two observations hold for the L-F as well as the second-order L-F type schemes in §3. First, in an unsteady calculation, if we wish to animate the solution, we should use the solutions at even time steps only (or odd ones only). Similarly, to check for convergence to a steady state, one must compare the solution at time level $n$ with that at time level $n - 2$. Second, when $\sigma = 0$, the L-F

scheme has no phase error while the damping error reaches a maximum. Consequently, in practical calculations, it is desirable to employ a time step as large as possible to minimize damping.

**1.4. Scheme blending Lax-Friedrichs and Euler forward.** The two drawbacks of the L-F scheme (odd-even decoupling and too much numerical dissipation) can be dealt with by blending it with the Euler forward method. In other words, in Fig. 1.1, instead of starting off at the point D, we can employ a start-off value which blends $D$ and $A$: the quantity $\frac{1}{2}(u_{j-1} + u_{j+1})$ in (1.8) is replaced by $(1 - |\sigma|)u_j + |\sigma|\frac{1}{2}(u_{j-1} + u_{j+1})$. The resulting scheme is

$$u_j^{n+1} = [(1 - |\sigma|)u_j + |\sigma|\frac{1}{2}(u_{j-1} + u_{j+1})]$$
$$- \sigma\frac{1}{2}(u_{j+1} - u_{j-1}). \qquad (1.10)$$

The above centered scheme turns out to be identical to the first-order upwind scheme. For the case of second-order accuracy, however, the blended scheme is different from the second-order upwind one.

Note that in extending (1.10) to the Euler equations, the quantity $|\sigma|$ takes the form $|A|\Delta t/h$ where $A$ is the Jacobian matrix. For simplicity, $|A|$ can be replaced by the sum of the magnitude of the local flow speed and the speed of sound.

**2. Upwind schemes via MUSCL approach.** The above schemes were derived by finite differencing. To ensure that shocks are captured, fluxes need to be balanced via the integral form of the equation (Lax 1954). Oscillations near shocks can be suppressed by using linear functions (or polynomials) to approximate the data in each cell and by limiting the slopes of these linear functions so that they are not too steep near a discontinuity (Van Leer 1977).

Integrating (1.1a) on the interval $[\alpha, \beta]$, one obtains

$$\frac{\partial}{\partial t}\int_\alpha^\beta u(x,t)\,dx + au(\beta,t) - au(\alpha,t) = 0, \quad (2.1)$$

where $au(\alpha, t)$ is the flux (or rate of flow) of $u$ across interface $\alpha$ at time $t$.

Next, let $x_{j+1/2} = (j + 1/2)h$ be the cell interfaces and $x_j = jh$ the cell centers of a uniform mesh, $j = 0, 1, 2, \ldots$ For each cell $[x_{j-1/2}, x_{j+1/2}]$, second-order accuracy is obtained by applying the midpoint rule to (2.1):

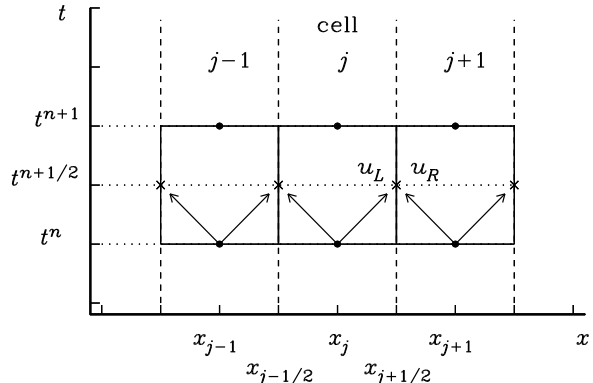$$u_j^{n+1} = u_j + \frac{\Delta t}{h}(au_{j-1/2}^{n+1/2} - au_{j+1/2}^{n+1/2}); \qquad (2.2)$$



Fig. 2.1. *Interface fluxes for second-order accurate upwind schemes. At each interface $j + 1/2$ time level $n + 1/2$, the reconstruction $r_j$ to the left of the interface yields $u_L$ while $r_{j+1}$ yields $u_R$. The upwind value is employed for the flux calculation.*

here, $u_j$ (i.e., $u_j^n$) approximates the average value of $u$ in the cell $j$ at time $t^n$, and $u_{j+1/2}^{n+1/2}$ approximates the value at the interface $j + 1/2$ at time level $n + 1/2$. Assume that we know the data $u_j$ (i.e., $u_j^n$) for all $j$; we wish to calculate the solution $u_j^{n+1}$.

**2.1. First-order upwind scheme.** For the first-order case, in each cell $j$, the data are approximated by a constant function: $r_j(x) = u_j$. Next, applying the FTCS (forward-time centered-space) approximation to (2.1), one obtains

$$u_j^{n+1} = u_j + \frac{\Delta t}{h}(au_{j-1/2} - au_{j+1/2}). \qquad (2.3)$$

At each interface $j + 1/2$, the data to the left is the constant function $u_j$, that to the right, $u_{j+1}$. The flux is given by the upwind choice

$$au_{j+1/2} = \begin{cases} au_j & \text{if } a \geq 0, \\ au_{j+1} & \text{otherwise.} \end{cases} \qquad (2.4)$$

For consistency with the case of the Euler equations later, at each interface, $a$ and $u$ are combined as above.

**2.2. Second-order upwind scheme.** For second-order accuracy, the data are approximated by a linear function in each cell. To calculate the flux $au_{j+1/2}^{n+1/2}$ in (2.2), suppose $\{(u_x)_j\}$ are known. They can be approximated by, e.g., the central difference (also called the average slope)

$$(u_x)_j = \frac{1}{2h}(u_{j+1} - u_{j-1}). \qquad (2.5)$$

(A weighted average for $(u_x)_j$ will be discussed later.) The time derivative $(u_t)_j$ follows from the

5

advection equation,

$$(u_t)_j = -a(u_x)_j. \qquad (2.6)$$

In the domain $[x_{j-1/2}, x_{j+1/2}] \times [t^n, t^{n+1}]$, the solution can then be approximated by a linear function $r_j$ called the *reconstruction*,

$$r_j(x,t) = u_j + (u_x)_j(x - x_j) + (u_t)_j(t - t^n). \quad (2.7)$$

Note that the term 'reconstruction' is typically used for $r_j(x, t^n)$. While $r_j(x, t)$ was employed by Harten et al. (1987) and Nessyahu and Tadmor (1990), it was not given a name. Here, we use the term 'reconstruction' in the extended sense above.

At the two interfaces of cell $j$, we can march to the half time step using $r_j$. At each interface $x_{j+1/2}$ and time $t^{n+1/2}$, we now have two values. The value from cell $j$ is denoted by $u_L$; that from $j + 1$, by $u_R$, as shown in Fig. 2.1,

$$u_L = u_j + \tfrac{h}{2}(u_x)_j + \tfrac{\Delta t}{2}(u_t)_j,$$

$$u_R = u_{j+1} - \tfrac{h}{2}(u_x)_{j+1} + \tfrac{\Delta t}{2}(u_t)_{j+1}.$$

The upwind choice for the flux is

$$a u_{j+1/2}^{n+1/2} = \begin{cases} a u_L & \text{if } a \ge 0, \\ a u_R & \text{otherwise.} \end{cases} \qquad (2.8)$$

This flux and (2.2) complete the description of the second-order upwind scheme.

The calculation (2.6) of the time partial derivative from the spatial one follows Hancock's observation (Van Albada et al. 1982). It simplifies the extension to systems of equations. Also notice that the Taylor series $r_j(x, t)$ yields a result identical to that by the method of characteristics (1.2) applied to the linear initial condition $r_j(x, t^n)$.

Note that if $a = 0$, then there is ambiguity in defining the interface value $u_{j+1/2}^{n+1/2}$ in (2.8), but since the flux is zero, this ambiguity causes no problem. In fact, the flux can be expressed without a conditional statement:

$$a u_{j+1/2}^{n+1/2} = \tfrac{1}{2}(a u_L + a u_R) - \tfrac{1}{2}|a|(u_R - u_L). \quad (2.9)$$

Similarly, for systems of equations, there is no ambiguity in splitting the flux, and the technique works well. Splitting the variables, however, may result in unwanted oscillations. Also note that the term $\tfrac{1}{2}|a|(u_R - u_L)$ produces the upwind-biased effect and thus numerical dissipation.

If $a \ge 0$, expressions (2.2) and (2.5–2.8) yield the following solution

$$u_j^{n+1} = u_j + \sigma(u_{j-1} - u_j) +$$
$$\tfrac{1}{4}\sigma(1-\sigma)[u_j + u_{j-1} - u_{j-2} - u_{j+1}].$$

This scheme was formulated via finite differencing by Fromm (1968). The above reconstruction formulation follows that of (Van Leer 1977) except that the method of characteristics there is replaced by the Taylor series (2.7) here. The scheme is the first and also the least accurate among a series of five schemes in the cited reference.

Fourier analysis yields the following amplification factor for the second-order upwind scheme,

$$\mathcal{A} = 1 + \sigma(e^{-Iw} - 1) +$$
$$\tfrac{1}{4}\sigma(1-\sigma)(1 + e^{-Iw} - e^{-2Iw} - e^{Iw}). \quad (2.10)$$

To avoid oscillations near a discontinuity, instead of the average (2.5), a weighted average (Van Albada et al. 1982) can be employed for the slope. For any two real variables $x$ and $y$, define

$$\text{wtav}(x,y) = \frac{x^2 y + y^2 x}{x^2 + y^2 + 10^{-20}}. \qquad (2.11)$$

If $x/y \approx 1$, the above expression yields essentially the average $\tfrac{1}{2}(x + y)$; on the other hand, if $x/y \approx 0$ or $\infty$, the result is essentially the value with smaller modulus. The slope formula takes the form

$$(u_x)_j = \tfrac{1}{h}\text{wtav}(u_{j+1} - u_j, u_j - u_{j-1}). \quad (2.12)$$

Note that between two slopes $\tfrac{1}{h}(u_j - u_{j-1})$ and $\tfrac{1}{h}(u_{j+1} - u_j)$, the above weighted average produces a result which is biased toward the less steep slope. Since such a result is closer to the zero slope of the first-order upwind scheme than the result by the average (2.5), the above weighted average adds a considerable amount of numerical dissipation. In fact, it yields a scheme which is only first-order accurate near an extremum—a well-known drawback of all standard weighted averages or limiter functions.

**2.3. Scheme II of Van Leer.** Next, we describe an upwind scheme which carries along (stores and updates) the interface values. The scheme is the second one presented in (Van Leer 1977) except, again, the method of characteristics employed there is replaced by the Taylor series (2.7) here. While this upwind scheme does not extend easily to the case of the Euler equations (more on this later), its centered counterpart, which turns out to be the CE/SE scheme with $\epsilon = 1/2$, does.

Suppose the data consist of the cell averages $\{u_j\}$ as well as the interface values $\{u_{j+1/2}\}$. We wish to calculate $\{u_j^{n+1}\}$ and $\{u_{j+1/2}^{n+1}\}$. First, instead of the average (2.5), the slope is defined via the interface values:

$$(u_x)_j = \tfrac{1}{h}(u_{j+1/2} - u_{j-1/2}). \qquad (2.13)$$

With the above slope formula, the cell averages $\{u_j^{n+1}\}$ are calculated as before, by (2.2) and (2.6)–(2.8). The interface values are updated by using the reconstruction in the upwind cell:

$$u_{j+1/2}^{n+1} = \begin{cases} r_j(x_{j+1/2}, t^{n+1}) & \text{if } a \geq 0, \\ r_{j+1}(x_{j+1/2}, t^{n+1}) & \text{otherwise.} \end{cases} \quad (2.14)$$

This completes the description of the scheme.

The above scheme has the same leading (phase) error compared with the scheme using (2.5). The dissipation error, however, is reduced by two thirds. Thus, carrying along the interface values improves the dissipation error considerably.

Note that extending (2.14) to the Euler equations appears to be difficult: when $a = 0$, there is no ambiguity in defining the flux, but the choice of the interface value becomes ambiguous.

### 3. Centered schemes via MUSCL appoach.
Each upwind scheme in §2 has a centered counterpart in this section. For centered schemes, the spatial region where the reconstruction is valid at time $t = t^n$, which is called the reconstruction cell here, is twice the size of the original cell. Together, these reconstruction cells cover the domain twice whereas the original cells cover the domain once only.

More precisely, on the same mesh of cells $[x_{j-1/2}, x_{j+1/2}]$ as in the upwind case, let the *reconstruction cell* $j$ be the interval $[x_{j-1}, x_{j+1}]$ of length $2h$ defined by the centroids of the two neighboring cells. Assume that we know the data $u_j$ (i.e., $u_j^n$) for all $j$; $u_j$ approximates the average value of $u$ on the reconstruction cell $j$. We wish to calculate the solution $u_j^{n+1}$. See Fig. 3.1.

### 3.1. First-order centered scheme (L-F).
Here, we rederive the L-F scheme (1.8) from the MUSCL perspective. In each reconstruction cell $j$, the data are approximated by a constant function: $r_j(x) = u_j$. Next, when updating the solution at $j$, *the reconstruction $r_j$ is ignored*, while $r_{j-1}$ and $r_{j+1}$ are employed. See Fig. 3.1. Thus, the start-off (average) value is

$$u_j^* = \tfrac{1}{2}(u_{j-1} + u_{j+1}). \quad (3.1)$$

Applying the FTCS approximation to (2.1) on the reconstruction cell $j$, one obtains

$$u_j^{n+1} = \tfrac{1}{2}(u_{j-1} + u_{j+1}) + \tfrac{\Delta t}{2h}(au_{j-1} - au_{j+1}). \quad (3.2)$$

Here, there is no need of upwinding because the flux $au_{j-1}$ is evaluated at the center of the reconstruction cell $j-1$, and the flux $au_{j+1}$, the center of the reconstruction cell $j+1$.
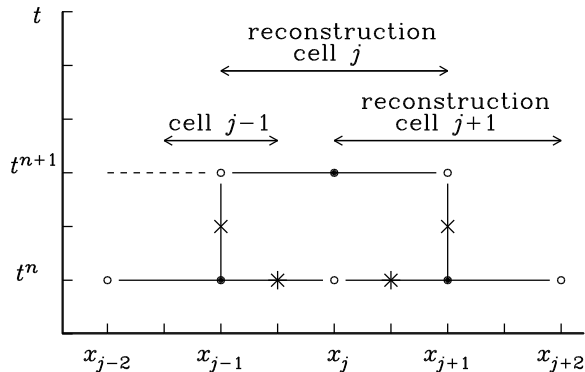


Fig. 3.1. *Centered schemes. The interval $[x_{j-1}, x_{j+1}]$ is called the reconstruction cell $j$ (the spatial region where the reconstruction is valid at time $t = t^n$). The two adjacent reconstruction cells $j$ and $j+1$ overlap on $[x_j, x_{j+1}]$. For each $j$, the reconstruction $r_j(x,t)$ is valid on the triangle defined by the three corners $(x_{j-1}, t^n)$, $(x_{j+1}, t^n)$, $(x_j, t^{n+1})$. When calculating $u_j^{n+1}$, the reconstruction $r_j$ is ignored, while $r_{j-1}$ and $r_{j+1}$ are employed. For second-order accuracy, the values at the points marked by $\times$ and $*$ are needed. They are $u_{j-1} + \frac{\Delta t}{2}(u_t)_{j-1}$, $u_{j+1} + \frac{\Delta t}{2}(u_t)_{j+1}$, and $u_{j-1} + \frac{h}{2}(u_x)_{j-1}$, and $u_{j+1} - \frac{h}{2}(u_x)_{j+1}$.*

The next remark concerns the initial values for the L-F scheme. The remark also holds for the N-T and the CE/SE schemes below. Suppose the initial condition is a step function, say, $u_0(x) = 1$ for $x \leq 0$, and $u_0(x) = 0$ otherwise. Suppose the cell centers are $x_j = j$, $j = -N, \ldots, N$ where $N = 50$. Then the solution at a later time $t^n$ in standard textbooks exhibits a staircase pattern: $u_1 = u_2$, and $u_3 = u_4$, etc. One way to avoid this problem is to define the initial data as the average of $u_0(x)$ on the reconstruction cell $j$, i.e., $u_0^0 = 1/2$.

### 3.2. Second-order centered scheme (N-T).
This scheme is the centered counterpart of the second-order accurate upwind scheme of §2.2. It was presented by Nessyahu and Tadmor (1990) using limiter functions of minmod type to define the slope and was named ORD. For consistency with the upwind scheme here, Van Albada's weighted average is employed. In addition, the average slope (2.5) makes Fourier analysis possible.

As in the upwind case, the slopes are given by

$$(u_x)_j = \tfrac{1}{2h}(u_{j+1} - u_{j-1}), \quad (3.3)$$

$$(u_t)_j = -a(u_x)_j. \quad (3.4)$$

At time $t^n$, on the reconstruction cell $j$, the data is

7

assumed to be the linear function $u_j + (u_x)_j (x - x_j)$. Since information propagates a distance no more than $h$ per time step, information from outside the reconstruction cell has not reached its center *regardless* of the wind direction. As a result, the above $u_t$ is valid in the triangle defined by the three corners $(x_{j-1}, t^n)$, $(x_{j+1}, t^n)$, $(x_j, t^{n+1})$. In this region, the solution can be approximated by, as in (2.7),

$$r_j(x,t) = u_j + (u_x)_j(x - x_j) + (u_t)_j(t - t^n). \quad (3.5)$$

When updating the solution at $j$, again, *the reconstruction $r_j$ is ignored*, while $r_{j-1}$ and $r_{j+1}$ are employed. The start-off (average) value is

$$u_j^* = \tfrac{1}{2}[u_{j-1} + \tfrac{h}{2}(u_x)_{j-1} + u_{j+1} - \tfrac{h}{2}(u_x)_{j+1}]. \quad (3.6)$$

See Fig. 3.1. The midpoint rule for (2.1) on the reconstruction cell $j$ yields

$$u_j^{n+1} = u_j^* + \tfrac{\Delta t}{2h}(au_{j-1}^{n+1/2} - au_{j+1}^{n+1/2}). \quad (3.7)$$

Note that the use of $u_j^*$ instead of $u_j$ results in numerical dissipation which stabilizes the scheme. The values $u_{j-1}^{n+1/2}$ and $u_{j+1}^{n+1/2}$ are evaluated by the reconstruction function: for any $j$,

$$u_j^{n+1/2} = u_j + \tfrac{\Delta t}{2}(u_t)_j. \quad (3.8)$$

The above completes the description of the N-T scheme.

With expression (3.3) for the slope, the result is,

$$\begin{aligned}
u_j^{n+1} = & \tfrac{1}{2}\big[u_{j-1} + \tfrac{1}{4}(u_j - u_{j-2}) + \\
& u_{j+1} - \tfrac{1}{4}(u_{j+2} - u_j)\big] + \\
& \tfrac{1}{2}\sigma\big\{[u_{j-1} - \tfrac{1}{4}\sigma(u_j - u_{j-2})] \\
& [u_{j+1} - \tfrac{1}{4}\sigma(u_{j+2} - u_j)]\big\}.
\end{aligned} \quad (3.9)$$

At first glance, since the right hand side above involves $u_j$, the N-T scheme appears to be (odd-even) coupled. This coupling, however, is not a strong coupling, as can be seen in the following example. Suppose the data at odd indices are 1 and those at even, $-1$. Then the solution by the N-T scheme after two time steps remains identical to the initial data for all time step sizes. It is desirable for a scheme to damp out this odd-even data. In spite of this observation, based on the author's experience, the N-T scheme produces solutions with no odd-even noise in practice provided that the initial condition is appropriate as discussed in the last paragraph of §3.1.

The amplification factor of the N-T scheme is given by, after some algebra,

$$\mathcal{A} = \cos(w) - I\sigma\sin(w) + \tfrac{1}{2}(1 - \sigma^2)\sin^2(w). \quad (3.10)$$

The next remark concerns the staggered-mesh version. At time level $n$, suppose we have data at even indices only, and we wish to calculate the solution at odd indices. Then, the slope at an even index $2j$ is $(u_{2j+2} - u_{2j-2})/(4h)$. For the nonstaggered version above, the slope at index $2j$ is $(u_{2j+1} - u_{2j-1})/(2h)$. Because the latter slope uses values at locations closer to $j$, the numerical dissipation of the nonstaggered version is only $1/3$ that of the staggered version. The phase errors, however, are essentially the same (for small wave numbers). Thus, the nonstaggered version costs twice as much as the staggered one, but its numerical dissipation improves by two thirds.

As in the first-order case, to reduce numerical dissipation for the nonstaggered N-T scheme when $\sigma$ is small, we can use a blended start-off value: with $u_j^*$ defined by (3.6)

$$\begin{aligned}
u_j^{n+1} = & \{(1 - |\sigma|)u_j + |\sigma| u_j^*\} + \\
& \tfrac{1}{2h}\Delta t(au_{j-1}^{n+1/2} - au_{j+1}^{n+1/2}).
\end{aligned} \quad (3.11)$$

The resulting scheme also damps the data of 1 at odd and $-1$ at even indices.

**3.3. CE/SE scheme.** The centered counterpart of Van Leer's second scheme in §2.3 turns out to be Chang's CE/SE scheme (the $\epsilon = 1/2$ member). To show this fact, we describe the CE/SE scheme as one which carries along (stores and updates) the slopes, and then as one which carries along the interface values.

At time $t^n$, suppose the data $\{u_j\}$ as well as the slopes $\{(u_x)_j\}$ are known and stored. We wish to calculate $\{u_j^{n+1}\}$ and $\{(u_x)_j^{n+1}\}$. The cell average quantities $\{u_j^{n+1}\}$ are updated exactly as in the N-T case: by (3.4)–(3.8). To update the slope, first, for each $j$, the point value (as opposed to cell average) $\widehat{u}_j^{n+1}$ is calculated by the reconstruction:

$$\widehat{u}_j^{n+1} = r_j(x_j, t^{n+1}) = u_j + (u_t)_j\Delta t. \quad (3.12)$$

The slope is updated by the two neighboring point values:

$$(u_x)_j^{n+1} = \tfrac{1}{2h}(\widehat{u}_{j+1}^{n+1} - \widehat{u}_{j-1}^{n+1}). \quad (3.13)$$

If the weighted average is employed, the cell average $u_j^{n+1}$ is also needed:

$$(u_x)_j^{n+1} = \tfrac{1}{h}\,\text{wtav}(\widehat{u}_{j+1}^{n+1} - u_j^{n+1}, u_j^{n+1} - \widehat{u}_{j-1}^{n+1}).$$

Clearly, instead of storing the slopes $\{(u_x)_j^{n+1}\}$, we could store the interface values $\{\widehat{u}_j^{n+1}\}$. Such a

scheme takes the following form: at time $t^n$, suppose the data $\{u_j\}$ as well as the point values $\{\widehat{u}_j\}$ are known and stored. To update the solution, first, the slope is given by

$$(u_x)_j = \tfrac{1}{h}\,\mathrm{wtav}(\widehat{u}_{j+1} - u_j, u_j - \widehat{u}_{j-1}). \quad (3.14)$$

The cell average $u_j^{n+1}$ is then updated by again (3.4)–(3.8). The point value $\widehat{u}_j^{n+1}$ is updated by (3.12).

Note that the CE/SE scheme is odd-even decoupled, as can be seen by the following staggered-mesh formulation. At time level $n$, suppose we have data $u_{2j}$ at even indices and $\widehat{u}_{2j+1}$ at odd indices only. Here, $\{\widehat{u}_{2j+1}\}$ are the values at the interfaces of the cells $[x_{2j-1}, x_{2j+1}]$. At time level $n+1$, we calculate the solution $u_{2j+1}^{n+1}$ and $\widehat{u}_{2j}^{n+1}$. The other staggered-mesh solution is obvious. Thus, the nonstaggered version is a result of overlaying two independent staggered-mesh solutions.

The above observation also shows that for 1D advection, in the form that carries along the interface point values, the CE/SE scheme is the centered counterpart of Van Leer's second scheme described in §2.3.

Notice that the CE/SE scheme in above form needs less storage: for the 3D case, instead of storing $u$, $u_x$, $u_y$, and $u_z$, we only need to store $u$ and (the point value) $\widehat{u}$. If we wish to adjust numerical dissipation ($\epsilon \neq 1/2$), however, we have to store the slopes.

A few remarks concerning the CE/SE and the N-T schemes are in order. First, if the slopes are carried along, the CE/SE scheme has a very compact stencil: when updating the solution at cell $j$, only the data at the immediate neighbors $j-1$ and $j+1$ are employed. But the slopes at $j-1$ and $j+1$ use the point values at $j-2$ and $j+2$. Thus, in the form of carrying along the interface point values, the stencil of the CE/SE scheme is the same as that of the N-T method.

Next, after one time step, the cell average update $u_j^{n+1}$ has an error of $O(h^3)$ (the solution is exact if the data are on a parabola). The point value update $\widehat{u}_j^{n+1}$, by (3.12), has an error of $O(h^2)$ (exact when the data are linear). Therefore, when calculating the slopes $(u_x)_j^{n+1}$, the cell averages $u_{j-1}^{n+1}$ and $u_{j+1}^{n+1}$ appear to be a better choice than the points values $\widehat{u}_{j-1}^{n+1}$ and $\widehat{u}_{j+1}^{n+1}$. It turns out, however, that the point values lead to a scheme with the same phase error but less dissipation error (more details later), i.e., the CE/SE scheme ($\epsilon = 1/2$) is less dissipative than the N-T scheme.

For Fourier stability analysis, set

$$\mathbf{u}_j = \begin{pmatrix} u_j \\ (u_x)_j \end{pmatrix}. \quad (3.16)$$

Then, the solution vector of the CE/SE scheme can be written as

$$\mathbf{u}_j^{n+1} = \mathbf{B}\,\mathbf{u}_{j-1} + \mathbf{C}\,\mathbf{u}_{j+1}, \quad (3.17)$$

where

$$\mathbf{B} = \begin{pmatrix} \tfrac{1}{2}(1+\sigma) & \tfrac{1}{4}(1-\sigma)(1+\sigma) \\ \tfrac{1}{2} & \tfrac{\sigma}{2} \end{pmatrix}, \quad (3.18)$$

and

$$\mathbf{C} = \begin{pmatrix} \tfrac{1}{2}(1-\sigma) & -\tfrac{1}{4}(1-\sigma)(1+\sigma) \\ -\tfrac{1}{2} & -\tfrac{\sigma}{2} \end{pmatrix}. \quad (3.19)$$

The two eigenvalues of the matrix $e^{-Iw}\mathbf{B} + e^{Iw}\mathbf{C}$ are the amplification factors of the CE/SE scheme. After some algebra,

$$\begin{aligned} \mathcal{A}^{\pm} &= \tfrac{1}{2}\cos(w) - I\sigma\sin(w) \pm \\ &\quad \tfrac{1}{2}\sqrt{1 + (1 - 2\sigma^2)\sin^2(w)}. \end{aligned} \quad (3.20)$$

Here $\mathcal{A}^+$, which approximates $e^{-I\sigma w}$ and determines the accuracy of the scheme, is the principal amplification factor, while $\mathcal{A}^-$ is the spurious one.

Our next remarks concern the start-off slopes and the reversible CE/SE scheme. This scheme involves looking both forward and backward in time (Chang 1995). It is described below using a start-off slope and a forward only time evolution. For the case of the Euler equations, such a forward marching scheme provides an approximation to the reversible scheme.

At time $t^n$, suppose $\{u_j\}$ and $\{(u_x)_j\}$ are known. Then if $\Delta t = 0$, the slope $(u_x)_j^{n+1}$ of the above CE/SE scheme ($\epsilon = 1/2$) reduces to the start-off slope

$$(u_x)_j^* = \tfrac{1}{2h}(u_{j+1} - u_{j-1}). \quad (3.21)$$

For $\Delta t \geq 0$, the slope update (again, $\epsilon = 1/2$) is given by

$$(u_x)_j^{n+1} = (u_x)_j^* + \tfrac{\Delta t}{2h}[(u_t)_{j+1} - (u_t)_{j-1}]. \quad (3.22)$$

The reversible CE/SE scheme ($\epsilon = 0$) can be described as a scheme with a different start-off slope. Instead of $(u_x)_j^*$, the start-off slope is defined by the values employed in the start-off average (3.6):

$$[(u_x)_j]_R^* = \tfrac{1}{h}\left\{ u_{j+1} - \tfrac{h}{2}(u_x)_{j+1} - [u_{j-1} + \tfrac{h}{2}(u_x)_{j-1}] \right\}. \quad (3.23)$$
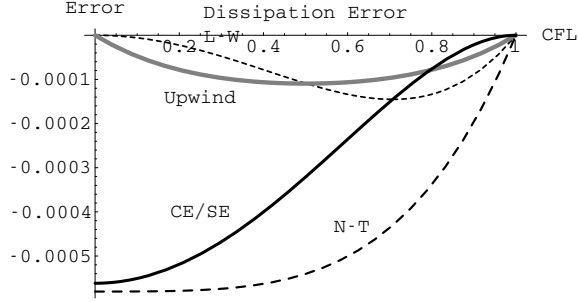
Fig. 3.2. *Dissipation errors (per time step) as functions of $\sigma$ for $w = \pi/12$ (relatively low frequency). For each scheme, this error is proportional to $w^4$. When $\sigma = 0$, the N-T and CE/SE schemes have a maximum amount of numerical dissipation. For reference, the error of the Lax-Wendroff scheme is also plotted. While the Lax-Wendroff scheme has relatively small errors, it causes oscillations near shocks.*

The slope update for the reversible scheme is

$$(u_x)_j^{n+1} = [(u_x)_j]_R^* + \frac{\Delta t}{2h}[(u_t)_{j+1} - (u_t)_{j-1}]. \quad (3.24)$$

To obtain the general CE/SE schemes, we simply blend the two start-off slopes (3.21) and (3.23) using a parameter $\epsilon$ by $[(u_x)_j]_R^* + 2\epsilon\{(u_x)_j^* - [(u_x)_j]_R^*\}$. The evolution for the slope, i.e., the quantity $\frac{\Delta t}{2h}[(u_t)_{j+1} - (u_t)_{j-1}]$, is the same for all $\epsilon$. As such, the CE/SE schemes have an *adjustable start-off slope*.

It is remarkable that the start-off average $u_j^*$ via (3.6) produces numerical dissipation, whereas the start-off slope can have the opposite effect: it can reduce dissipation. Note that one can derive the upwind counterparts of the reversible and the general CE/SE schemes ($0 \le \epsilon \le 1$) as well. These schemes, however, are beyond the scope of this paper.

Also note that the reconstruction cell here is identical to the spatial part (or spatial projection) of the solution element in Chang (1995), and the control volume on which fluxes are balanced here, the conservation element there.

**3.4. Accuracy comparison.** The upwind and centered schemes described above are stable for $|\sigma| \le 1$. To compare accuracy, recall that the amplification factor $\mathcal{A}$ approximates the exact amplification factor $e^{-I\sigma w}$. The phase error per time step is $\text{Arg}(\mathcal{A}) - (-\sigma w)$. For the second-order schemes discussed here, this quantity is proportional to $w^3$ when the wave number $w$ is small. The dissipation



Fig. 3.3. *Phase errors (per time step) as functions of $\sigma$ for $w = \pi/12$. For each scheme, this error is proportional to $w^3$ and is the leading error. Note that the phase errors of the two centered schemes are nearly identical. (They are identical on a similar plot for $w = \pi/32$.)*

error per time step is $|\mathcal{A}| - 1$, which is proportional to $w^4$.

Figure 3.2 shows the dissipation errors as functions of $\sigma$ for $w = \pi/12$ (relatively smooth data). For ease of reference, the error of the Lax-Wendroff scheme is also plotted.

Figure 3.3 shows phase errors as functions of $\sigma$ for again, $w = \pi/12$. The phase errors of the centered schemes are nearly identical. Notice that when $\sigma = 1/2$, the upwind scheme has no phase error, while its dissipation error reaches a maximum.

Observe that for 1D advection, the upwind scheme is considerable more accurate than the non-staggered N-T and CE/SE schemes. The reason is that the cells of the upwind scheme are half the size of the reconstruction cells of the two centered schemes.

Between the N-T and the CE/SE schemes, Fig. 3.2 shows that the latter has less numerical dissipation. For a small CFL number, however, both schemes have about the same amount of dissipation. In practice, the weighted average adds additional dissipation and the results by these two schemes are nearly the same, as will be shown.

**4. Schemes for the Euler equations.** The second-order upwind, N-T, and CE/SE schemes are extended to the case of the Euler equations in this section. For this case, the N-T and CE/SE schemes remain simple; in fact, the equations are identical to those in the case of advection except that the symbols are boldfaced. As for the upwind scheme, we need additional techniques, namely, Roe's splitting with an entropy correction. Roe's method consists of a diagonalization and an upwind side selection

in the characteristic frame. The entropy correction serves the purpose of excluding non-physical solutions and avoiding kinks (glitches) in the solution near sonic points. The entropy correction in Huynh (1995a), which amounts to adding dissipation when the wave speeds spread past zero, is employed. While these concepts for the upwind step are somewhat involved, the coding remains concise and economical. The derivation for this upwind step is carried out below to show the complexity— or simplicity—of upwinding. Also note that the presentation of Roe's scheme here is simpler than most presentations in the literature.

The one-dimensional flow of an inviscid and compressible gas obeys the conservation laws for mass, momentum, and energy:

$$\mathbf{U}_t + \mathbf{F}(\mathbf{U})_x = 0, \qquad (4.1)$$

$$\mathbf{U} = \begin{pmatrix} \rho \\ \rho u \\ e \end{pmatrix}, \qquad \mathbf{F} = \begin{pmatrix} \rho u \\ \rho u^2 + p \\ (e+p)u \end{pmatrix},$$

where $t$ is time, $x$ distance, $\rho$ density, $e$ total energy per unit volume, $u$ velocity, and $p$ pressure. Let $\gamma$ be the ratio of specific heats. Then for a perfect gas,

$$p = (\gamma - 1)(e - \tfrac{1}{2}\rho u^2). \qquad (4.2)$$

Integrating (4.1) on the interval $[\alpha, \beta]$, one obtains the integral form

$$\frac{\partial}{\partial t}\int_\alpha^\beta \mathbf{U}(x,t)\,dx + \mathbf{F}(\mathbf{U}(\beta,t)) - \mathbf{F}(\mathbf{U}(\alpha,t)) = 0, \qquad (4.3)$$

where $\mathbf{F}(\mathbf{U}(\alpha,t))$ is the flux across interface $\alpha$ at time $t$.

In regions where $\mathbf{U}$ is smooth, (4.1) is equivalent to the non-conservation form

$$\mathbf{U}_t + \mathbf{A}_c \mathbf{U}_x = 0, \qquad (4.4)$$

where the Jacobian matrix $\mathbf{A}_c$ is

$$\mathbf{A}_c = \frac{\partial \mathbf{F}}{\partial \mathbf{U}}. \qquad (4.5)$$

With $\mathbf{F}^{(k)}$ denoting the $k$-th component of $\mathbf{F}$,

$$(\mathbf{A}_c)_{k,l} = \frac{\partial \mathbf{F}^{(k)}}{\partial \mathbf{U}^{(l)}}.$$

By rewriting $\mathbf{F}$ in terms of $\rho$, $m$ ($m = \rho u$), and $e$, after some algebra,

$$\mathbf{A}_c = \begin{pmatrix} 0 & 1 & 0 \\ (\gamma-3)u^2/2 & (3-\gamma)u & \gamma-1 \\ (\gamma-1)u^3 - \gamma u e/\rho & a_{32} & \gamma u \end{pmatrix}, \qquad (4.6)$$

where

$$a_{32} = -3(\gamma-1)u^2/2 + \gamma e/\rho.$$

For the advection equation, the speed is $a$; here, there are three wave speeds; they can be calculated by diagonalizing $\mathbf{A}_c$.

Since the centered schemes are very simple, they are described first. Also note that we consider only the interior points; boundary conditions are standard, and are omitted.

**4.1. The N-T and CE/SE schemes.** Let $x_{j+1/2} = (j + 1/2)h$ be the cell interfaces and $x_j = jh$ the cell centers, $j = 0, 1, 2, \ldots$ Recall that the reconstruction cell $j$ is the interval $[x_{j-1}, x_{j+1}]$ of length $2h$.

At time $t^n$, assume that we know $\mathbf{U}_j$ (i.e. $\mathbf{U}_j^n$) for all $j$; $\mathbf{U}_j$ appoximates the average of $\mathbf{U}$ on the reconstruction cell $j$. We wish to calculate $\mathbf{U}_j^{n+1}$. In the CE/SE case, the slopes $\{(\mathbf{U}_x)_j\}$ are also stored, and $\{(\mathbf{U}_x)_j^{n+1}\}$ must also be calculated.

For stability, the time step $\Delta t$ is required to satisfy the CFL condition

$$\left[\max_j(|u_j| + a_j)\right]\frac{\Delta t}{h} \leq 1, \qquad (4.7)$$

where $a$ denotes the speed of sound: $a = (\gamma p/\rho)^{1/2}$. Notice that the time step $\Delta t$ here corresponds to the quantity $\Delta t/2$ in Chang (1995) and Wang and Chang (1999).

To update the cell average $\mathbf{U}_j^{n+1}$, the data $\mathbf{U}_j$ is ignored. The midpoint rule for (4.3) on the reconstruction cell $[x_{j-1}, x_{j+1}]$ takes the form

$$\mathbf{U}_j^{n+1} = \mathbf{U}_j^* + \frac{\Delta t}{2h}(\mathbf{F}_{j-1}^{n+1/2} - \mathbf{F}_{j+1}^{n+1/2}), \qquad (4.8)$$

where the start-off value $\mathbf{U}_j^*$ takes the place of $\mathbf{U}_j$ and will be defined below.

Next, we describe the reconstruction function. In the CE/SE case, the stored slope is employed; for the N-T case, the slope is estimated by Van Albada's weighted average:

$$(\mathbf{U}_x)_j = \frac{1}{h}\mathrm{wtav}(\mathbf{U}_{j+1} - \mathbf{U}_j, \mathbf{U}_j - \mathbf{U}_{j-1}). \quad (4.9)$$

The time derivative is given by (4.4):

$$(\mathbf{U}_t)_j = -(\mathbf{A}_c)_j (\mathbf{U}_x)_j. \qquad (4.10)$$

The reconstruction function takes the form

$$\mathbf{r}_j(x,t) = \mathbf{U}_j + (\mathbf{U}_x)_j(x - x_j) + (\mathbf{U}_t)_j(t - t^n). \quad (4.11)$$

When updating the solution at $j$, $\mathbf{r}_j$ is ignored while $\mathbf{r}_{j-1}$ and $\mathbf{r}_{j+1}$ are employed as shown in

Fig. 3.1. The start-off value is evaluated at $t = t^n$,

$$\mathbf{U}_j^* = \frac{1}{2}\left[\mathbf{U}_{j-1} + \frac{h}{2}(\mathbf{U}_x)_{j-1} + \right.$$
$$\left. \mathbf{U}_{j+1} - \frac{h}{2}(\mathbf{U}_x)_{j+1}\right]. \tag{4.12}$$

To estimate the flux $\mathbf{F}_j^{n+1/2}$ in (4.8) for an arbitrary $j$, one first obtains

$$\mathbf{r}_j(x_j, t^{n+1/2}) = \mathbf{U}_j + \frac{\Delta t}{2}(\mathbf{U}_t)_j, \tag{4.13}$$

and then evaluates the flux at this state. The above completes the algorithm of the N-T scheme.

Instead of (4.13), which was utilized by Nessyahu and Tadmor (1990), Chang (1995) employed

$$(\mathbf{F}_t)_j = (\mathbf{A}_c)_j(\mathbf{U}_t)_j, \tag{4.14a}$$

and

$$\mathbf{F}_j^{n+1/2} = \mathbf{F}_j + \frac{\Delta t}{2}(\mathbf{F}_t)_j . \tag{4.14b}$$

This calculation is costlier than (4.13), but it is useful in the derivation of the reversible scheme.

Finally, to update the slope for the CE/SE scheme, the point value $\widehat{\mathbf{U}}_j^{n+1}$ is calculated by the reconstruction,

$$\widehat{\mathbf{U}}_j^{n+1} = \mathbf{U}_j + (\mathbf{U}_t)_j \Delta t, \tag{4.15}$$

and

$$(\mathbf{U}_x)_j^{n+1} = \frac{1}{h}\text{wtav}(\widehat{\mathbf{U}}_{j+1}^{n+1} - \mathbf{U}_j^{n+1}, \mathbf{U}_j^{n+1} - \widehat{\mathbf{U}}_{j-1}^{n+1}). \tag{4.16}$$

The above completes the description of the CE/SE method.

**4.2. Second-order upwind scheme.** Here, for reasons of economy, we interpolate the primitive variable $\mathbf{V}$. Equation (4.4) implies

$$\mathbf{V}_t + \mathbf{A}_p\mathbf{V}_x = 0, \tag{4.17}$$

where

$$\mathbf{V} = \begin{pmatrix} \rho \\ u \\ p \end{pmatrix}, \qquad \mathbf{A}_p = \begin{pmatrix} u & \rho & 0 \\ 0 & u & 1/\rho \\ 0 & \gamma p & u \end{pmatrix} . \tag{4.18}$$

Note that $\mathbf{A}_p$ is simpler than $\mathbf{A}_c$. As a result, (4.17) is slightly more economical than (4.4).

By applying the midpoint rule to the integral form (4.3) on the cell $j$,

$$\mathbf{U}_j^{n+1} = \mathbf{U}_j + \frac{\Delta t}{h}(\mathbf{F}_{j-1/2}^{n+1/2} - \mathbf{F}_{j+1/2}^{n+1/2}). \tag{4.19}$$

The problem reduces to obtaining $\{\mathbf{F}_{j+1/2}^{n+1/2}\}$. To this end, first $\{\mathbf{V}_j\}$ are calculated in a straightforward manner from their definitions. Next, the spatial derivative is estimated by Van Albada's weighted average:

$$(\mathbf{V}_x)_j = \frac{1}{h}\text{wtav}(\mathbf{V}_{j+1} - \mathbf{V}_j, \mathbf{V}_j - \mathbf{V}_{j-1}). \tag{4.20}$$

The time derivative then follows from (4.17):

$$(\mathbf{V}_t)_j = -(\mathbf{A}_p)_j(\mathbf{V}_x)_j. \tag{4.21}$$

The linear reconstruction takes the form

$$\mathbf{r}_j(x, t) = \mathbf{V}_j + (\mathbf{V}_x)_j(x - x_j) + (\mathbf{V}_t)_j(t - t^n). \tag{4.22}$$

At each interface $j + 1/2$ time $t^{n+1/2}$, as shown in Fig. 2.1, the reconstruction $\mathbf{r}_j$ yields a value denoted by $\mathbf{V}_L$, and the reconstruction $\mathbf{r}_{j+1}$, $\mathbf{V}_R$:

$$\mathbf{V}_L = \mathbf{r}_j(x_{j+1/2}, t^{n+1/2}), \qquad \text{and}$$
$$\mathbf{V}_R = \mathbf{r}_{j+1}(x_{j+1/2}, t^{n+1/2}). \tag{4.23}$$

With the above $\mathbf{V}_L$ and $\mathbf{V}_R$, the flux is calculated by upwinding as described below.

**4.3. Upwind flux.** The upwind flux employed here is Roe's splitting (1986) with an entropy correction. Roe's scheme amounts to picking the upwind side in the characteristic frame depending on the wave speed.

Recall that for the case of advection, the upwind flux is

$$f_U = \frac{1}{2}(au_L + au_R) - \frac{1}{2}|a|(u_R - u_L),$$

where $a$ is the wave speed. In other words, with

$$f_R - f_L = a(u_R - u_L),$$

dissipation is added by replacing $a$ with $|a|$. For the Euler equations, $\mathbf{A}_c = \partial\mathbf{F}/\partial\mathbf{U}$. The question is: does the mean value theorem hold for these equations? More precisely, with $\mathbf{V}_L$ and $\mathbf{V}_R$ given by (4.23), denote $\Delta\mathbf{F} = \mathbf{F}_R - \mathbf{F}_L$ and $\Delta\mathbf{U} = \mathbf{U}_R - \mathbf{U}_L$. Then, is there a state $\tilde{\mathbf{U}}$ which satisfies

$$\Delta\mathbf{F} = \tilde{\mathbf{A}}_c\Delta\mathbf{U} ?$$

If such a state exists, the upwind flux is given by

$$\mathbf{F}_U = \frac{1}{2}(\mathbf{F}_L + \mathbf{F}_R) - \frac{1}{2}|\tilde{\mathbf{A}}_c|\Delta\mathbf{U},$$

where $|\tilde{\mathbf{A}}_c|$ is calculated via a diagonalization.

Another way to explain the motivation for Roe's state is as follows. With $\mathbf{V}_L$ and $\mathbf{V}_R$ given by (4.23), which state should we use for the diagonalization? We could use the average state $\frac{1}{2}(\mathbf{V}_L + \mathbf{V}_R)$, but when the wave speed is zero, such a state

leads to ambiguity in the upwind side selection. This ambiguity can be avoided by using Roe's tilde state.

To diagonalize $\mathbf{A}_c$, denote the Jacobian matrix of the transformation between the primitive and conservative variables by $\mathbf{M}$ (Hirsch 1990):

$$\mathbf{M} = \frac{\partial \mathbf{U}}{\partial \mathbf{V}}. \qquad (4.24)$$

Then

$$\mathbf{M} = \begin{pmatrix} 1 & 0 & 0 \\ u & \rho & 0 \\ u^2/2 & \rho u & 1/(\gamma - 1) \end{pmatrix}. \qquad (4.25)$$

By applying the chain rule to (4.4),

$$\mathbf{A}_p = \mathbf{M}^{-1}\mathbf{A}_c\mathbf{M}. \qquad (4.26)$$

Thus the diagonalization of $\mathbf{A}_c$ reduces to that of $\mathbf{A}_p$. The eigenvalues of $\mathbf{A}_p$ are

$$\lambda^{(1)} = u - a, \quad \lambda^{(2)} = u, \quad \lambda^{(3)} = u + a. \quad (4.27)$$

Let $\mathbf{R}_p$ be the matrix of the right eigenvectors of $\mathbf{A}_p$; $\mathbf{L}_p$, that of the left; then, $\mathbf{L}_p = \mathbf{R}_p^{-1}$, and

$$\mathbf{L}_p = \begin{pmatrix} 0 & -\rho/(2a) & 1/(2a^2) \\ 1 & 0 & -1/a^2 \\ 0 & \rho/(2a) & 1/(2a^2) \end{pmatrix}, \qquad (4.28)$$

$$\mathbf{R}_p = \begin{pmatrix} 1 & 1 & 1 \\ -a/\rho & 0 & a/\rho \\ a^2 & 0 & a^2 \end{pmatrix}. \qquad (4.29)$$

Denote by $\mathbf{\Lambda}$ be the diagonal matrix whose diagonal entries are $\lambda^{(1)}$, $\lambda^{(2)}$, and $\lambda^{(3)}$. Then

$$\mathbf{L}_p \mathbf{A}_p \mathbf{R}_p = \mathbf{\Lambda} \quad \text{or} \quad \mathbf{A}_p = \mathbf{R}_p \mathbf{\Lambda} \mathbf{L}_p. \qquad (4.30a, b)$$

The diagonalization of $\mathbf{A}_c$ follows from the above and (4.26):

$$\mathbf{L}_c \mathbf{A}_c \mathbf{R}_c = \mathbf{\Lambda} \quad \text{or} \quad \mathbf{A}_c = \mathbf{R}_c \mathbf{\Lambda} \mathbf{L}_c, \qquad (4.31a, b)$$

where

$$\mathbf{L}_c = \mathbf{L}_p \mathbf{M}^{-1}, \qquad \mathbf{R}_c = \mathbf{M} \mathbf{R}_p. \qquad (4.32a, b)$$

Let $H$ be the total enthalpy,

$$H = (e + p)/\rho. \qquad (4.33)$$

Then, (4.32b), (4.29), and (4.25) lead to

$$\mathbf{R}_c = \begin{pmatrix} 1 & 1 & 1 \\ u - a & u & u + a \\ H - ua & u^2/2 & H + ua \end{pmatrix}. \qquad (4.34)$$

Note that $\mathbf{L}_c$ is complicated, but we only need $\mathbf{R}_c$ and $\mathbf{L}_p$ for the final expression of Roe's scheme.

Roe's tilde state is a state $\tilde{\mathbf{U}}$ which satisfies

$$\Delta \mathbf{F} = \tilde{\mathbf{A}}_c \Delta \mathbf{U}, \qquad (4.35)$$

and

$$\Delta \mathbf{U} = \tilde{\mathbf{M}} \Delta \mathbf{V}. \qquad (4.36)$$

Here, (4.35) provides consistency to the picking process, while (4.36) leads to the use of $\mathbf{L}_p$, which is more economical than $\mathbf{L}_c$. Solving the above two equations, we obtain

$$\tilde{\rho} = \sqrt{\rho_L \rho_R}, \qquad (4.37)$$

and with

$$\beta_L = \rho_L/(\rho_L + \tilde{\rho}), \qquad \beta_R = 1 - \beta_L, \qquad (4.38)$$

the other two quantities are given by

$$\tilde{u} = \beta_L u_L + \beta_R u_R, \qquad (4.39)$$

$$\tilde{H} = \beta_L H_L + \beta_R H_R. \qquad (4.40)$$

With the above tilde state, we can multiply by $\tilde{\mathbf{L}}_c$ to switch to the characteristic frame:

$$\mathbf{G}_L = \tilde{\mathbf{L}}_c \mathbf{F}_L, \quad \text{and} \quad \mathbf{G}_R = \tilde{\mathbf{L}}_c \mathbf{F}_R. \qquad (4.41)$$

The $i$-th component of either $\mathbf{G}_L$ or $\mathbf{G}_R$ is chosen by the sign of $\tilde{\lambda}^{(i)}$ for $i = 1, 2$ and 3. The result is an upwind flux in the characteristic frame

$$\mathbf{G}_U = \tfrac{1}{2}(\mathbf{G}_L + \mathbf{G}_R) - \tfrac{1}{2}\text{sign}(\tilde{\mathbf{\Lambda}})\Delta\mathbf{G}. \qquad (4.42)$$

Next, we switch back to the regular frame by multiplying by $\tilde{\mathbf{R}}_c$: $\mathbf{F}_U = \tilde{\mathbf{R}}_c \mathbf{G}_U$,

$$\mathbf{F}_U = \tfrac{1}{2}(\mathbf{F}_L + \mathbf{F}_R) - \tfrac{1}{2}\tilde{\mathbf{R}}_c \, \text{sign}(\tilde{\mathbf{\Lambda}})\, \tilde{\mathbf{L}}_c \Delta\mathbf{F}.$$

Using (4.35) and (4.31b), one obtains

$$\mathbf{F}_U = \tfrac{1}{2}(\mathbf{F}_L + \mathbf{F}_R) - \tfrac{1}{2}\tilde{\mathbf{R}}_c \, |\tilde{\mathbf{\Lambda}}| \, \tilde{\mathbf{L}}_c \Delta\mathbf{U}. \qquad (4.43)$$

Note that the above expression involves no conditional statement; i.e., due to (4.35), if $\tilde{\lambda}^{(i)} = 0$ the choice of the left or the right value in (4.42) leads to the same final upwind flux. Next, due to (4.36), $\tilde{\mathbf{L}}_c \Delta\mathbf{U} = \tilde{\mathbf{L}}_p \Delta\mathbf{V}$. As a result,

$$\mathbf{F}_U = \tfrac{1}{2}(\mathbf{F}_L + \mathbf{F}_R) - \tfrac{1}{2}\tilde{\mathbf{R}}_c \, |\tilde{\mathbf{\Lambda}}| \, \tilde{\mathbf{L}}_p \Delta\mathbf{V}. \qquad (4.44)$$

We now discuss the entropy correction. Set

$$\Delta\mathbf{W} = \tilde{\mathbf{L}}_p \Delta\mathbf{V}. \qquad (4.45)$$

For the first characteristic, the spreading rate can be estimated by (for details, see (Huynh 1995a)),

$$\Delta\lambda = \Delta\lambda^{(1)} = \frac{(\gamma+1)\tilde{a}\Delta w^{(1)}}{2\tilde{\rho}}; \qquad (4.46a)$$

or if it is the third one,

$$\Delta\lambda = \Delta\lambda^{(3)} = \frac{(\gamma+1)\tilde{a}\Delta w^{(3)}}{2\tilde{\rho}}. \qquad (4.46b)$$

For any real number $z$, define the negative part of $z$ by $z^- = \min(z,0)$, and the positive parts of $z$ by $z^+ = \max(z,0)$. To see how far the wave speeds spread past zero, evaluate

$$\eta = (-|\tilde{\lambda}| + \tfrac{1}{2}\Delta\lambda)^+. \qquad (4.47)$$

For an approximation of Osher's flux (1984), set

$$\eta = (-|\tilde{\lambda}| + \tfrac{1}{2}|\Delta\lambda|)^+. \qquad (4.48)$$

Roe's splitting with an entropy correction can be coded as follows.

Given the left and right states, calculate $H_L$ and $H_R$ by (4.33) and the tilde state by (4.37)–(4.40). Next, if $\tilde{u} \geq 0$, obtain $\Delta w^{(1)}$ via (4.45), $\Delta\lambda^{(1)}$ (4.46a), $\eta^{(1)}$ (4.47) and, with $\mathbf{R}_c^1$ denoting the first column of $\mathbf{R}_c$,

$$\mathbf{F}_U = \mathbf{F}_L + [(\tilde{u}-\tilde{a})^- + \tfrac{1}{2}\eta^{(1)}]\Delta w^{(1)}\tilde{\mathbf{R}}_c^1. \quad (4.49a)$$

Otherwise, obtain $\Delta w^{(3)}$ via (4.45), $\Delta\lambda^{(3)}$ (4.46b), $\eta^{(3)}$ (4.47), and

$$\mathbf{F}_U = \mathbf{F}_R - [(\tilde{u}+\tilde{a})^+ + \tfrac{1}{2}\eta^{(3)}]\Delta w^{(3)}\tilde{\mathbf{R}}_c^3. \quad (4.49b)$$

**5.    Two-dimensional extensions on a quadrilateral mesh.** The second-order upwind, N-T, and CE/SE schemes are extended to the 2D case in this and the next section.

The 2D Euler equations take the form

$$\mathbf{U}_t + \mathbf{F}(\mathbf{U})_x + \mathbf{G}(\mathbf{U})_y = 0, \qquad (5.1)$$

where

$$\mathbf{U} = \begin{pmatrix} \rho \\ \rho u \\ \rho v \\ e \end{pmatrix},$$

$$\mathbf{F} = \begin{pmatrix} \rho u \\ \rho u^2 + p \\ \rho uv \\ (e+p)u \end{pmatrix}, \text{ and } \mathbf{G} = \begin{pmatrix} \rho v \\ \rho vu \\ \rho v^2 + p \\ (e+p)v \end{pmatrix}.$$

At regions where $\mathbf{U}$ is smooth, (5.1) is equivalent to the non-conservation form

$$\mathbf{U}_t + \mathbf{A}_c\mathbf{U}_x + \mathbf{B}_c\mathbf{U}_y = 0; \qquad (5.2)$$

here, with $\mathbf{F}^{(k)}$ denoting the $k$-th component of $\mathbf{F}$,

$$(\mathbf{A}_c)_{k,l} = \frac{\partial\mathbf{F}^{(k)}}{\partial\mathbf{U}^{(l)}} \quad \text{and} \quad (\mathbf{B}_c)_{k,l} = \frac{\partial\mathbf{G}^{(k)}}{\partial\mathbf{U}^{(l)}}.$$

Let $\mathbf{V}$ be the vector of primitive variables. Then (5.2) can be put in the form

$$\mathbf{V}_t + \mathbf{A}_p\mathbf{V}_x + \mathbf{B}_p\mathbf{V}_y = 0. \qquad (5.3)$$

As in the 1D case, $\mathbf{A}_p$ and $\mathbf{B}_p$ are simpler than $\mathbf{A}_c$ and $\mathbf{B}_c$.

Let $\Omega$ be a spatial control volume whose boundary is $\partial\Omega$. At each point on the boundary, denote by $\vec{n} = (n_x, n_y)$ the outward unit normal and

$$\vec{\mathbf{F}} = (\mathbf{F}, \mathbf{G}).$$

Then the normal flux vector is

$$\vec{\mathbf{F}} \cdot \vec{n} = n_x\mathbf{F} + n_y\mathbf{G}. \qquad (5.4)$$

The Euler equations in integral form for the control volume $\Omega$ take the form

$$\frac{\partial}{\partial t}\int\int_\Omega \mathbf{U}\,dx\,dy + \oint_{\partial\Omega}\vec{\mathbf{F}}\cdot\vec{n}\,ds = 0. \qquad (5.5)$$

A discussion on the integral form for a space-time domain can be found in, e.g., (Roe 1983).

Below, $\Omega$ is a polygon of index $j$ with $n_e$ edges. Denote by $\text{area}_j$ the area of $\Omega$, and by $\mathbf{U}_j$ the average value of $\mathbf{U}$ on $\Omega$ at time $t^n$. For each $i$-th edge, denote by $l_i$ its length, $\vec{n}_i$ its outward unit normal, and $\vec{\mathbf{F}}_i^{n+1/2}$ the flux vector $\vec{\mathbf{F}}$ at the midpoint of the edge at time $t^{n+1/2}$, $i = 1,\ldots,n_e$. Then, for second-order accuracy, (5.5) can be approximated by

$$\mathbf{U}_j^{n+1} = \mathbf{U}_j - \frac{\Delta t}{\text{area}_j}\sum_{i=1}^{n_e} l_i\,(\vec{\mathbf{F}}_i^{n+1/2}\cdot\vec{n}_i). \quad (5.6)$$

Next, let the $(x,y)$-plane be divided into nonoverlapping quadrilaterals called cells; two adjacent quadrilaterals share a common edge (Fig. 5.1). The mesh can be structured: the cells are indexed by $i$ and $j$; it can be unstructured: the cells are indexed by $j$ and, for each cell $j$, the four neighboring cells are located by pointers. Without loss of generality, we assume the mesh is unstructured.
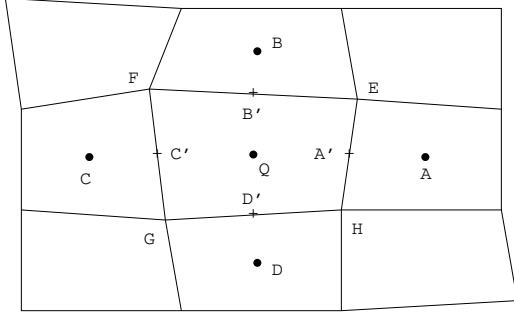
14

Fig. 5.1. *A quadrilateral cell and its four neighbors. The centroids of the cells are denoted by $Q$, $A$, $B$, $C$, and $D$. For the second-order upwind scheme, to update the solution at $Q$, we need the fluxes at time $t^{n+1/2}$ at the midpoints of the cell edges, denoted by $A'$, $B'$, $C'$, and $D'$.*

**5.1. Second-order upwind scheme.** Suppose the data $\{\mathbf{U}_j\}$ (i.e., $\{\mathbf{U}_j^n\}$) are known; $\mathbf{U}_j$ approximates the average of $\mathbf{U}$ on the cell $j$. We wish to calculate $\{\mathbf{U}_j^{n+1}\}$.

Let $EFGH$ be a typical quadrilateral whose centroid is $Q$; this quadrilateral is also identified as the cell $Q$. Let $A$, $B$, $C$, and $D$ be the centroids of the four neighboring cells as shown by Fig. 5.1.

Applying (5.6) to the cell $Q$, the problem reduces to obtaining the fluxes at time $t^{n+1/2}$ across the four edges. These fluxes are evaluated at the midpoint $A'$, $B'$, $C'$, and $D'$ of the edges by using the reconstruction functions described below.

For reasons of economy, we interpolate $\mathbf{V}$. The spatial derivatives of $\mathbf{V}$ at $Q$ are estimated by Van Albada's weighted average in the following manner. Set

$$\vec{e}_\xi = \tfrac{1}{2}(\vec{FE} + \vec{GH}) \quad \text{and} \quad \vec{e}_\eta = \tfrac{1}{2}(\vec{GF} + \vec{HE}). \quad (5.7a)$$

We can also use

$$\vec{e}_\xi = \tfrac{1}{2}\vec{CA} \quad \text{and} \quad \vec{e}_\eta = \tfrac{1}{2}\vec{DB}. \quad (5.7b)$$

Next, let $\vec{e}_\xi$ and $\vec{e}_\eta$ be the basis vectors for the $(\xi, \eta)$ coordinates:

$$(x, y) = \xi\,\vec{e}_\xi + \eta\,\vec{e}_\eta. \quad (5.7c)$$

Then

$$\begin{aligned}
(\mathbf{V}_\xi)_Q &= \text{wtav}(\mathbf{V}_A - \mathbf{V}_Q, \mathbf{V}_Q - \mathbf{V}_C), \\
(\mathbf{V}_\eta)_Q &= \text{wtav}(\mathbf{V}_B - \mathbf{V}_Q, \mathbf{V}_Q - \mathbf{V}_D).
\end{aligned} \quad (5.8)$$

From $\mathbf{V}_\xi$ and $\mathbf{V}_\eta$, the chain rule yields $\mathbf{V}_x$ and $\mathbf{V}_y$. The time derivative follows by applying (5.3):

$$(\mathbf{V}_t)_Q = -(\mathbf{A}_p)_Q (\mathbf{V}_x)_Q - (\mathbf{B}_p)_Q (\mathbf{V}_y)_Q. \quad (5.9)$$

The linear reconstruction takes the form

$$\begin{aligned}
\mathbf{r}_Q(x, y, t) &= \mathbf{V}_Q + (\mathbf{V}_x)_Q (x - x_Q) + \\
&\quad (\mathbf{V}_y)_Q (y - y_Q) + (\mathbf{V}_t)_Q (t - t^n).
\end{aligned} \quad (5.10)$$

Note that the term 'reconstruction' is typically used for $\mathbf{r}_Q(x, y, t^n)$. Here, we use the term 'reconstruction' in the extended sense above.

At time $t^{n+1/2}$, at the midpoint of each of the four edges, say, at $A'$, the reconstruction $\mathbf{r}_Q$ yields a value denoted by $\mathbf{V}_L$; the reconstruction $\mathbf{r}_A$, a value denoted by $\mathbf{V}_R$:

$$\begin{aligned}
\mathbf{V}_L &= \mathbf{r}_Q(x_{A'} - x_Q, y_{A'} - y_Q, t^{n+1/2}), \\
\mathbf{V}_R &= \mathbf{r}_A(x_{A'} - x_A, y_{A'} - y_A, t^{n+1/2}).
\end{aligned} \quad (5.11)$$

With the above $\mathbf{V}_L$ and $\mathbf{V}_R$, the flux is calculated by upwinding as described in the next subsection. The other three fluxes at $B'$, $C'$, and $D'$ are similar.

Observe that for this upwind scheme, there is a trade-off between computing time and storage if the quadrilateral mesh is unstructured. This observation also holds for an unstructured triangular mesh, but if the mesh is structured, it does not apply. For each cell $Q$, we can carry out five reconstructions (per cell) at $Q$, $A$, $B$, $C$, and $D$, and four upwind fluxes to update the solution at $Q$. If the number of cells is $N$, this approach needs to evaluate $5N$ (vector) reconstructions and $4N$ upwind fluxes. The number of values $\{\mathbf{U}_j^n\}$ and $\{\mathbf{U}_j^{n+1}\}$ to be stored is $2 \times (4N)$. To reduce computing time, we can carry out only one reconstruction per cell and store the left and right interface values at time $t^{n+1/2}$ for all edges, and then, loop over all edges to get the fluxes. If the number of cells is $N$, the number of edges is roughly $2N$. This approach needs to evaluate only $N$ reconstructions and $2N$ upwind fluxes. However, the number of additional values to be stored is $4 \times (4N)$.

**5.2. Upwind flux.** Again the upwind flux employed here is Roe's splitting (1986) with the entropy correction in Huynh (1995).

Let $\vec{n} = (n_x, n_y)$ be the outward unit normal at $A'$. If $\mathbf{U}$ is a state at $A'$, then, because the flux is homogeneous of degree one in $\mathbf{U}$,

$$\vec{\mathbf{F}} \cdot \vec{n} = (n_x, n_y) \cdot (\mathbf{F}, \mathbf{G}) = (n_x \mathbf{A}_c + n_y \mathbf{B}_c)\mathbf{U}. \quad (5.12)$$

To obtain the upwind flux, we need to diagonalize the matrix $(n_x \mathbf{A}_c + n_y \mathbf{B}_c)$. (Notice that because the tangential component is ignored, this diagonalization is sometimes said to be dimensional

splitting; but the flux for the centered schemes is also given by (5.12) where the tangential component is ignored. As such, the upwind and centered schemes here are equally dimensional splitting or non-splitting.) Next, set

$$q = n_x u + n_y v. \tag{5.13}$$

Then the eigenvalues of the matrix $(n_x \mathbf{A}_c + n_y \mathbf{B}_c)$ are

$$\lambda^{(1)} = q - a, \quad \lambda^{(2)} = q, \quad \lambda^{(3)} = q, \quad \lambda^{(4)} = q + a; \tag{5.14}$$

and, the matrix $\mathbf{R}_c$ of the right eigenvectors is

$$\begin{pmatrix} 1 & 0 & 1 & 1 \\ u & n_x a & n_y a & u & u + n_x a \\ v & n_y a & n_x a & v & v + n_y a \\ H & qa & a(n_y u + n_x v) & \frac{u^2 + v^2}{2} & H + qa \end{pmatrix}. \tag{5.15}$$

As for the matrix of the left eigenvectors of $n_x \mathbf{A}_p + n_y \mathbf{B}_p$, where the subscript $p$ stands for 'primitive',

$$\mathbf{L}_p = \begin{pmatrix} 0 & n_x \rho/(2a) & n_y \rho/(2a) & 1/(2a^2) \\ 0 & n_y \rho/a & n_x \rho/a & 0 \\ 1 & 0 & 0 & 1/a^2 \\ 0 & n_x \rho/(2a) & n_y \rho/(2a) & 1/(2a^2) \end{pmatrix}. \tag{5.16}$$

We only need the first and last columns of (5.15) and the first and last rows of (5.16).

Roe's splitting with an entropy correction can be coded as follows.

Let the left and right states and the unit normal $\vec{n}$ be given. First, calculate $H_L$ and $H_R$ by an expression similar to (4.33) and the tilde state by (4.37)–(4.40) ($\tilde{v}$ is similar to $\tilde{u}$ and $\tilde{H}$). Next, if $\tilde{q} \geq 0$, obtain $\Delta w^{(1)}$ via (4.45), $\Delta \lambda^{(1)}$ (4.46a), $\eta^{(1)}$ (4.47) and, with $\mathbf{R}_c^1$ denoting the first column of $\mathbf{R}_c$,

$$(\vec{\mathbf{F}} \cdot \vec{n})_U = (\vec{\mathbf{F}} \cdot \vec{n})_L + [(\tilde{q} - \tilde{a})^- + \tfrac{1}{2}\eta^{(1)}]\Delta w^{(1)} \tilde{\mathbf{R}}_c^1. \tag{5.17a}$$

Otherwise, obtain $\Delta w^{(4)}$ via (4.45), $\Delta \lambda^{(4)}$ (4.46b), $\eta^{(4)}$ (4.47), and

$$(\vec{\mathbf{F}} \cdot \vec{n})_U = (\vec{\mathbf{F}} \cdot \vec{n})_R - [(\tilde{q} + \tilde{a})^+ + \tfrac{1}{2}\eta^{(4)}]\Delta w^{(4)} \tilde{\mathbf{R}}_c^4. \tag{5.17b}$$

**5.3. Second-order centered schemes.** The CE/SE schemes were extended to the case of a quadrilateral mesh by Zhang et al. (2002). This approach to extension is also applied to the N-T scheme below. The key difference here, however, is that the slope estimates are simplified, and they only use the data at the current time level. The



Fig. 5.2. *Reconstruction cells for the two centered schemes. The domain where the reconstruction is valid at time $t = t^n$ for the cell $Q$ is the dotted line octagon, namely, $AEBFCGDH$. This octagon is called the reconstruction cell $Q$ here. It is also the control volume on which fluxes are balanced when the solution at $Q$ is updated. The centroid of the octagon is marked by the gray dot and denoted by $\bar{Q}$.*

simplified slope estimate together with an observation on distributing the fluxes to the neighboring cells result in an extended N-T scheme which is faster and requires considerably less storage than the CE/SE method for the case of an unstructured mesh. The extended CE/SE scheme here is also different from that by Zhang et al. (2002) in that essentially, the scheme which carries along the interface values explained after (3.14) is employed. Finally, the presentation below is simpler.

Consider a typical quadrilateral $EFGH$ whose centroid is $Q$ shown in Fig. 5.2. Let the centroid of the four neighboring cells be denoted by $A$, $B$, $C$, and $D$. For the two centered schemes, let the *reconstruction cell* $Q$ be defined as the (dotted line) octagon $AEBFCGDH$. At time $t = t^n$, the reconstruction for these two schemes is assumed to be valid on this octagon. This octagon is roughly twice as big as the original cell $EFGH$. The centroid of this octagon is marked by a gray dot and denoted by $\bar{Q}$. It is called a *solution point* by Wang and Chang (1999). Note that the two adjacent reconstruction cells $Q$ and $A$ overlap on the quadrilateral $QHAE$ (see also Fig. 5.3). Also note that the reconstruction cell is the spatial part (or spatial projection) of the solution element in Zhang et al. (2002).

At time $t^n$, assume that we know $\mathbf{U}_j$ (i.e. $\mathbf{U}_j^n$) for all $j$; $\mathbf{U}_j$ approximates the average of $\mathbf{U}$ on the reconstruction cell $j$. We wish to calculate $\mathbf{U}_j^{n+1}$.

If $\mathbf{U}_Q$ approximates the average of $\mathbf{U}$ on the reconstruction cell $Q$, it is considered to be the value

**U** at the centroid $\bar{Q}$, not at $Q$. A more precise notation would be $\mathbf{U}_{\bar{Q}}$, but for simplicity, we use the notation $\mathbf{U}_Q$ below.

In the CE/SE case, the spatial slopes $\{(\mathbf{U}_x)_j\}$ and $\{(\mathbf{U}_y)_j\}$ are also stored; in addition, $\{(\mathbf{U}_x)_j^{n+1}\}$ and $\{(\mathbf{U}_y)_j^{n+1}\}$ must be calculated.

To update the cell average $\mathbf{U}_Q^{n+1}$, the data $\mathbf{U}_Q$ is ignored. Denote by $\mathrm{area}(\mathrm{rc}(Q))$ the area of the reconstruction cell $Q$. Then the midpoint rule for (5.5) on this reconstruction cell takes the form

$$\mathrm{area}(\mathrm{rc}(Q))\,\mathbf{U}_Q^{n+1} = \mathrm{area}(\mathrm{rc}(Q))\,\mathbf{U}_Q^*$$
$$\Delta t \sum_{i=1}^{8} l_i\,(\vec{\mathbf{F}}_i^{n+1/2} \cdot \vec{n}_i)\,, \tag{5.18}$$

where the start-off value $\mathbf{U}_Q^*$ and the eight fluxes are defined below.

First, we describe the reconstruction function for an arbitrary cell, say, cell $Q$ shown in Fig. 5.2. In the CE/SE case, the stored slopes are employed; for the N-T case, they are estimated by using the neighboring values as follows. Let the $(\xi, \eta)$ coordinates be defined as in (5.7c) with $\vec{e}_\xi = \frac{1}{2}\overrightarrow{CA}$ and $\vec{e}_\eta = \frac{1}{2}\overrightarrow{DB}$, or we can also use (5.7a). Then

$$\begin{aligned}
(\mathbf{U}_\xi)_Q &= \mathrm{wtav}(\mathbf{U}_A \quad \mathbf{U}_Q, \mathbf{U}_Q \quad \mathbf{U}_C),\\
(\mathbf{U}_\eta)_Q &= \mathrm{wtav}(\mathbf{U}_B \quad \mathbf{U}_Q, \mathbf{U}_Q \quad \mathbf{U}_D).
\end{aligned} \tag{5.19}$$

From $\mathbf{U}_\xi$ and $\mathbf{U}_\eta$, the chain rule yields $\mathbf{U}_x$ and $\mathbf{U}_y$. The time derivative follows by applying (5.2):

$$(\mathbf{U}_t)_Q = (\mathbf{A}_c)_Q\,(\mathbf{U}_x)_Q \quad (\mathbf{B}_c)_Q\,(\mathbf{U}_y)_Q. \tag{5.20}$$

Since $\mathbf{U}_Q$ approximates the average of $\mathbf{U}$ on the reconstruction cell, it is considered to be the value at the centroid $\bar{Q}$. The linear reconstruction takes the form

$$\begin{aligned}
\mathbf{r}_Q(x,y,t) &= \mathbf{U}_Q + (\mathbf{U}_x)_Q\,(x \quad x_{\bar{Q}}) +\\
&\quad (\mathbf{U}_y)_Q\,(y \quad y_{\bar{Q}}) + (\mathbf{U}_t)_Q\,(t \quad t^n).
\end{aligned} \tag{5.21}$$

Note $x_{\bar{Q}}$ and $y_{\bar{Q}}$ above take the place of $x_Q$ and $y_Q$ in the upwind case.

Next, when updating the solution at $Q$, $\mathbf{r}_Q$ is ignored, and $\mathbf{r}_A$, $\mathbf{r}_B$, $\mathbf{r}_C$, and $\mathbf{r}_D$ are employed (Fig. 5.3). The start-off (average) value at time $t = t^n$ is calculated as follows. Denote by $A^*$ the centroid of $QHAE$; $B^*$, $QEBF$; $C^*$, $QFCG$; and, $D^*$, $QGDH$. At time $t = t^n$, the value at $A^*$ is evaluated using $\mathbf{r}_A$; at $B^*$, using $\mathbf{r}_B$; $C^*$, $\mathbf{r}_C$; and



Fig. 5.3. *Centered schemes. When updating the solution at $Q$, $\mathbf{r}_Q$ is ignored, while the reconstructions at the neighboring cells $\mathbf{r}_A$, $\mathbf{r}_B$, $\mathbf{r}_C$, and $\mathbf{r}_D$ are employed. Denote the centroid of $QHAE$ by $A^*$, and that of $QEBF$, $B^*$. The values needed from $\mathbf{r}_A$ are: at time $t^n$ at $A^*$; at time $t^{n+1/2}$ at the midpoint of $AH$ and the midpoint of $AE$ marked by $(+)$. The values needed from $\mathbf{r}_B$ are: at time $t^n$ at $B^*$; at time $t^{n+1/2}$ at the midpoint of $BE$ and the midpoint of $BF$ marked by $(+)$. Similar statements hold for the reconstructions at $C$ and $D$.*

$D^*$, $\mathbf{r}_D$. The start-off value $\mathbf{U}_Q^*$ is given by

$$\begin{aligned}
\mathrm{area}(\mathrm{rc}(Q))\,\mathbf{U}_Q^* =&\\
\mathrm{area}(QHAE)\,\mathbf{r}_A(x_{A^*}, y_{A^*}, t^n) +&\\
\mathrm{area}(QEBF)\,\mathbf{r}_B(x_{B^*}, y_{B^*}, t^n) +&\\
\mathrm{area}(QFCG)\,\mathbf{r}_C(x_{C^*}, y_{C^*}, t^n) +&\\
\mathrm{area}(QGDH)\,\mathbf{r}_D(x_{D^*}, y_{D^*}, t^n);&
\end{aligned} \tag{5.22}$$

here, again, $\mathrm{area}(\mathrm{rc}(Q))$ is the area of the reconstruction cell (octagon) $AEBFCGDH$.

As for the fluxes $\vec{\mathbf{F}}_i^{n+1/2} \cdot \vec{n}_i$ across the eight edges, the fluxes across $AH$ and $AE$ are calculated by the reconstruction $\mathbf{r}_A$; those across $BE$ and $BF$, by $\mathbf{r}_B$; $CF$ and $CG$, by $\mathbf{r}_C$; and $DG$ and $DH$, by $\mathbf{r}_D$. Denote by $M$ the midpoint of $AE$. Then one calculates the conservative variables at $M$ at time $t^{n+1/2}$ via $\mathbf{r}_A(x_M, y_M, t^{n+1/2})$, and the flux across $AE$ by (5.4). The other seven fluxes are obtained in the same manner. This completes the quadrilateral-mesh extension of the N-T scheme.

Next, the slope update for the CE/SE scheme below is a simplification of that in Zhang et al. (2002). First, the point value at $\bar{A}$ at time $t^{n+1}$, denoted by $\widehat{\mathbf{U}}_{\bar{A}}^{n+1}$, is calculated by the reconstruction,

$$\widehat{\mathbf{U}}_{\bar{A}}^{n+1} = \mathbf{r}_A(x_{\bar{A}}, y_{\bar{A}}, t^{n+1})\,. \tag{5.23}$$

The other point values $\widehat{\mathbf{U}}_{\bar{B}}^{n+1}$, $\widehat{\mathbf{U}}_{\bar{C}}^{n+1}$, and $\widehat{\mathbf{U}}_{\bar{D}}^{n+1}$ are
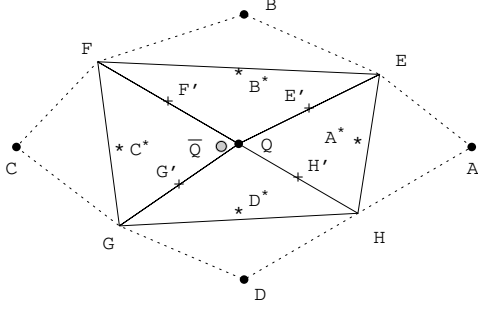
Fig. 5.4. *Distributing the fluxes results in a faster algorithm with no penalty in storage for the extended N-T scheme. After calculating* $\mathbf{U}_x$, $\mathbf{U}_y$, *and* $\mathbf{U}_t$ *at* $\bar{Q}$, *distribute the quantity* area$(QHAE)\,\mathbf{r}_Q(x_{A^*}, y_{A^*}, t^n)$ *and the two fluxes at* $H'$ *and* $E'$ *to cell* $A$; *distribute the quantity* area$(QEBF)\,\mathbf{r}_Q(x_{B^*}, y_{B^*}, t^n)$ *and the two fluxes at* $E'$ *and* $F'$ *to cell* $B$; *the distribution to cells* $C$ *and* $D$ *is similar.*

obtained in a similar manner, and

$$(\mathbf{U}_\xi)_Q = \text{wtav}(\widehat{\mathbf{U}}_{\bar{A}}^{n+1} \quad \mathbf{U}_Q^{n+1}, \mathbf{U}_Q^{n+1} \quad \widehat{\mathbf{U}}_{\bar{C}}^{n+1}),$$

$$(\mathbf{U}_\eta)_Q = \text{wtav}(\widehat{\mathbf{U}}_{\bar{C}}^{n+1} \quad \mathbf{U}_Q^{n+1}, \mathbf{U}_Q^{n+1} \quad \widehat{\mathbf{U}}_{\bar{D}}^{n+1}).$$
(5.24)

The spatial derivatives $(\mathbf{U}_x)_Q^{n+1}$ and $(\mathbf{U}_y)_Q^{n+1}$ can then be calculated via the chain rule. This calculation completes the quadrilateral-mesh extension of the CE/SE method.

The following observations are in order.

(a) Contrary to the upwind case, for an unstructured quadrilateral (or triangular) mesh, the above extended N-T scheme can be coded in a manner which results in a faster algorithm with no penalty in storage. Indeed, instead of gathering fluxes to update the solution, fluxes can be distributed to the neighboring cells in the following manner.

First, set all $\mathbf{U}_j^{n+1}$ to zero. Next, evaluate area$(\text{rc}(j))\,\mathbf{U}_j^{n+1}$ for all $j$ as follows. In a Do loop with index $j$, when $j$ equals $Q$, obtain $\mathbf{r}_Q$. Calculate the values of $\mathbf{r}_Q$ at time $t = t^n$ at $A^*$, $B^*$, $C^*$, and $D^*$ (see Fig. 5.4). Calculate the values of $\mathbf{r}_Q$ at time $t = t^{n+1/2}$ at the spatial locations marked by $(+)$ in Fig. 5.4. These locations are: $E'$, the midpoint of $QE$; $F'$, that of $QF$; $G'$, $QG$; and $H'$, $QH$. Distribute the quantity area$(QHAE)\,\mathbf{r}_A(x_{A^*}, y_{A^*}, t^n)$ and the two fluxes at $H'$ and $E'$ to cell $A$, and store the sum in $\mathbf{U}_A^{n+1}$; distribute the quantity area$(QEBF)\,\mathbf{r}_A(x_{B^*}, y_{B^*}, t^n)$ and the two fluxes at $E'$ and $F'$ to cell $B$, and store the sum in $\mathbf{U}_B^{n+1}$; the distribution to cells $C$ and $D$ is similar. When the



Fig. 5.5. *Two staggered meshes in Arminjon et al. (1995) overlay to form a nonstaggered mesh in Zhang et al. (2002) provided that the indices* $(k, l)$ *are assigned along the two diagonal directions.*

loop is completed, each cell will have received all the fluxes it needs to update its cell average value.

(b) The above extension of the nonstaggered N-T scheme relates to the extension of the staggered version in Arminjon et al. (1995) and Jiang and Tadmor (1998) as follows.

For the staggered version, the scheme alternates between the black and gray dots shown in Fig. 5.5. More precisely, let a mesh of squares whose edges are the black lines, centers black dots, and vertices gray dots be given. At time $t^n$, the values $\mathbf{U}_{i,j}^n$ are known at the black dots. The reconstruction at each black dot is assumed to be valid on the corresponding square (black lines). At time $t^{n+1}$, calculate the solutions at the gray dots by balancing fluxes on the squares whose edges are the gray lines. At time $t^{n+2}$, obtain the solutions at the black dots by balancing fluxes on the squares whose edges are black lines.

This extension of the staggered scheme can be made nonstaggered in the following manner. At time $t^n$, assume the data are known at both the black and gray dots. At time $t^{n+1}$, obtain the solutions at black dots by using the reconstructions at gray dots, and the solutions at gray dots by using the reconstructions at black dots. If we assign indices $(k, l)$ along the two diagonal directions (observed by Drs. Ananda Himansu and Xiao-Yen Wang), the resulting scheme is very similar to the nonstaggered extension of the N-T scheme presented here.

(c) Compared with the CE/SE scheme discussed in Zhang et al. (2002), we can save storage by carrying along the point values $\widehat{\mathbf{U}}_{\bar{Q}}$ instead of the
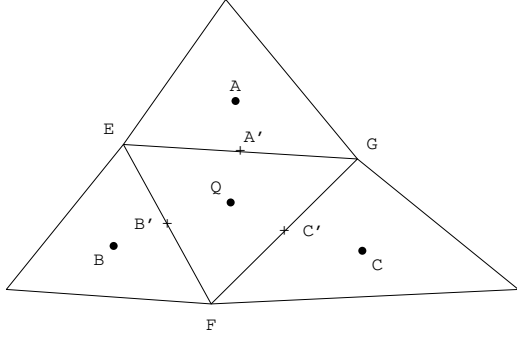
**Fig. 6.1.** *The centroid of the triangle $EFG$ is denoted by $Q$. The centroids of the three neighboring cells are denoted by $A$, $B$, and $C$. For the second-order upwind scheme, to update the solution at $Q$, we need the fluxes at time $t^{n+1/2}$ at the midpoints of the cell edges, denoted by $A'$, $B'$, and $C'$.*
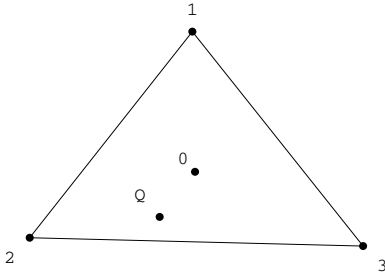


**Fig. 6.2.** *To derive the triangular-mesh version for the weighted average formula, it is convenient to denote $A$, $B$, $C$ by 1, 2, and 3 respectively.*

slopes. In addition, we can save computing time by applying the observation on distributing the fluxes so that each cell is visited only once instead of four times. Also note that for a quadrilateral mesh, since there is no reversible scheme, it is not clear how to adjust numerical dissipation even if the slopes are carried along.

(d) For the CE/SE method, the flux at time $t^{n+1/2}$ can be obtained via $(\mathbf{F}_t)_A$ and $(\mathbf{G}_t)_A$ in a manner similar to (4.14). Such a calculation, however, is costlier.

**6. Two-dimensional extensions on a triangular mesh.** For the schemes discussed here, the extension to a triangular mesh differs from that to a quadrilateral mesh described in the previous section only in the slope evaluation, which employs three neighboring data, not four.

Let the $(x,y)$-plane be divided into nonoverlapping triangles called cells; two adjacent triangles

share a common edge (Fig. 6.1). Assume that the mesh is unstructured: the cells are indexed by $j$ and, for each cell $j$, the three neighboring cells are located by pointers.

Let $EFG$ be a typical triangle whose centroid is denoted by $Q$. This triangle is also identified as the cell $Q$. Let $A$, $B$, and $C$ be the centroids of the three neighboring cells as shown in Fig. 6.1.

Next, we describe the triangular-mesh extension of Van Albada's weighted average. The following version is a modification of the extension by Wang and Chang (1999). It is designed for the upwind and the N-T schemes where only data at the current time level are available. It also makes the observation on distributing the fluxes applicable.

For the purpose of deriving this weighted average, as will be seen in (6.1), it is more convenient to use numbers 1, 2, and 3 to denote the neighbors; e.g., instead of $(x_A, y_A)$, we use $(x_1, y_1)$ (see Fig. 6.2). Let $u_Q$ be a scalar value at the the point $(x_Q, y_Q)$, and similarly, $u_1$, $u_2$, and $u_3$, at the three neighboring cell centers. Next, the values at any three points determine a plane in the $(x, y, u)$ space; for example, the values $u_1$, $u_2$, and $u_3$ at $(x_1, y_1)$, $(x_2, y_2)$, and $(x_3, y_3)$ respectively, determine a plane. Denote the slopes of this plane by $(u_x)_{123}$ and $(u_y)_{123}$.

The plane we wish to obtain is biased toward the least steep one among the three planes $12Q$, $23Q$, and $31Q$, but if $(x_Q, y_Q)$ lies on one of the edges of the triangle 123, then one of the three planes is not well defined. To avoid this problem, let $(x_0, y_0)$ be the centroid of the triangle 123. The value $u_0$ is defined by linear extrapolation using the plane 123:

$$u_0 = u_Q + (u_x)_{123}(x_0 - x_Q) + (u_y)_{123}(y_0 - y_Q).$$

The final slopes are obtained by the values at the points 0, 1, 2, and 3 as follows. After calculating the slopes of the three planes $012$, $023$, and $031$, set

$$\begin{aligned}
\theta_1 &= [(u_x)_{023}]^2 + [(u_y)_{023}]^2, \\
\theta_2 &= [(u_x)_{031}]^2 + [(u_y)_{031}]^2, \qquad (6.1) \\
\theta_3 &= [(u_x)_{012}]^2 + [(u_y)_{012}]^2,
\end{aligned}$$

and

$$\theta_{123} = \theta_1\theta_2 + \theta_2\theta_3 + \theta_3\theta_1.$$

Note the less steep the plane $012$, the smaller the quantity $\theta_3$. The final slope at $Q$ is given by biasing toward the least steep plane among the above three:

$$\begin{aligned}
(u_x)_Q &= [(\theta_1\theta_2)(u_x)_{012} + (\theta_2\theta_3)(u_x)_{023} + \\
&\quad (\theta_3\theta_1)(u_x)_{031}]/\theta_{123}.
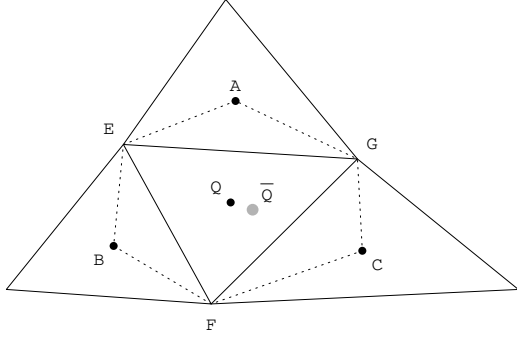\end{aligned}$$

$$(6.2)$$

Fig. 6.3. *Reconstruction cells for the two centered schemes. The domain where the reconstruction is valid at time $t = t^n$ for the cell $Q$ is the dotted line hexagon, namely, $AEBFCG$. This hexagon is called the reconstruction cell $Q$. It is also the control volume on which fluxes are balanced when the solution at $Q$ is updated. The centroid of the hexagon is marked by the gray dot and denoted by $\bar{Q}$.*

A similar expression holds for $(u_y)_Q$. This completes the triangular-mesh extension of the weighted average.

**6.1. Second-order upwind scheme.** Suppose the data $\{\mathbf{U}_j\}$ (i.e., $\{\mathbf{U}_j^n\}$) are known; $\mathbf{U}_j$ approximates the average of $\mathbf{U}$ on the cell $j$. We wish to calculate $\{\mathbf{U}_j^{n+1}\}$.

To update the solution at the cell $Q$, we need to calculate the fluxes at the midpoint $A'$, $B'$, and $C'$ at time $t^{n+1/2}$ (Fig. 6.1). These fluxes are obtained by using the reconstructions described below.

Consider an arbitrary cell, say, cell $Q$. The values $\mathbf{V}$ at $Q$ and at the three neighboring cells determine the slopes $\mathbf{V}_x$ and $\mathbf{V}_y$ at $Q$ via (6.2). The time derivative then follows from (5.9), and the reconstruction, (5.10). At time $t^{n+1/2}$, at say, $A'$, the reconstruction $\mathbf{r}_Q$ yields $\mathbf{V}_L$, and $\mathbf{r}_A$ yields $\mathbf{V}_R$. The flux at $A'$ is calculated by the upwind step (5.17). The other two fluxes are similar. These fluxes complete the upwind algorithm.

Observe that for this upwind scheme, there is a trade-off between computing time and storage. For each cell $Q$, we can carry out four reconstructions (per cell) at $Q$, $A$, $B$, and $C$, and three upwind fluxes to update the solution at $Q$. To reduce computing time, we can carry out only one reconstruction per cell and store the left and right interface values at time $t^{n+1/2}$ for all edges, and then, loop over all edges to get the fluxes.

**6.2. Second-order centered schemes.** The CE/SE schemes were extended to the case of an unstructured triangular mesh by Wang and Chang (1999). This approach to extension is also applied to the N-T scheme below. The key difference here, however, is that the slope estimates are modified so that only the data at the current time level are needed. As in the quadrilateral-mesh case, the modified slope estimate made the observation on distributing the fluxes to the neighboring cells applicable. This observation, in turn, results in an extended N-T scheme which is faster and requires considerably less storage than the CE/SE method (slopes need not be stored). The extended CE/SE scheme here is also different from that by Wang and Chang (1999) in that essentially, the scheme which carries along the interface values explained after (3.14) is employed. Finally, the presentation below is simpler.

For the two centered schemes, let the *reconstruction cell $Q$* be defined as the (dotted line) hexagon $AEBFCG$ (Fig. 6.3). At time $t = t^n$, the reconstruction for these two schemes is assumed to be valid on this reconstruction cell. The centroid of the reconstruction cell is marked by a gray dot and denoted by $\bar{Q}$. The two adjacent reconstruction cells $Q$ and $A$ overlap on the quadrilateral $QGAE$ (Fig. 6.4). Note that the reconstruction cell is the spatial part (or spatial projection) of the solution element in Wang and Chang (1999).

At time $t^n$, assume that we know $\mathbf{U}_j$ (i.e. $\mathbf{U}_j^n$) for all $j$; $\mathbf{U}_j$ approximates the average of $\mathbf{U}$ on the reconstruction cell $j$. We wish to calculate $\mathbf{U}_j^{n+1}$.

In the CE/SE case, the spatial slopes $\{(\mathbf{U}_x)_j\}$ and $\{(\mathbf{U}_y)_j\}$ are also stored; in addition, $\{(\mathbf{U}_x)_j^{n+1}\}$ and $\{(\mathbf{U}_y)_j^{n+1}\}$ must also be calculated.

Next, we describe the reconstruction function for an arbitrary cell, say, cell $Q$ shown in Fig. 6.3. In the CE/SE case, the stored spatial slopes are employed; for the N-T case, the value $\mathbf{U}_Q$, which is considered to be at $\bar{Q}$, and those at $\bar{A}$, $\bar{B}$, and $\bar{C}$ determined $\mathbf{U}_x$ and $\mathbf{U}_y$ via the weighted average (6.2). The time derivative is given by (5.20), and the reconstruction, (5.21).

The solution at $Q$ is calculated by balancing fluxes over the reconstruction cell $Q$. Denote by $\mathrm{area}(\mathrm{rc}(Q))$ the area of the reconstruction cell $Q$. Then, similar to (5.18), by the midpoint rule,

$$\mathrm{area}(\mathrm{rc}(Q))\,\mathbf{U}_Q^{n+1} = \mathrm{area}(\mathrm{rc}(Q))\,\mathbf{U}_Q^*$$
$$\Delta t \sum_{i=1}^{6} l_i \left(\vec{\mathbf{F}}_i^{n+1/2} \cdot \vec{n}_i\right). \tag{6.3}$$

When updating the solution at $Q$, $\mathbf{r}_Q$ is ignored, and $\mathbf{r}_A$, $\mathbf{r}_B$, and $\mathbf{r}_C$ are employed (Fig. 6.4). Denote by $A^*$ the centroid of $QGAE$; $B^*$, $QEBF$; and $C^*$,
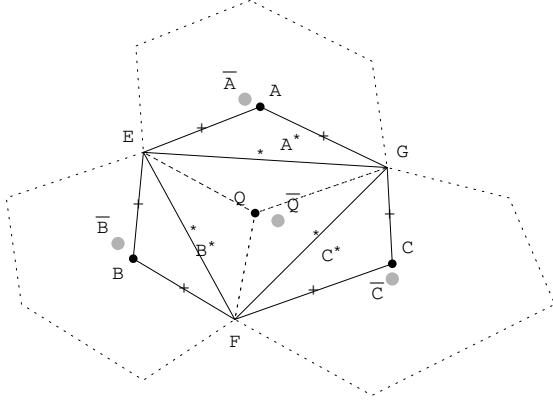
Fig. 6.4. *Centered schemes. When updating the solution at $Q$, $\mathbf{r}_Q$ is ignored, while the reconstructions at the neighboring cells $\mathbf{r}_A$, $\mathbf{r}_B$, and $\mathbf{r}_C$ are employed. Denote the centroid of $QGAE$ by $A^*$; that of $QEBF$, $B^*$; and that of $QFCG$, $C^*$. The values needed from $\mathbf{r}_A$ are: at time $t^n$ at $A^*$; at time $t^{n+1/2}$ at the midpoint of $AE$ and the midpoint of $AG$ marked by $(+)$. The values needed from $\mathbf{r}_B$ are: at time $t^n$ at $B^*$; at time $t^{n+1/2}$ at the midpoint of $BE$ and the midpoint of $BF$ marked by $(+)$. A similar statement holds for the reconstruction at $C$.*

$QFCG$. The start-off (average) value is given by

$$\begin{aligned}
\text{area}(AEBFCG)\,\mathbf{U}_Q^* &= \\
\text{area}(QGAE)\,\mathbf{r}_A(x_{A^*}, y_{A^*}, t^n) &+ \\
\text{area}(QEBF)\,\mathbf{r}_B(x_{B^*}, y_{B^*}, t^n) &+ \\
\text{area}(QFCG)\,\mathbf{r}_C(x_{C^*}, y_{C^*}, t^n) &\,.
\end{aligned} \tag{6.4}$$

As for the fluxes $\vec{\mathbf{F}}_i^{n+1/2} \cdot \vec{n}_i$ across the six edges at time $t^{n+1/2}$, those across $AG$ and $AE$ are calculated by the reconstruction $\mathbf{r}_A$; those across $BE$ and $BF$, by $\mathbf{r}_B$; and $CF$ and $CG$, by $\mathbf{r}_C$. Denote by $M$ the midpoint of $AE$. Then one calculates the conservative variables at $M$ via $\mathbf{r}_A(x_M, y_M, t^{n+1/2})$, and the flux across $AE$ by (5.4). The other five fluxes are calculated in the same manner. This completes the triangular-mesh extension of the N-T scheme.

Next, the slope update for the CE/SE scheme below is a simplification of that by Wang and Chang (1999). First, the point value at $\bar{A}$ at time $t^{n+1}$, denoted by $\widehat{\mathbf{U}}_{\bar{A}}^{n+1}$, is calculated via the reconstruction,

$$\widehat{\mathbf{U}}_{\bar{A}}^{n+1} = \mathbf{r}_A(x_{\bar{A}}, y_{\bar{A}}, t^{n+1}) \,. \tag{6.5}$$

The point values $\widehat{\mathbf{U}}_{\bar{B}}^{n+1}$ and $\widehat{\mathbf{U}}_{\bar{C}}^{n+1}$ are obtained in a similar manner. These three values and $\mathbf{U}_Q^{n+1}$



Fig. 6.5. *For the staggered version, the scheme can alternate between any two of the three sets of data: at cell vertices, at cell centers, and at midpoints of edges.*

(at $\bar{Q}$) determined $(\mathbf{U}_x)_Q^{n+1}$ and $(\mathbf{U}_y)_Q^{n+1}$ via the weighted average (6.2). This calculation completes the triangular-mesh extension of the CE/SE method.

The following remarks are in order.

(a) Distributing the fluxes results in a faster algorithm with no penalty in storage for the extended N-T scheme. Indeed, let $E'$ denote the midpoint of $QE$; $F'$, $QF$; and $G'$, $QG$. After calculating $\mathbf{U}_x$, $\mathbf{U}_y$, and $\mathbf{U}_t$ at $\bar{Q}$, distribute the quantity $\text{area}(QGAE)\,\mathbf{r}_Q(x_{A^*}, y_{A^*}, t^n)$ and the two fluxes at $G'$ and $E'$ to cell $A$; distribute the quantity $\text{area}(QEBF)\,\mathbf{r}_Q(x_{B^*}, y_{B^*}, t^n)$ and the two fluxes at $E'$ and $F'$ to cell $B$; finally, distribute the quantity $\text{area}(QFCG)\,\mathbf{r}_Q(x_{C^*}, y_{C^*}, t^n)$ and the two fluxes at $F'$ and $G'$ to cell $C$.

(b) The above observation can also be applied to the CE/SE scheme: instead of storing the slopes, we can store the point values $\widehat{\mathbf{U}}_{\bar{Q}}$ and, when visiting the cell $Q$, after distributing the fluxes, we update $\widehat{\mathbf{U}}_{\bar{Q}}^{n+1}$.

(c) The next remark concerns the extension of the staggered version. Let an unstructured triangular mesh of $N$ vertices be given. Then the number of cells is roughly $2N$, and the number of edges, roughly $3N$. The data can be stored at the vertices, the cell centers, or the midpoint of the edges. We can alternate between any two of these three sets of data, say, between data at vertices and those at edges. Then, at one time level, the flow field is resolved by $N$ points (vertices), and at the next time level, by $3N$ points (edges). Such a solution can only be as accurate as resolving the flow field by $N$ points only; thus, the scheme is not optimal, and there is no gain in accuracy compared with the nonstaggered version (see below). In addition, the reconstruction, the solution procedure, and the bookkeeping for such a method are quite involved.

The triangular-mesh extension of the N-T scheme here has $2N$ pieces of data at each time level, but because the reconstruction cell is twice as large as the original triangular cell, the flow field is effectively resolved by $N$ points (not $2N$). However, the scheme is simple and has numerous nice features as will be shown in the next section.

**7. Fourier analysis for the 2D case.** For convenience, the 2D extension of the N-T scheme presented above will be called *the centered scheme* from here on—as opposed to the upwind scheme. It is indeed the centered counterpart of the upwind scheme. When more details are needed, we refer to it as the Lax-Friedrichs type second-order accurate centered scheme. It can also be considered as a coupled version of the CE/SE method. (Since the upwind scheme is sometimes called MUSCL-Roe, this scheme could have been named LF-MUSCL-NT-CE/SE, but such a name seems too unwieldy.) If the slopes are carried along ($\epsilon = 1/2$), the corresponding 2D extension is called the CE/SE scheme with the contributions by MUSCL, N-T, and the various authors explained in the previous sections understood. Thus, in this section, we carry out the stability and accuracy analyses for the upwind, centered, and CE/SE schemes.

For Fourier analysis, consider the 2D scalar advection equation:

$$u_t + au_x + bu_y = 0. \tag{7.1}$$

The mesh consists of square cells with cell widths $\Delta x = \Delta y = 1$ and cell centers at $(x, y) = (i, j)$. Next, since $\Delta x = \Delta y = 1$, set

$$\sigma_x = a\Delta t, \quad \text{and} \quad \sigma_y = b\Delta t. \tag{7.2}$$

For each time step, the data are advected by a spatial vector $(\sigma_x, \sigma_y)$, which is called the displacement vector (per time step) here.

In this section, as is typical for Fourier analysis, the slope employed is the average slope.

**7.1 Fourier analysis for a square mesh.** Assume that $a \geq 0$ and $b \geq 0$. Then the upwind scheme reproduces the exact solution when $(\sigma_x, \sigma_y)$ equals $(0, 0)$, $(1, 0)$, and $(0, 1)$. Loosely put, as can be seen from Fig. 7.1, when $(\sigma_x, \sigma_y) = (1, 0)$ the square centered at $B$ slides onto the square centered at $Q$.

The centered and the CE/SE schemes, on the other hand, do not reproduce the exact solution for any $(\sigma_x, \sigma_y)$ even when $(\sigma_x, \sigma_y) = (0, 0)$. Indeed,



Fig. 7.1. *The upwind scheme reproduces the exact solution when $(\sigma_x, \sigma_y)$ equals $(0, 0)$, $(1, 0)$, or $(0, 1)$.*



Fig. 7.2. *The centered and the CE/SE schemes do not reproduce the exact solution for any $(\sigma_x, \sigma_y)$.*

suppose $(\sigma_x, \sigma_y) = \overrightarrow{BQ} = (1, 0)$ (Fig. 7.2). Then, when updating the solution for the cell $Q$, at the beginning of the time step, the contribution from the cell $B$ is the reconstruction on $BKQE$ shown in Fig. 7.2. By the end of the time step, the reconstruction on $BEGHIK$ flows into the cell $Q$ via the fluxes across $BE$ and $BK$. Thus, when updating the solution at $Q$, the contribution from the cell $B$ is the reconstruction on $QEGHIKQ$, which is not the whole reconstruction cell. Therefore, we do not recover the value at $B$. Another way to see this fact is by setting all slopes to be zero. Using indices $i$ and $j$ and assuming $(\sigma_x, \sigma_y) = (1, 0)$, one obtains, after some algebra,

$$u_{0,0}^{n+1} = \tfrac{1}{4}(3u_{-1,0} + u_{0,-1} + u_{0,1} - u_{1,0}).$$

That is, we do not recover $u_{-1,0}$.

The fact that the solution recovers the exact value, loosely speaking, helps keep the error from becoming too big. For a triangular mesh, the situation is reversed as will be shown in the next subsection: the centered and CE/SE schemes recover the exact solution for certain values of $(\sigma_x, \sigma_y)$ while the upwind scheme does not.

The upwind scheme yields the following solution

at $i = 0$ and $j = 0$:

$$u_{0,0}^{n+1} = u_{0,0} +$$
$$\tfrac{1}{4}\sigma_x(\ u_{-2,0} + 5u_{-1,0} - 3u_{0,0} - u_{1,0}) +$$
$$\tfrac{1}{4}\sigma_y(\ u_{0,-2} + 5u_{0,-1} - 3u_{0,0} - u_{0,1}) +$$
$$\tfrac{1}{4}\sigma_x^2(u_{-2,0} - u_{-1,0} - u_{0,0} + u_{1,0}) +$$
$$\tfrac{1}{4}\sigma_y^2(u_{0,-2} - u_{0,-1} - u_{0,0} + u_{0,1}) +$$
$$\tfrac{1}{4}\sigma_x\sigma_y\,(2u_{-1,-1} - u_{-1,0} - u_{-1,1}$$
$$- u_{0,-1} + u_{0,1} - u_{1,-1} + u_{1,0}).$$
$$(7.3)$$

The amplification factor $\mathcal{A}$ is obtained by replacing $u_{i,j}$ in the above expression by $e^{I\,(iw_x + jw_y)}$, where $w_x$ and $w_y$ are the wave numbers in the $x$ and $y$ directions respectively, and $I = \sqrt{-1}$. The exact amplification factor is $e^{I(-\sigma_x w_x - \sigma_y w_y)}$.

The derivation of the amplification factors for the CE/SE scheme is similar to that of the 1D case except for the more involved algebra. If we carry along the interface values, we need to calculate the eigenvalues of a $2 \times 2$ matrix; if we carry along the slopes, we need to calculate the eigenvalues of a $3 \times 3$ matrix, but one of the three eigenvalues turns out to be identically zero.

The expressions for the amplification factors of the three schemes are lengthy and are omitted. However, the stability regions and the plots on accuracy comparison, which are generated by Mathematica, will be shown.

Figure 7.3 shows the stability regions for the upwind, centered, and CE/SE schemes. Note that the upwind scheme has the largest stability region, and the CE/SE, the smallest. Here, the two amplification factors of the CE/SE scheme yield the same stability region.

To compare errors among the three schemes, first, observe that the amplification factors are functions of $\sigma_x$, $\sigma_y$, $w_x$, and $w_y$. We wish to plot the errors for relatively smooth data; therefore, we fix small wave numbers $w_x$ and $w_y$ and plot the errors as functions of

$$\sigma = |(\sigma_x, \sigma_y)|$$

along the rays $\sigma_y = \text{constant} * \sigma_x$. (That is, for a fixed angle $\alpha$, $\sigma_x = \sigma\cos(\alpha)$ and $\sigma_y = \sigma\sin(\alpha)$.)

The phase error per time step is $\text{Arg}(\mathcal{A}) + (\sigma_x w_x + \sigma_y w_y)$. For the second-order schemes discussed here, this quantity is proportional to $O(w_x^3) + O(w_y^3)$ and is the leading error. The dissipation error per time step is $|\mathcal{A}| - 1$, which is proportional to $O(w_x^4) + O(w_y^4)$. The total error is $|\mathcal{A} - e^{I(-\sigma_x w_x - \sigma_y w_y)}|$.



Fig. 7.3. *Stability regions for a rectangular mesh. If the displacement vector based at the origin lies inside the stability region, the corresponding scheme is stable. The upwind scheme has the largest stability region, and the CE/SE, the smallest.*

For the next four plots, the wave numbers are $w_x = \pi/16$ and $w_y = \pi/8$. The results by the standard one-step Lax-Wendroff scheme are also included for ease of reference. Note that this scheme has difficulties dealing with shocks.

In Figs. 7.4, the flow is along the diagonal $y = x$, i.e., $\sigma_x = \sigma_y$. Fig. 7.4(a) shows the phase errors as functions of $\sigma$ (denoted by CFL on the plot). Note that the phase errors of the centered and the CE/SE schemes are nearly identical.

Figure 7.4(b) shows the dissipation errors as functions of $\sigma$. The upwind scheme has the smallest dissipation error, and the centered scheme, the largest.

Figure 7.4(c) shows the total errors as functions of $\sigma$. Here, the centered and CE/SE schemes have essentially the same errors while the upwind scheme has the smallest error.

Figure 7.5 shows the total errors averaged over all flow directions; again, $w_x = \pi/16$ and $w_y = \pi/8$.

The next four plots in Figs. 7.6 and 7.7 are identical to the previous four except that $w_x = \pi/32$ and $w_y = \pi/16$.

Figures 7.5 and 7.7 show that the total error of upwind scheme is considerably less than those of the centered and CE/SE schemes, and the latter two have essentially identical total errors.
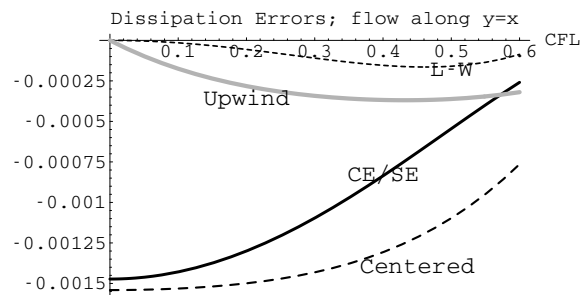
Note that the order of accuracy of the schemes can be verified using these figures. For example, the maximum phase error for the CE/SE scheme in Fig. 7.4(a) is about 0.0016; that in Fig. 7.6(a) is about 0.0002, which is 1/8 of 0.0016. In other
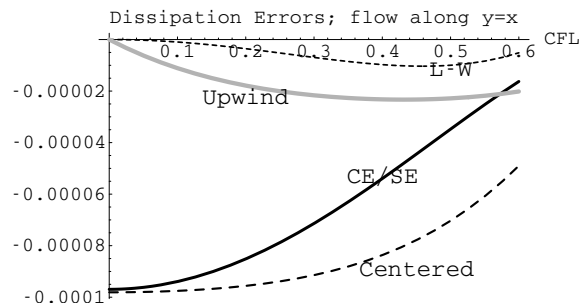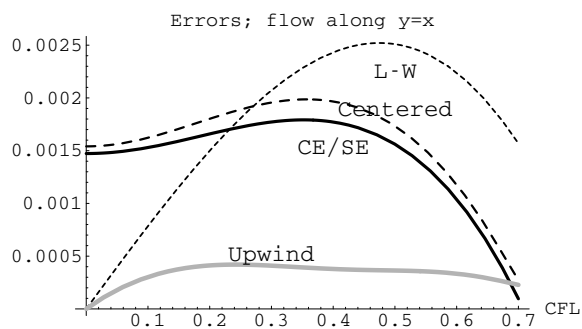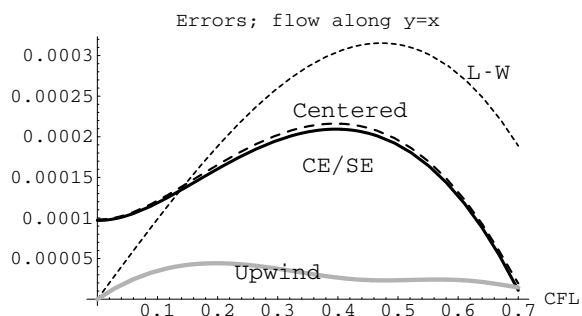
(a) *Phase errors*



(b) *Dissipation errors*



(c) *Total errors*

Fig. 7.4. *Errors as functions of $\sigma$ (denoted by CFL); flow along $y = x$; $w_x = \pi/16$ and $w_y = \pi/8$.*



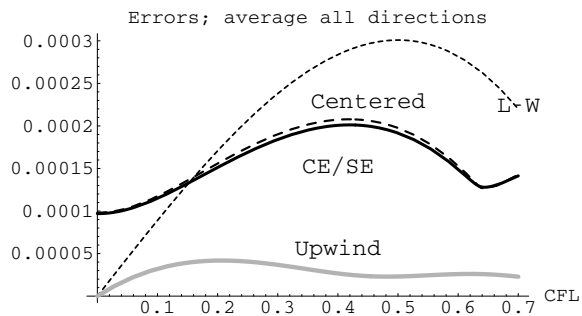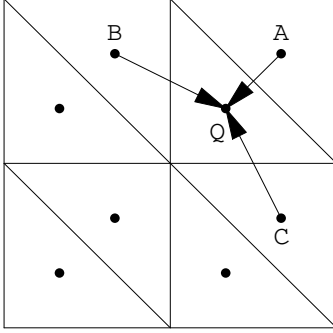Fig. 7.5. *Total errors; average for all flow directions; again $w_x = \pi/16$ and $w_y = \pi/8$.*



(a) *Phase errors*



(b) *Dissipation errors*



(c) *Total errors*

Fig. 7.6. *Errors as functions of $\sigma$; flow along $y = x$; $w_x = \pi/32$ and $w_y = \pi/16$.*



Fig. 7.7. *Total errors; average for all flow directions; again $w_x = \pi/32$ and $w_y = \pi/16$.*

Fig. 7.8. *The upwind scheme does not reproduce the exact solution for any $(\sigma_x, \sigma_y)$ other than the obvious value of $(0,0)$.*



Fig. 7.9. *The centered and the CE/SE schemes reproduce the exact solution if $(\sigma_x, \sigma_y)$ equals one of the three vectors $\overrightarrow{AQ}$, $\overrightarrow{BQ}$, or $\overrightarrow{CQ}$.*

words, when we double the number of mesh points in each direction (quadruple the total number of points), the wave numbers reduce by a factor of 2, the phase error after one time step reduces by a factor of 8, and the dissipation error, a factor of 16 (as can be seen in Figs. 7.4(b) and 7.6(b)). Similar observations also hold for the other schemes.

Note that the errors for the square-mesh case here are similar to those for the 1D case in §3.

## 7.2 Fourier analysis on a triangular mesh.
Here, each square in Fig. 7.1 is cut into two triangles along the diagonal from the northwest to the southeast corners as shown in Fig. 7.8.

The first question is when do the schemes recover the exact solution? Here, the upwind scheme recovers the exact solution only for the obvious case of $(\sigma_x, \sigma_y) = (0,0)$. The reason is that under the translation $\overrightarrow{AQ}$, the triangle $A$ does not match with



Fig. 7.10. *The numbering of the triangles for Fourier analysis.*

the triangle $Q$ as can be seen from Fig. 7.8. Similar observations hold for the translations $\overrightarrow{BQ}$ and $\overrightarrow{CQ}$.

The centered and CE/SE schemes, on the other hand, do not recover the exact solution when $(\sigma_x, \sigma_y) = (0,0)$, but they recover the exact solution when $(\sigma_x, \sigma_y)$ equals one of the three displacement vectors $\overrightarrow{AQ}$, $\overrightarrow{BQ}$, or $\overrightarrow{CQ}$, as shown in Fig. 7.9. Indeed, suppose $(\sigma_x, \sigma_y) = \overrightarrow{AQ}$. Then, loosely put, the reconstruction cell $A$ slides onto the reconstruction cell $Q$. More precisely, at the end of the time step, through the start-off value on $AEQG$ and the two fluxes across $AE$ and $AG$, the contribution from the cell $A$ is the reconstruction on the whole reconstruction cell $A$. Therefore, after one time step, the solution at $Q$ is identical to the data at $A$: we recover the exact solution. This property of recovering the exact solution helps the centered and the CE/SE schemes gain back some accuracy compared with the upwind scheme as will be shown later.

The next few examples using the first-order schemes are simple, but they convey the behavior of schemes on a triangular mesh.

For first-order accuracy, the slopes are set to zero, and both the centered and CE/SE schemes reduce to the L-F scheme. If the data are on a plane, this scheme recovers the exact solution. Indeed, the solution at $Q$ is obtained by applying the displacement to the plane through the data at $A$, $B$ and $C$ shown in Fig. 7.9.

The upwind scheme, on the other hand, does not recover the plane. Indeed, suppose the flow angle is between 0 and $\pi/4$ so that the upwind cell to cell $Q$ is cell $B$ (Fig. 7.8). Then, when updating the solution at $Q$, the solution involves the data at only $Q$ and $B$. These two pieces of data cannot recover

a plane (we need three). What happens is that typically there is a cancellation of errors (for the fluxes) so that a scheme using a piecewise constant reconstruction reproduces the exact solution if the data are on a plane. Here, the cancellation of errors did not occur.

The next question is: what does the piecewise-constant upwind scheme do to a data which is on a plane? It turns out that the scheme propagates the plane and turns it into a 'bumpy plane'. As an example, suppose $\sigma_y = 0$ and $\sigma_x > 0$. Then (see Fig. 7.10), for the piecewise-constant upwind scheme,

$$u_{L,i,j}^{n+1} = u_{L,i,j} + 2\sigma_x(u_{R,i-1,j} - u_{L,i,j}), \quad (7.4a)$$

and

$$u_{R,i,j}^{n+1} = u_{R,i,j} + 2\sigma_x(u_{L,i,j} - u_{R,i,j}). \quad (7.4b)$$

If the data are on the plane $u = x$, then (see Fig. 7.10),

$$u_{L,i,j} = i - 1/6, \quad \text{and} \quad u_{R,i,j} = i + 1/6. \quad (7.5)$$

Next, consider the following data of odd-even type noise

$$e_{L,i,j} = -1/12, \quad \text{and} \quad e_{R,i,j} = 1/12. \quad (7.6)$$

Suppose the data is a 'bumpy plane' obtained by superimposing (7.6) on (7.5):

$$u_{L,i,j} = (i - 1/6) - 1/12, \qquad u_{R,i,j} = (i+1/6)+1/12. \quad (7.7)$$

Then the solution by the piecewise-constant upwind scheme for this initial data after one time step is, by (7.4),

$$u_{L,i,j}^{n+1} = u_{L,i,j} - \sigma_x, \quad (7.8a)$$

and

$$u_{R,i,j}^{n+1} = u_{R,i,j} - \sigma_x. \quad (7.8b)$$

For the initial data (7.5), i.e., the data that are on the plane $u = x$, the exact solution is also given by (7.8).

Thus, the bumpy initial data (7.7) is preserved by the piecewise-constant upwind scheme in the sense that the solution is given by first propagating the plane $u = x$ exactly, and then superimposing the odd-even noise (7.6) back on. In fact, if we start with the data on the plane, the solution by the piecewise-constant upwind scheme turns it into a 'bumpy plane' of type (7.7) after 8 iterations (for $\sigma_x = .2$). Note that if we calculate the value at the center of each square by averaging the 'bumpy'

values at the centroids of the two triangles, we get back the data on the plane, i.e., if the triangles are combined in pairs to form squares, then the cancellation effect again takes place.

The above observation is consistent with the result of Fourier analysis. To carry out this analysis for the triangular mesh case, we must pair up the downward and upward pointing triangles so that the solution at each square looks the same as the solution at any other square. Consequently, each scheme has two amplification factors: principal and spurious. The principal eigenfunctions of the (piecewise constant and piecewise linear) upwind schemes turn out to be somewhat bumpy, while those for the centered and CE/SE schemes remain smooth. For the piecewise-constant upwind scheme, the spurious component is of order $O(h)$, but this error does not accumulate as we march to a fixed final time, and the piecewise constant upwind scheme retains its first-order accuracy. A similar remark also holds for the second-order upwind scheme (the spurious component is of order $O(h^2)$).

To calculate the amplification factor, define (see Fig. 7.8)

$$\mathbf{u}_{i,j} = \begin{pmatrix} u_{L,i,j} \\ u_{R,i,j} \end{pmatrix}. \quad (7.9)$$

Then, the solution vector can be written as

$$\mathbf{u}_{i,j}^{n+1} = \mathbf{C}_{-1,0}\mathbf{u}_{i-1,j} + \mathbf{C}_{0,-1}\mathbf{u}_{i,j-1} + \mathbf{C}_{0,0}\mathbf{u}_{i,j}$$
$$+\mathbf{C}_{1,0}\mathbf{u}_{i+1,j} + \mathbf{C}_{0,1}\mathbf{u}_{i,j+1} + \cdots$$
$$(7.10)$$

Here each $\mathbf{C}_{i,j}$ is a $2 \times 2$ matrix. The two eigenvalues of the matrix $e^{-Iw_x}\mathbf{C}_{-1,0}+e^{-Iw_y}\mathbf{C}_{0,-1}+\cdots$ are the amplification factors of the corresponding scheme. The eigenvalue which approximates the exact amplification factor $e^{-I(\sigma_x w_x - \sigma_y w_y)}$ is the principal amplification factor, and the other, the spurious.

Note that because the CE/SE scheme is odd-even decoupled, instead of pairing up the triangles, we only need to take two time steps. After calculating the eigenvalues (associated with two time steps), we need to take their square roots to obtain the amplification factor (per time step).

The expressions for the amplification factors of the three schemes are very lengthy and are omitted here. The stability regions and accuracy comparisons are shown below, however.

Figures 7.11(a), (b), and (c) show the stability regions for the upwind, centered, and CE/SE schemes respectively. Here, the time step size for the upwind scheme is reduced considerably due to the stability

limit along the two axes shown in Fig. 7.11(a). In practice, however, the maximum time step size is not restricted too much compared with those of the centered and CE/SE schemes.
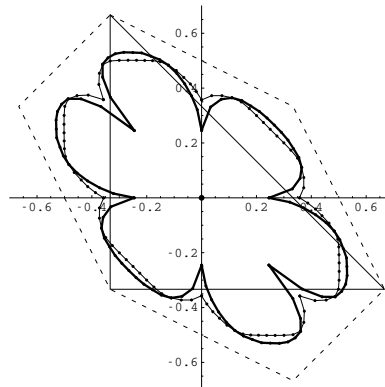
For the next four plots, Figs. 7.12(a)–(c) and 7.13, the wave numbers are $w_x = \pi/16$ and $w_y = \pi/8$. The four plots in Figs. 7.14(a)–(c) and 7.15 are the same as the previous four except that $w_x = \pi/32$ and $w_y = \pi/16$. Observe that as in the rectangular-mesh case, the errors are of appropriate order.

Note that as shown by Figs. 7.13 and 7.15, the advantage in accuracy of the upwind scheme over the centered and CE/SE schemes is considerably less here compared to the rectangular-mesh case in Figs. 7.5 and 7.7. The upwind scheme also has the drawbacks of a somewhat large spurious component and a relatively small time step size (due to the flower-shape stability region).

## 7.3. Accuracy comparison between rectangular and triangular meshes.

We can also compare the errors between the rectangular- and the triangular-mesh cases. First, observe that because the triangular mesh is obtained by slicing each square into two triangles, the number of triangles are twice that of squares. To have the same number of cells as the triangles, the squares must have a width of $1/\sqrt{2}$. Thus, for the same number of cells, if the scheme is equally accurate between the square and the triangular mesh, then the phase errors in Figs. 7.5–7.7 must be $(\sqrt{2})^3 \approx 2.8$ times the corresponding ones in Figs. 7.13–7.15, and the dissipation errors, $(\sqrt{2})^4 = 4$ times.

For the upwind scheme, the maximum error in Fig. 7.7 is about 0.00004; that in Fig. 7.15 is also about 0.00004; i.e., there is no improvement in the triangular mesh case even though the number of cells is twice that of the quadrilateral mesh. Thus, for the same number of cells, the upwind scheme is more efficient on a rectangular mesh.

Another drawback for the upwind scheme in the case of a triangular mesh is a rather large spurious component. For a smooth data the spurious component is of order $O(h^2)$ whereas the error of the principal amplification factor per time step is $O(h^3)$; here, $h$ is the smallest edge length. To reach a fixed final time, we need a number of time steps proportional to $1/h$, and the error by the principal amplification factor accumulates to $O(h^2)$. The spurious component is eventually damped out, and the corresponding error $O(h^2)$ does not accumulate. As a result, the (piecewise linear) upwind scheme



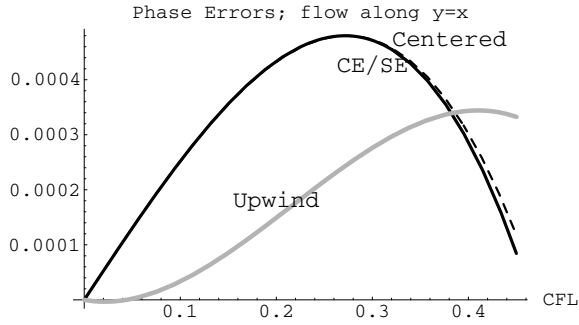(a) *Upwind scheme, stability region, triangular mesh.*



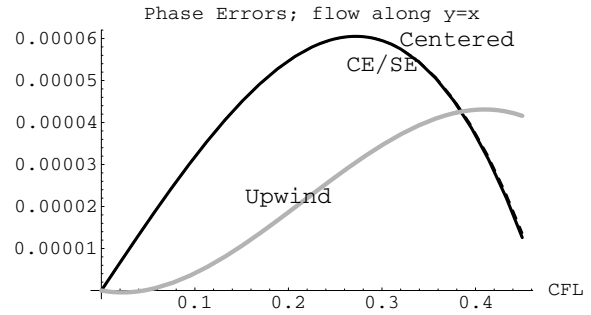(b) *Centered scheme (i.e., extended N-T or coupled CE/SE).*
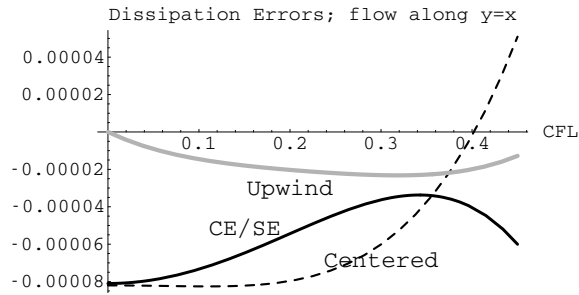


(c) *CE/SE scheme.*

Fig. 7.11. *Triangular-mesh stability regions. If the displacement vector based at the origin lies inside the stability region, the scheme is stable. The dark curve is produced by the principal amplification factor; the light curve, the spurious one. The triangular cell and the reconstruction cell (hexagon) for the centered and CE/SE schemes are shown.*
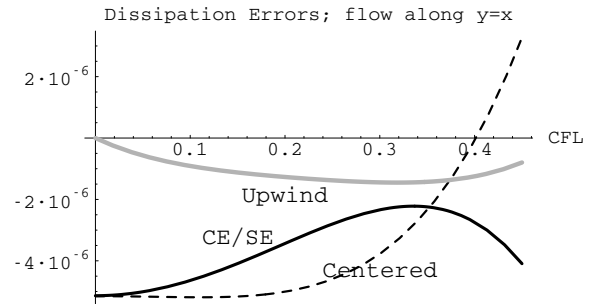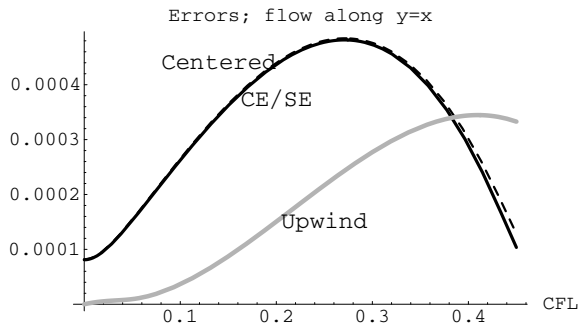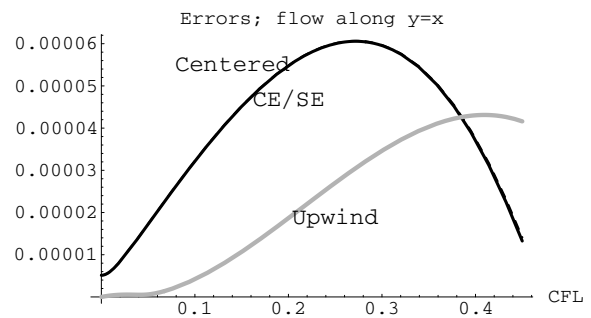
(a) *Phase errors; triangular mesh*



(b) *Dissipation errors*



(c) *Total errors*

Fig. 7.12. *Errors for the triangular-mesh case; flow along* $y = x$; $w_x = \pi/16$ *and* $w_y = \pi/8$.



(a) *Phase errors; triangular mesh*



(b) *Dissipation errors*



(c) *Total errors*

Fig. 7.14. *Errors for the triangular-mesh case; flow along* $y = x$; $w_x = \pi/32$ *and* $w_y = \pi/16$.



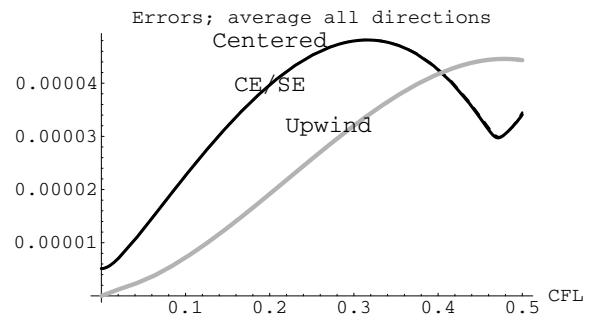Fig. 7.13. *Total errors; average for all flow directions; again* $w_x = \pi/16$ *and* $w_y = \pi/8$.



Fig. 7.15. *Total errors; average for all flow directions; again* $w_x = \pi/32$ *and* $w_y = \pi/16$.

28

retains its second-order accuracy for the case of a triangular mesh.

For the centered and CE/SE schemes, the maximum error in Fig. 7.7 is about 0.0002; that in Fig. 7.15 is about 0.00005; the improvement is by a factor of 4, more than the expected factor of 2.8. The maximum dissipation error in Fig. 7.7 is the value at $\sigma = \text{CFL} = 0$, which is about 0.0001; that in Fig. 7.15 is about 0.000005, which is an improvement by a factor of 20, more than the expected factor of 4. Thus, between a quadrilateral and a triangular mesh with the same number of cells, the centered and CE/SE schemes are more efficient on the triangular mesh.

Notice that the above analysis is linear. It provides useful information on stability and accuracy. However, the analysis no longer holds when there is a solid wall or when a limiter function is employed. In such cases, we resort to numerical experiments.

**8. Numerical results.** Results for the 1D as well as the 2D quadrilateral and triangular meshes are shown below.

**8.1 Results for the 1D case.** In the following 1D numerical examples, $\gamma = 1.4$ and the CFL number is 0.9. Unless otherwise stated, the number of mesh points is 200.

The first numerical test, used by Sod (1978), is the Riemann problem

$$(\rho_L, u_L, p_L) = (1, 0, 1),$$

$$(\rho_R, u_R, p_R) = (0.125, 0, 0.1).$$

The final time is $t = 0.2$. The solid line represents the exact solution. The results are shown in Fig. 8.1.

The second problem, due to Shu and Osher (1989), has several extrema in the smooth regions. In the interval $\;5 \le x \le 5$, a moving Mach 3 shock interacts with sine waves in density as described by the following initial conditions:

$$(\rho, u, p) = \begin{cases} (3.857, 2.629, 10.333) & \text{if } x < \;4, \\ (1 + 0.2 \sin 5x, 0, 1) & \text{otherwise.} \end{cases}$$

The final time is $t = 1.8$. Figure 8.2 shows the results with 800 mesh points. Since the exact solution is not known, the solid line represents the solution for 1600 cells via a uniformly second-order accurate upwind scheme described in Huynh (1995a).

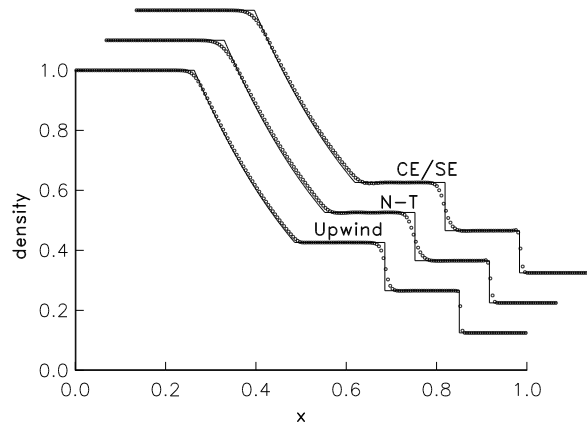Note that for these two problems, the upwind results are more accurate than the N-T and CE/SE
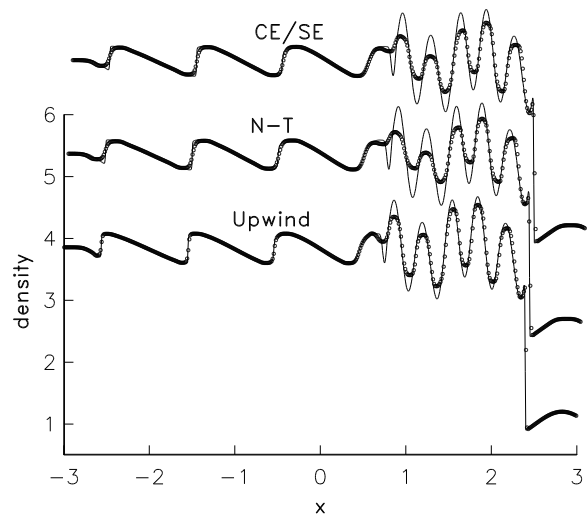


Fig. 8.1 *Sod's problem.*



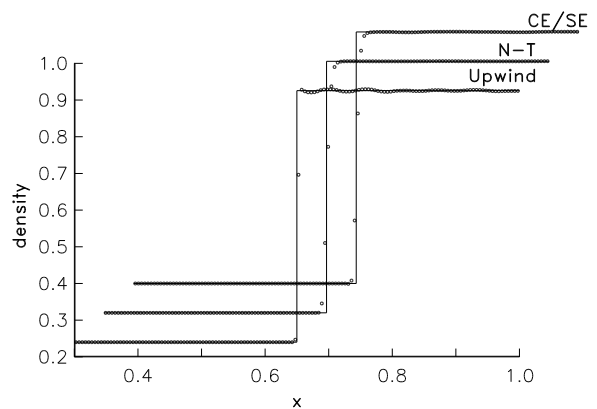Fig. 8.2 *Shu and Osher's problem (800 mesh points).*



Fig. 8.3 *Slow-moving shock problem. Final time = 2.2*

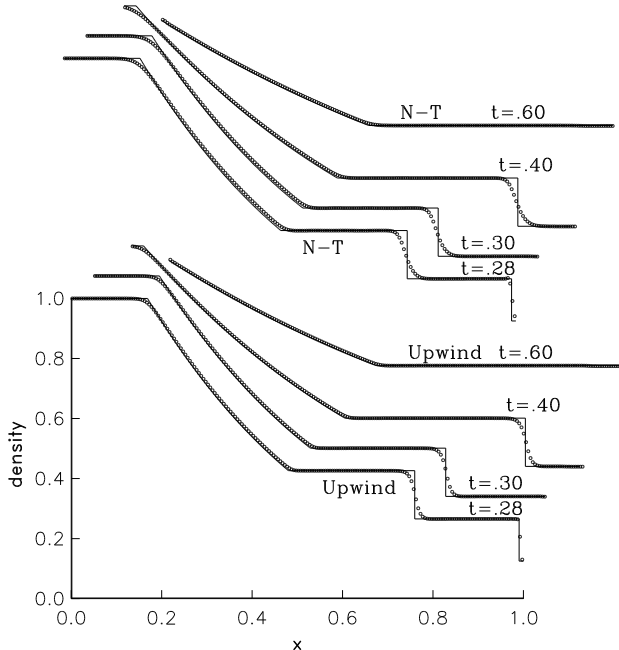Fig. 8.4 *Straightforward nonreflecting boundary condition works for the CE/SE as well as the upwind and N-T schemes.*



Fig. 8.5. *Oblique shock;* $80 \times 20$ *rectangular and* $40 \times 10 \times 4$ *triangular meshes.*

results. Between the N-T and CE/SE solutions, the formers are very slightly more dissipative.

The next test is the slow-moving shock problem for which the upwind scheme generates oscillations. The initial condition is a velocity of 0.05 superimposed on a Mach 3 normal shock:

$$(\rho_L, u_L, p_L) = (0.24, 2.2238, 0.09),$$

$$(\rho_R, u_R, p_R) = (0.9256, 0.6137, 0.9299).$$

Figure 8.3 shows the results for the final time $t = 3$. Here, the upwind solution oscillates. Note that the entropy condition (4.49) employed helps reduce oscillations considerably compared with a typical implementation of Roe's scheme or (4.48). The N-T and CE/SE solutions have no oscillations. An upwind step that generates no oscillations for this problem can be found in Wada and Liou (1997).

The last 1D test concerns nonreflecting boundary conditions, which let waves leave the domain of computation. The boundary condition employed below is an obvious one: the value at the ghost cell is set to be that of the immediate (interior) neighbor, and the slope is set to zero (see also Chang et al. 1999). The results for Sod's problem at time .28, .30, .40, and .60 are shown in Fig. 8.4. Here, since the CE/SE solutions are again essentially the
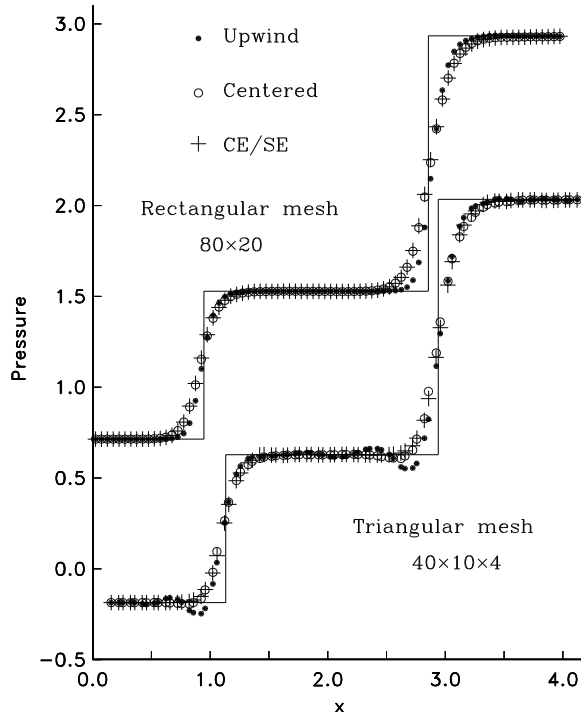
same as the N-T solutions, they are omitted. The nonreflecting boundary condition works well for all schemes: there are no reflecting waves.

**8.2 Results for the 2D case.** As explained at the beginning of §7, the extended N-T scheme using the CE/SE approach to extension is called the centered scheme here—as opposed to the upwind scheme. If the slopes are carried along, the corresponding extension is called the CE/SE scheme.

For 2D test problems, the solutions are obtained by the author's codes except in the case of the CE/SE scheme on a triangular mesh where the CE/SE code, kindly provided by Dr. Xiao-Yen Wang, is employed. Because a unified CFL number for both the triangular and quadrilateral meshes is not readily available, we provide information on the time steps, which are generally chosen to be close to the largest time step possible.

The schemes discussed here are designed to solve unsteady problems. The 2D test cases chosen below, however, are steady-state problems because the exact solutions are either known or have certain features which are useful for the purpose of comparison.

The first 2D test is the oblique shock problem on the domain $[0,4] \times [0,1]$ (Yee et al. 1983). The

30

condition at the inflow boundary is $(\rho, u, v, p) = (1., 2.9, 0., 1./1.4)$ and that at the top boundary, $(1.7, 2.6193, \quad .5063, 1.5282)$. The rectangular mesh is $80 \times 20$. The triangular mesh is obtained by slicing each rectangle (along the two diagonals) of a $40 \times 10$ rectangular mesh into four triangles. Thus, each mesh has 1600 cells. The solutions are shown in Fig. 8.5; here, pressure at cell centers at a fixed $y$ location are plotted; for the rectangular-mesh case $y = .475$; for the triangular-mesh case $y = .45$. Note that all schemes smear the oblique shocks. The upwind solution is slightly less smeared than the centered and CE/SE solutions for the rectangular-mesh case, but it oscillates for the triangular-mesh case, while the other two solutions do not. Also note that again, the solutions by the centered and CE/SE schemes are nearly identical. For these two schemes, the triangular-mesh solutions are also nearly the same as the rectangular-mesh solutions.

Notice that the CE/SE solutions shown in Fig. 8.5 here are considerably different from the solutions by Chang et al. (1999) and Zhang et al. (2002), where the shocks are resolved by only one point. The reason is that the mesh is different. The result for the $80 \times 20$ rectangular mesh here has also been verified by Dr. Xiao-Yen Wang using her CE/SE code (private communication).

Next, we discuss the time step sizes and convergence. In the quadrilateral-mesh case, for the upwind scheme, with $\Delta t = 0.0106$, the solution converges to machine accuracy; it blows up with $\Delta t = 0.0107$. For the centered scheme, with $\Delta t = 0.0100$, the solution converges to machine accuracy; it does not converge for $0.0101 \le \Delta t \le 0.0104$; and it blows up with $\Delta t = 0.0105$. For the CE/SE scheme, the solution blows up at $\Delta t = 0.0108$, and does not converge well for most $\Delta t$, e.g., with $\Delta t = 0.006$, the maximum residual goes down from .35E0 to .79E-4 after 3000 iterations (the other two schemes converge to machine accuracy of E-15 under the same conditions). It appears that because the scheme is odd-even decoupled, it does not converge well. Note that for the centered and CE/SE schemes, the residual is calculated using the data at even time levels only (as discussed at the end of §1.3). In the triangular-mesh case, for the upwind and centered schemes, with $\Delta t = 0.0078$, the solutions converge to machine accuracy. For the CE/SE scheme, with $\Delta t = 0.0075$, the maximum residual goes down three orders of magnitude.

Figure 8.6(a), (b), and (c) show the solutions on an unstructured triangular mesh. (The mesh for this case as well as that in Fig. 8.9 were kindly provided by Dr. Philip Jorgenson.) Here, the upwind solution is slightly less smeared compared with the other two. Also note that this upwind solution (Fig. 8.6(a)) improves considerably compared with the oscillatory solution in Fig. 8.5; the minimum pressure here is .698 whereas that for the $40 \times 10 \times 4$ triangular mesh, .612 (and that for the $y = .45$ slice in Fig. 8.5, .653); the exact minimum pressure is .714.
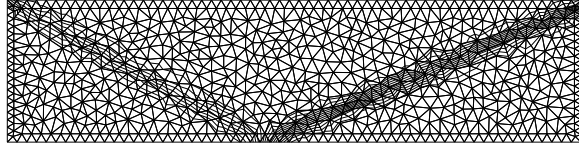
Since solutions by the CE/SE schemes are essentially the same as those by the centered schemes, they are omitted from here on.

The next problem is the nozzle problem (Verhoff 1985). Here, the domain is $[0, 3] \times [0, 1]$ except that for $x$ in $[1, 2]$, the bottom wall is defined by $y = .1 \sin^2(x \quad 1)\pi$. The flow is subsonic, and the boundary conditions are standard. At inflow, we set the total pressure $p_0 = 1$; stagnation speed of sound $a_0 = 1$; and $v_0 = 0$; pressure, however, is extrapolated from the interior. At outflow, pressure is set to .843, which results in an inflow Mach number of .5; the other three variables are extrapolated from the interior. Since there are no shocks in the solution, the average slope is employed.
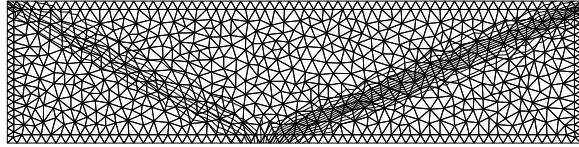
The $64 \times 16$ quadrilateral mesh and the solutions are shown in Fig. 8.7(a) and (b) respectively. Here, in Fig. 8.7(b), the line is the solution by the centered scheme on a $128 \times 32$ mesh. Observe that the solution by the upwind scheme is slightly more accurate: the peak is higher, and the valleys are lower than those of the centered scheme.

For the same problem, the $32 \times 8 \times 4$ triangular mesh and the solutions are shown in Fig. 8.8(a) and (b) respectively. Here, in Fig. 8.8(b), the solution by centered scheme is slightly more accurate and maintains symmetry better than that of the upwind scheme. Also note the accuracy at the lower wall near the outflow boundary.
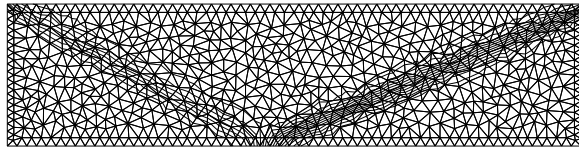
The last test is the transonic flow over a circular bump. The boundary conditions are the same as those of the above nozzle problem except that the outflow pressure is .736. The mesh has 2466 cells. For the upwind scheme, the time step is .011. After 5000 iterations, the residual decreases from .938E-1 to .918E-4, but after another 24000 iterations, the residual is still .532E-4, i.e., the solution does not converge well. A time step of 0.008 does not improve convergence. For the centered scheme, the time step is .013 (.014 does not converge). After 10000 iterations, the residual decreases from .157E0 to .822E-8, and after another 10000 iterations, the residual reaches machine error of E-14.

(a) *Upwind scheme*



(b) *Centered scheme*



(c) *CE/SE scheme*

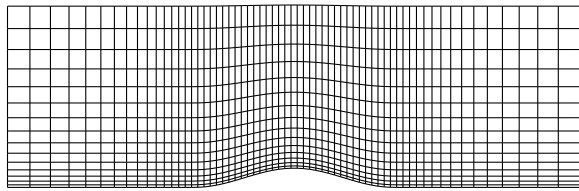Fig. 8.6. *Pressure contours, oblique shock.*



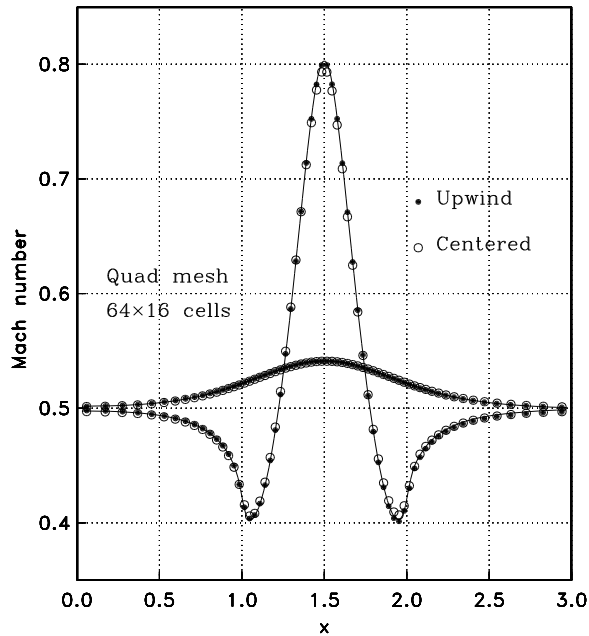Fig. 8.7(a) *Quad mesh (64 × 16); nozzle problem.*



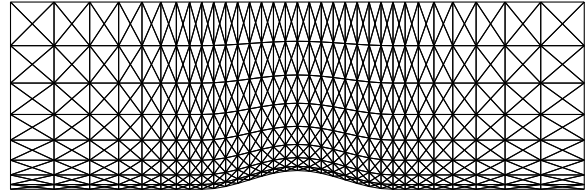Fig. 8.7(b) *Quad-mesh results; nozzle problem.*



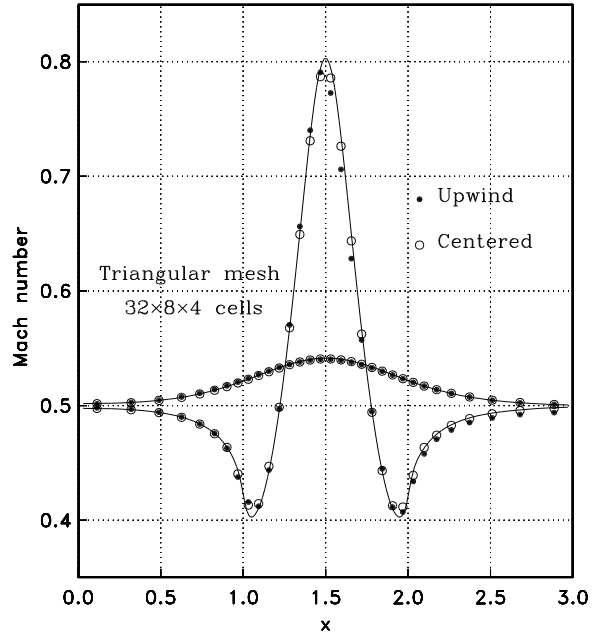Fig. 8.8(a) *Triangular mesh (32 × 8 × 4); nozzle problem.*
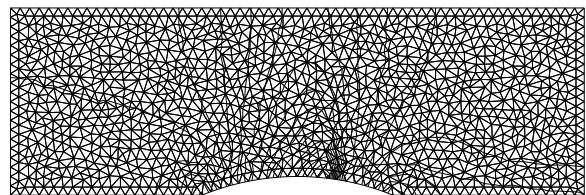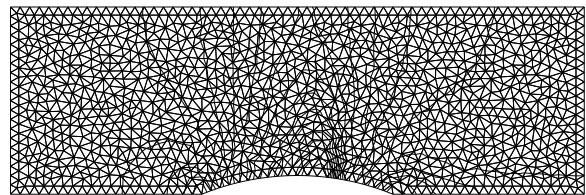


Fig. 8.8(b) *Triangular-mesh results; nozzle problem.*



(a) *Upwind scheme*



(b) *Centered scheme*

Fig. 8.9. *Mach contours, circular bump on a wall.*

**9. Conclusions and discussion.** The numerical tests here are clearly preliminary; more tests need to be performed.

In summary, the second-order accurate upwind, the Lax-Friedrichs type second-order accurate centered (or the extended N-T scheme using the CE/SE approach to extension), and the CE/SE ($\epsilon = 1/2$) schemes were presented in a framework that facilitates their comparison. These schemes all use piecewise-linear reconstructions (MUSCL interpolants). The key difference is that the centered and CE/SE schemes avoid the upwind step by employing reconstruction cells that are roughly twice as big as the original cells.

The schemes employed have several properties in common. They are all of finite-volume type, conserve fluxes locally and globally, are equally dimensional splitting (or non-splitting), can capture shocks, are only first-order accurate near extrema when the slope is defined by a weighted average or a limiter function, and can handle nonreflecting boundary conditions.

Fourier stability and accuracy analyses were carried out for these schemes for the standard 1D and the 2D quadrilateral mesh cases. In the nonstandard case of a triangular mesh, the triangles were paired up for the analyses of the upwind and the centered schemes. Among the three schemes, on a quadrilateral mesh, the upwind scheme is more accurate. On a triangular mesh, however, accuracy among the three schemes appears to be comparable. Comparing the same scheme on the two meshes with the same number of cells, the upwind scheme is more efficient on a quadrilateral mesh, and the centered and CE/SE schemes, on a triangular mesh.

The centered scheme has numerous advantages over the CE/SE scheme: it produces essentially the same solution, runs faster, converges better, and requires only one third of the storage for flow variables. The capability of adjusting dissipation of the CE/SE schemes, however, has not been fully explored.

Compared with the upwind scheme, the centered scheme is more stable, conceptually simpler (the geometry involved is slightly more complicated, however), and requires less storage (for an unstructured mesh). Besides providing an alternative to the upwind choice, the centered scheme can be very useful when the equations become complicated, and the upwind model either has not been derived or is no longer appropriate (the equations are no longer hyperbolic).

The development, analysis, and comparison of these schemes with viscous terms remain to be explored. (A discussion of the CE/SE scheme for viscous flows can be found in Chang et al. (1995).)

### References.

P. Arminjon, D. Stanescu, and M.C. Viallon, *A two-dimensional finite volume extension of the Lax-Friedrichs and Nessyahu-Tadmor schemes for compressible flows*, Proc. of the 6th Int. Symp. on CFD, Lake Tahoe, Nevada, September 4–8, 1995.

P. Arminjon, M.C. Viallon, A. Madrane, *A Finite Volume Extension of the Lax-Friedrichs and Nessyahu-Tadmor Schemes for Conservation Laws on Unstructured Grids*, IJCFD, Vol. 9, (1997) pp1–22.

T. J. Barth and D. C. Jespersen., *The Design and Application of Upwind Schemes on Unstructured Meshes*, AIAA-89-0366, 27th AIAA Aerospace Sciences Meeting, January 09–12, 1989, Reno, NV.

S.-C. Chang, *The method of space-time conservation element and solution element—A new approach for solving the Navier-stokes and Euler equations*, J. Comp. Phys., 119 (1995), pp. 295-324.

S.-C. Chang, X.-Y. Wang, C.-Y. Chow, and A. Himansu, *The method of space-time conservation element and solution element—development of a new implicit solver*, Proceedings of the Ninth International Conference on Numerical Methods in Laminar and Turbulent Flows, July 10–14, 1995, Atlanta, GA.

S.-C. Chang, X.-Y. Wang, C.-Y. Chow, *The space-time conservation element and solution element method: A new high resolution and genuinely multidimensional paradigm for solving conservation laws*, J. Comp. Phys., 156 (1999), pp. 89-136.

G. Cook, Jr., *High Accuracy Capture of Shock Fronts Using the Method of Space-Time Conservation Element and Solution Element*, AIAA-99-1008, 37th AIAA Aerospace Sciences Meeting and Exhibit, January 11–14, 1999, Reno, NV.

R. Courant, E. Isaacson, and M. Rees, *On the solution of nonlinear hyperbolic differential equations bu finite differences*, Comm. Pure Appl. Math. 5 (1952), pp. 243–49.

J. E. Fromm, *A method for reducing dispersion in convective difference schemes*, J. Comp. Phys., 3 (1968), pp. 176–189.

S. K. Godunov, *A finite difference method for the numerical computation of discontinuous solutions*

*of the equations of fluid dynamics*, Mat. Sb., 47 (1959), pp. 357–393.

A. Harten, B. Engquist, S. Osher, and S. R. Chakravarthy, *Uniformly high-order accurate essentially nonoscillatory schemes.* III, J. Comp. Phys., 71 (1987), pp. 231–303.

C. Hirsch, *Numerical Computation of Internal and External Flows*, Vol. 2, John Wiley & Sons, New York, 1990, 691 pp.

H. T. Huynh, *Accurate upwind methods for the Euler equations*, SIAM J. Numer. Anal., 32 (1995a), pp. 1565–1619.

H. T. Huynh, *Accurate upwind schemes for the Euler equations*, AIAA 95-1737, AIAA 12th Computational Fluid Dynamics Conference, San Diego, CA (1995b).

A. Jameson, W. Schmidt, and E. Turkel, *Numerical solutions of the Euler equations by finite-volume methods using Runge-Kutta time-stepping*, AIAA Paper 81-1259

G.-S. Jiang, D. Levy, C.-T. Lin, S. Osher, and E. Tadmor, *High-Resolution Non-oscillatory Central Schemes with Non-Staggered Grids for Hyperbolic Conservation Laws*, SIAM J. Numer. Anal., 35 (1998), pp. 2147–2168.

G.-S. Jiang and E. Tadmor, *Nonoscillatory Central Schemes for Multidimensional Hyperbolic Conservation Laws*, SIAM J. Sci. Comput., 19 (1998), pp. 1892–1917.

P. D. Lax, *Weak solutions of nonlinear hyperbolic equations and their numerical computation*, Commun. Pure Appl. Math., 7 (1954) 159–193.

M.-S. Liou and C. J. Steffen, Jr., *A new flux splitting scheme*, J. Comp. Phys., 107 (1993), pp. 23–39.

H. Nessyahu and E. Tadmor, *Non-oscillatory central differencing for hyperbolic conservation laws*, J. Comp. Phys., 87 (1990), pp. 408–463.

S. Osher, *Riemann solvers, the entropy condition, and difference approximations*, SIAM J. Numer. Anal., 21 (1984), pp. 217–235.

P. L. Roe, *Approximate Riemann Solvers, Parameter Vectors, and Difference Schemes*, J. Comp. Phys., 43 (1981), pp. 357–72.

P. L. Roe, *An Introduction to Numerical Methods Suitable for the Euler Equations*, Introduction to Computational Fluid Dynamics, January 24–28, 1983, Lecture Series 1983–01, von Karman Institute for Fluid Dynamics.

P. L. Roe, *Characteristic-based schemes for the Euler equations*, Ann. Rev. Fluid Mech., 18 (1986), pp. 337–365.

P. L. Roe, *A survey of upwind differencing techniques*, Lecture Notes in Physics, Vol. 323 (Springer-Verlag, New York/Berlin, 1989), pp. 69–78.

C-W. Shu and S. Osher, *Efficient implementation of essentially non-oscillatory shock-capturing schemes*, II, J. Comp. Phys., 83 (1989), pp. 32–78.

G. A. Sod, *A survey of several finite difference methods for systems of non-linear hyperbolic conservation laws*, J. Comp. Phys., 27 (1978), pp. 1–31.

G. D. van Albada, B. van Leer, and W. W. Roberts, Jr., *A comparative study of computational methods in cosmic gas dynamics*, Astronom. and Astrophys., 108 (1982), pp. 76–84.

B. van Leer, *Towards the ultimate conservative difference scheme. IV. A new approach to numerical convection*, J. Comp. Phys., 23 (1977), pp. 276–298.

V. Venkatakrishnan, *A Perspective on Unstructured Grid Flow Solvers*, ICASE Report No. 95–3, February 1995, pp. 1–37.

A. Verhoff, *Modeling of computational and solid surface boundary conditions for fluid dynamics calculations*, AIAA 85-1496, AIAA 7th Computational Fluid Dynamics Conference, (1985).

Y. Wada and M.-S. Liou, *An accurate and robust flux splitting scheme for shock and contact discontinuities*, SIAM J. Sci. Comput., 18 (1997), pp. 633–657.

X.-Y. Wang and S.-C. Chang, *A 2D non-splitting unstructured triangular mesh Euler solver based on the space-time conservation element and solution element method*, Comp. Fluid Dynamics Journal, Vol. 8, no. 2 (1999) pp. 309–325.

H. C. Yee, R. F. Warming, and A. Harten, *Implicit Total Variation Diminishing (TVD) Schemes for Steady-State Calculations*, AIAA Paper 83-1902 (1983).

Z.-C. Zhang, S. T. Yu, and S.-C. Chang, *A space-time conservation element and solution element method for solving the two- and three-dimensional unsteady Euler equations using quadrilateral and hexahedral meshes*, J. Comp. Phys., 175 (2002), pp. 168–199.