

On a Large Time-Step High Resolution Scheme*

By Ami Harten

Abstract. This paper presents a class of new second-order accurate $(2K + 3)$ -point explicit schemes for the computation of weak solutions of hyperbolic conservation laws, that are total-variation-diminishing under a CFL restriction of K . These highly nonlinear schemes are obtained by applying a nonoscillatory first-order accurate $(2K + 1)$ -point scheme to a modified flux. The so-derived second-order accurate schemes achieve high resolution, while retaining the robustness of the original first-order accurate scheme.

1. Introduction. In this paper, we discuss numerical approximations to weak solutions of the initial value problem (IVP) for hyperbolic systems of conservation laws

$$(1.1) \quad u_t + f(u)_x = 0, \quad u(x, 0) = \phi(x), \quad -\infty < x < \infty,$$

where $f(u)$ is continuously differentiable and $\phi(x)$ is a BV function of compact support.

We consider finite-difference approximations that are obtained by $(2K + 1)$ -point explicit schemes in conservation form

$$(1.2a) \quad v_j^{n+1} = v_j^n - \lambda(\bar{f}_{j+1/2} - \bar{f}_{j-1/2}), \quad \lambda = \tau/h,$$

where

$$(1.2b) \quad \bar{f}_{j+1/2} = \bar{f}(v_{j-K+1}^n, \dots, v_{j+K}^n).$$

Here $v_j^n = v(jh, n\tau)$ and \bar{f} , the numerical flux, is consistent with the flux $f(u)$ in the following sense:

$$(1.2c) \quad \bar{f}(u, \dots, u) = f(u).$$

The Courant-Friedrichs-Lewy (CFL) theorem states that the maximal time-step for a stable $(2K + 1)$ -point explicit scheme is restricted by

$$(1.3) \quad \lambda \max_{u, k} |a_k(u)| \leq K,$$

where $a_k(u)$ are the eigenvalues of the Jacobian matrix $A(u) = \partial f / \partial u$.

Recently, LeVeque [5], [6] has experimented with a large time-step first-order accurate scheme for the computation of discontinuous solutions of (1.1). His results show that using the scheme with $K \leq 6$ yields a rather adequate description of

Received March 8, 1983; revised June 26, 1985.

1980 *Mathematics Subject Classification.* Primary 65P05, 35L65, 76L05.

*This work was supported under the National Aeronautics and Space Administration under NASA Contract NAS1-15810 while the author was in residence at the Institute for Computer Applications in Science and Engineering, NASA Langley Research Center, Hampton, Virginia 23665.

propagation, and even interaction of discontinuities. Similar results in the scalar case were obtained by Brennier [1].

In [2] we describe a rather general technique to convert a nonoscillatory first-order accurate scheme into a second-order accurate one. Applying this technique to a 3-point first-order accurate scheme which is stable under a CFL restriction of 1 ($K = 1$ in (1.3)), results in a 5-point second-order accurate scheme that is stable under the same CFL restriction.

Encouraged by the results of LeVeque and Brennier, we apply the technique of [2] to a nonoscillatory $(2K + 1)$ -point first-order accurate scheme which has the maximal CFL restriction K , with $K > 1$, to obtain a nonoscillatory $(2K + 3)$ -point second-order accurate scheme that is stable under the same CFL restriction. The primary motivation for the development of schemes of this nature is the possible gain in computational efficiency due to the fact that a single application of such a large time-step scheme is less costly than several applications of the scheme with $K = 1$.

The enlargement of the numerical domain of dependence requires a special treatment of boundaries. Fortunately, for the type of schemes considered in this paper, this can be accomplished in a rather simple manner.

2. First-Order Accuracy (Scalar Case). In this section we consider the IVP (1.1) for a *single* conservation law. In the scalar case, the total variation in x of the solution to the IVP (1.1) is diminishing in time. Therefore, we consider finite-difference approximations in conservation form (1.2) that are also total-variation diminishing (TVD), i.e.,

$$(2.1a) \quad \text{TV}(v^{n+1}) \leq \text{TV}(v^n),$$

where

$$(2.1b) \quad \text{TV}(v) = \sum_{j=-\infty}^{\infty} |\Delta_{j+1/2} v|;$$

here, and throughout this paper, we use the notation

$$(2.2) \quad \Delta_{j+1/2} b = b_{j+1} - b_j$$

for any mesh-function b .

First let us consider the constant coefficient case $f(u) = au$, $a = \text{constant}$. A *linear* $(2K + 1)$ -point scheme is consistent with (1.1), if of the form

$$(2.3a) \quad v_j^{n+1} = \sum_{k=-K}^K B_k(\nu) v_{j+k}^n, \quad \nu = \lambda a,$$

$$(2.3b) \quad \sum_{k=-K}^K B_k(\nu) \equiv 1.$$

If (2.3a) is a first (or higher)-order accurate scheme, then

$$(2.3c) \quad \sum_{k=-K}^K k B_k(\nu) \equiv -\nu.$$

Subtracting (2.3a) at j from (2.3a) at $j + 1$, we get

$$(2.4a) \quad \Delta_{j+1/2} v^{n+1} = \sum_{k=-K}^K B_k(\nu) \Delta_{j+k+1/2} v^n.$$

Hence, (2.3) is TVD if, and only if,

$$(2.4b) \quad B_k(\nu) \geq 0 \quad \text{for } |k| \leq K.$$

Condition (2.4b) also implies that (2.3a) is a monotone scheme; consequently, a linear TVD scheme is only first-order accurate (see [2]). Second-order accurate TVD schemes are essentially nonlinear in the sense that they are *nonlinear* even in the constant coefficient case.

We turn now to describe a $(2K + 1)$ -point scheme that is TVD under the maximal CFL restriction (1.3).

To construct this scheme we operate K times with the 3-point upstream-differencing scheme

$$(2.5a) \quad v_j^{n+1} = [L(\nu)v^n]_j = v_j^n - \nu^- \Delta_{j+1/2} v^n - \nu^+ \Delta_{j-1/2} v^n,$$

where

$$(2.5b) \quad \nu^\pm = \frac{1}{2}(\nu \pm |\nu|),$$

and then consider it to be a single step of a $(2K + 1)$ -point scheme, i.e.,

$$(2.6a) \quad v^{n+1} = [L(\nu/K)]^K v^n.$$

Expressing L in (2.5a) in terms of translation operators $T^k v_j \equiv v_{j+k}$, we get in (2.6a)

$$(2.6b) \quad v_j^{n+1} = [1 - |\nu|/K + (\nu^+/K)T^{-1} - (\nu^-/K)T]^K v_j^n.$$

Using the binomial expansion in (2.6b), and the fact that $\nu^+ \nu^- \equiv 0$, we rewrite (2.6) as (2.3) with

$$(2.6c) \quad B_{\pm k}(\nu) = \begin{cases} b_k(\mp \nu^\mp / K) & \text{for } K \geq k \geq 1, \\ b_0(|\nu|/K) & \text{for } k = 0, \end{cases}$$

where

$$(2.6d) \quad b_k(x) = \begin{cases} \binom{K}{k} x^k (1-x)^{K-k} & \text{for } k \geq 1, \\ (1-x)^K & \text{for } k = 0. \end{cases}$$

Obviously, $B_k(\nu) \geq 0$ for $|\nu| \leq K$.

The scheme (2.6) can also be written in the conservation form (1.2):

$$(2.7a) \quad v_j^{n+1} = v_j^n - \lambda(\bar{f}_{j+1/2} - \bar{f}_{j-1/2}),$$

with

$$(2.7b) \quad \lambda \bar{f}_{j+1/2} = \frac{\lambda}{2}(f_j + f_{j+1}) - \sum_{k=-K+1}^{K-1} C_k(\nu) \Delta_{j+k+1/2} v^n,$$

where

$$(2.7c) \quad C_{\pm k}(\nu) = \begin{cases} c_k(\mp \nu^\mp / K) & \text{for } 1 \leq k \leq K - 1, \\ c_0(|\nu|/K) & \text{for } k = 0, \end{cases}$$

and

$$(2.7d) \quad c_l(x) = \begin{cases} -\binom{K}{l} x^l \sum_{k=1}^{K-l} \binom{K-l}{k} \frac{kl}{(k+l-1)(k+l)} (-x)^k & \text{for } l \geq 1, \\ \frac{K}{2} x & \text{for } l = 0; \end{cases}$$

f_j denotes $f(v_j)$.

The relation between b_l (2.6d) and c_l (2.7d) is given by

$$(2.8a) \quad b_l = c_{l-1}(x) + c_{l+1}(x) - 2c_l(x) + \begin{cases} 0 & \text{for } K \geq l \geq 2, \\ -\frac{K}{2}x & \text{for } l = 1, \\ 1 & \text{for } l = 0, \end{cases}$$

where we have used the convention that $c_{-1}(x) = c_K(x) = 0$.

It follows from (2.8a) and (2.6d) that

$$(2.8b) \quad 2 \sum_{l=1}^{K-1} c_l(x) = \sum_{l=1}^{K-1} l^2 b_l(x) = K(K-1)x^2.$$

Next we extend the first-order accurate scheme (2.7)–(2.8) to nonlinear conservation laws by

$$(2.9a) \quad v_j^{n+1} = v_j^n - \lambda(\bar{f}_{j+1/2} - \bar{f}_{j-1/2}),$$

$$(2.9b) \quad \lambda \bar{f}_{j+1/2} = \frac{\lambda}{2}(f_j + f_{j+1}) - \sum_{k=-K+1}^{K-1} C_k(v_{j+k+1/2}) \Delta_{j+k+1/2} v^n,$$

where $v_{j+1/2}$ is the following mean-value CFL number

$$(2.9c) \quad v_{j+1/2} = \lambda a_{j+1/2}, \quad a_{j+1/2} = \Delta_{j+1/2} f / \Delta_{j+1/2} v,$$

and $C_k(x)$ are (2.7c).

We now show that the scheme (2.9) is TVD under the CFL restriction

$$(2.10) \quad \max_j |v_{j+1/2}| \leq K.$$

To see that, we subtract (2.9) at j from (2.9) at $j + 1$ to get

$$(2.11a) \quad \Delta_{j+1/2} v^{n+1} = \sum_{k=-K}^K B_k(v_{j+k+1/2}) \Delta_{j+k+1/2} v^n.$$

Here $B_k(x)$ are (2.6c); this follows directly from the relation (2.8a). Hence $B_k(v_{j+k+1/2}) \geq 0$ under the CFL restriction (2.10), and

$$(2.11b) \quad \sum_{k=-K}^K B_k(v_{j+1/2}) \equiv 1.$$

Taking absolute value of (2.11a) and using the triangle inequality, we get

$$(2.11c) \quad |\Delta_{j+1/2} v^{n+1}| \leq \sum_{k=-K}^K B_k(v_{j+k+1/2}) |\Delta_{j+k+1/2} v^n|.$$

Summing (2.11c) from $j = -\infty$ to $j = +\infty$, we get by shifting indices and using (2.11b) that

$$\begin{aligned} \text{TV}(v^{n+1}) &= \sum_{j=-\infty}^{\infty} |\Delta_{j+1/2} v^{n+1}| \leq \sum_{j=-\infty}^{\infty} \sum_{k=-K}^K B_k(v_{j+k+1/2}) |\Delta_{j+k+1/2} v^n| \\ &= \sum_{j=-\infty}^{\infty} |\Delta_{j+1/2} v^n| \sum_{k=-K}^K B_k(v_{j+1/2}) = \sum_{j=-\infty}^{\infty} |\Delta_{j+1/2} v^n| = \text{TV}(v^n). \end{aligned}$$

The scheme (2.9) with $C_k(x)$ defined by (2.7c) admits a stationary *expansion* shock as its steady solution (see [4]). To prevent this violation of the entropy inequality, and at the same time make the numerical flux a smooth function of its arguments, we replace $|v|$ in (2.5b) and (2.7c) by

$$(2.12a) \quad Q(v) = \begin{cases} \frac{1}{2}(v^2/\epsilon + \epsilon) & \text{for } |v| < \epsilon, \\ |v| & \text{for } |v| \geq \epsilon, \end{cases} \quad \epsilon = \lambda\delta,$$

where δ has the dimension of velocity (see [4] for more details). Denoting

$$(2.12b) \quad \mu_{\pm}(v) = \frac{1}{2}[Q(v/K) \pm v/K],$$

we take $C_k(v)$ in (2.9b) to be

$$(2.12c) \quad C_{\pm k}(v) = \begin{cases} c_k(\mu_{\mp}(v)) & \text{for } 1 \leq k \leq K-1, \\ \frac{K}{2}Q\left(\frac{v}{K}\right) & \text{for } k = 0. \end{cases}$$

The coefficients $B_k(v)$ in (2.11a) become

$$(2.13) \quad B_{\pm k}(v) = \begin{cases} b_k(\mu_{\mp}(v)) & \text{for } 1 \leq k \leq K, \\ b_0(\mu_{-}(v)) + b_0(\mu_{+}(v)) - 1 & \text{for } k = 0, \end{cases}$$

where $b_k(x)$ are (2.6d). We note that if $\epsilon = 0$ in (2.12) then $B_k(v)$ in (2.13) become identical to (2.6c). Hence, for ϵ sufficiently small (namely $\epsilon \leq 2K(1 - 2^{-1/K})$), $B_{\pm k}(v) \geq 0$ for $|v| \leq K$, and the scheme (2.9) with (2.12) is also TVD under the CFL restriction (2.10).

The scheme (2.9) is a first-order accurate approximation to (1.1). Its modified equation, i.e., the equation it approximates to second-order accuracy, is

$$(2.14a) \quad u_t + f(u)_x = \frac{1}{\lambda} \hat{g}_x,$$

where

$$(2.14b) \quad \hat{g} = h\sigma(v)u_x,$$

$$(2.14c) \quad \sigma(v) = \frac{K}{2} \left\{ Q\left(\frac{v}{K}\right) \left[1 + \frac{K-1}{2} Q\left(\frac{v}{K}\right) \right] - \frac{K+1}{2} \left(\frac{v}{K}\right)^2 \right\};$$

clearly, $\sigma(v) \geq 0$ for $|v| \leq K$.

In the next section, we convert the $(2K + 1)$ -point first-order accurate TVD scheme (2.9) into a $(2K + 3)$ -point second-order accurate TVD scheme. Our technique is based on the following observations:

(i) The scheme (2.9) is TVD for any flux, provided that the CFL restriction (2.10) is satisfied.

(ii) Consider the application of the scheme (2.9) to the flux $f + g/\lambda$, where $g = \hat{g} + O(h^2)$, and \hat{g} is defined in (2.14). By the definition of a modified equation, the numerical scheme (2.9) is a second-order accurate approximation to its solutions. However, the modified equation (2.14) corresponding to (2.9) with $f + g/\lambda$ satisfies $u_t + f_x = O(h^2)$. Hence (2.9), with the flux $f + g/\lambda$, is a second-order accurate approximation to (1.1).

3. Second-Order Accuracy (Scalar Case). In this section, we convert the scheme (2.9) into a second-order accurate one, by applying it to a modified flux $f + g/\lambda$, as follows:

$$(3.1a) \quad v_j^{n+1} = v_j^n - \lambda(\bar{f}_{j+1/2} - \bar{f}_{j-1/2}),$$

$$(3.1b) \quad \begin{aligned} \lambda \bar{f}_{j+1/2} &= \frac{\lambda}{2}(f_j + f_{j+1}) + \frac{1}{2}(g_j + g_{j+1}) \\ &\quad - \sum_{k=-K+1}^{K-1} C_k(\nu + \gamma)_{j+k+1/2} \Delta_{j+k-1/2} v^n, \end{aligned}$$

where $C_k(x)$ is (2.12),

$$(3.1c) \quad g_j = s \cdot \max\left\{0, \min\left(\sigma_{j+1/2} |\Delta_{j+1/2} v|, s \cdot \sigma_{j-1/2} \Delta_{j-1/2} v\right)\right\}$$

and

$$(3.1d) \quad \gamma_{j+1/2} = \Delta_{j+1/2} g / \Delta_{j+1/2} v;$$

here $s = \text{sgn}(\Delta_{j+1/2} v)$, and $\sigma_{j+1/2} = \sigma(\nu_{j+1/2})$; cf., (2.14c).

In [2], we study g_j and show that wherever v is smooth

$$(3.2a) \quad (i) \quad \frac{1}{2}(g_j + g_{j+1}) = h\sigma(\nu) v_{x|_{x_{j+1/2}}} + O(h^2) = g_{j+1/2} + O(h^2),$$

$$(3.2b) \quad (ii) \quad \gamma_{j+1/2} \Delta_{j+1/2} v = g_{j+1} - g_j = O(h^2)$$

and that if $|\Delta_{j+1/2} v| > 0$, then

$$(3.2c) \quad (iii) \quad |\gamma_{j+1/2}| = |\Delta_{j+1/2} g / \Delta_{j+1/2} v| \leq \sigma(\nu_{j+1/2}).$$

We turn now to analyze the order of accuracy of the scheme (3.1). To do so, let us assume that v is smooth in a neighborhood of $x_{j+1/2}$ and expand the numerical flux (3.1b) around $x_{j+1/2}$ up to $O(h^2)$ terms. Since $C_k(x)$ in (2.12) are Lipschitz-continuous functions, we get from (3.2b) that

$$\left| [C_k(\nu_{i+1/2} + \gamma_{i+1/2}) - C_k(\nu_{i+1/2})] \Delta_{i+1/2} v \right| \leq K |\gamma_{i+1/2} \Delta_{i+1/2} v| = O(h^2);$$

therefore,

$$\begin{aligned} C_k(\nu_{j+k+1/2} + \gamma_{j+k+1/2}) \Delta_{j+k+1/2} v &= C_k(\nu_{j+k+1/2}) \Delta_{j+k+1/2} v + O(h^2) \\ &= C_k(\nu) h v_{x|_{x_{j+k+1/2}}} + O(h^2) = h C_k(\nu) v_{x|_{j+1/2}} + O(h^2). \end{aligned}$$

Using the above relation and (3.2a), we get

$$(3.3a) \quad \lambda \bar{f}_{j+1/2} = \left[\lambda f + h\sigma(\nu) v_x - h \sum_{k=-K+1}^{K-1} C_k(\nu) v_x \right]_{j+1/2} + O(h^2).$$

Since by (2.7c) and (2.12c)

$$(3.3b) \quad \begin{aligned} &\sum_{k=-K+1}^{K-1} C_k(\nu) \\ &= \frac{1}{2} Q(\nu) + \frac{1}{2} K(K-1) \left\{ \left(\frac{1}{2K} [Q(\nu) - \nu] \right)^2 + \left(\frac{1}{2K} [Q(\nu) + \nu] \right)^2 \right\}, \end{aligned}$$

we get by (2.14c) that

$$\sigma(\nu) - \sum_{k=-K+1}^{K-1} C_k(\nu) = -\frac{1}{2} \nu^2,$$

and, therefore,

$$(3.3c) \quad \lambda \bar{f}_{j+1/2} = [\lambda f - \frac{1}{2} h v^2 v_x]_{j+1/2} + O(h^2).$$

Let us express the $O(h^2)$ term in (3.3c) as

$$\beta_{j+1/2} h^2 + O(h^3)$$

and examine the smoothness of the coefficient β . We see that except for the critical points $v_x = 0$ where g_j in (3.1c) is not differentiable, the coefficient β is a Lipschitz function of x ,

$$\beta_{j+1/2} - \beta_{j-1/2} = O(h),$$

and, therefore,

$$\begin{aligned} v_j^{n+1} &= v_j^n - \lambda (\bar{f}_{j+1/2} - \bar{f}_{j-1/2}) = \left[v - \tau f_x + \frac{\tau^2}{2} (a^2 v_x)_x \right]_j^n + O(h^3) \\ &= \left[v + \tau v_t + \frac{\tau^2}{2} v_{tt} \right]_j^n + O(h^3) = v(x_j, t_n + \tau) + O(h^3). \end{aligned}$$

This shows that, except at critical points, the scheme (3.1) is second-order accurate in the sense of local truncation error.

Next, we show that the scheme (3.1) is TVD under the CFL restriction (2.10). Since (3.1) is (2.9) applied to $f + g/\lambda$, it follows that it is TVD, provided that the CFL restriction

$$(3.4a) \quad \max_j |v_{j+1/2} + \gamma_{j+1/2}| \leq K$$

is satisfied. The inequality (3.2c) enables us to replace the condition (3.4a) by a CFL restriction on v alone, i.e.,

$$(3.4b) \quad \max_j [|v_{j+1/2}| + \sigma(v_{j+1/2})] \leq K.$$

It is easy to see that $|v| + \sigma(v)$ does not have a local maximum in $(-K, K)$, consequently, if $|v| \leq K$, then

$$(3.4c) \quad |v| + \sigma(v) \leq K + \sigma(K) = K.$$

We conclude that (3.1) is TVD subject to the original CFL restriction (2.10).

Remarks. (1) $g_j = g(v_{j-1}, v_j, v_{j+1})$ and consequently

$$\gamma_{j+1/2} = \gamma(v_{j-1}, v_j, v_{j+1}, v_{j+2}).$$

Although the numerical flux (3.1b) has the same functional form as (2.9b), its dependence on the numerical characteristic speed γ enlarges its support by 2 points. Therefore, a $(2K + 1)$ -point first-order accurate scheme (2.9) is converted into a $(2K + 3)$ -point second-order accurate scheme (3.1).

(2) g_j in (3.1c), which is $O(h)$, can be modified by the addition of $\tilde{g}_j = O(h^2)$; this leaves the scheme second-order accurate, and \tilde{g} can be so chosen to leave the CFL restriction unchanged. The purpose of such a modification is to introduce artificial compression for contact-discontinuities and shocks (see Section 5 and [2] for more details).

(3) We have shown that the scheme (3.1) is second-order accurate in the sense of local truncation error, except at critical points where it remains first-order accurate. Numerical experiments in [10] show that the cumulative error in smooth problems

behaves in a similar manner: It is second-order except at local extrema where it degenerates locally to first-order. Consequently, the global cumulative error in smooth problems is second-order in the L_1 -norm, but only first-order in the maximum norm.

4. Systems of Conservation Laws. In this section we describe how to extend the new second-order accurate scheme (3.1) to hyperbolic systems of conservation laws. Our extension technique is a somewhat generalized version of the procedure suggested by Roe in [8].

Let $A(u)$ be the Jacobian matrix of $f(u)$ in (1.1),

$$(4.1a) \quad A(u) = f_u(u),$$

and denote its eigenvalues and right-eigenvectors (columns) by $a^k(u)$, $r^k(u)$, $1 \leq k \leq m$, respectively. Since the set of eigenvectors is complete (by the hyperbolicity assumption), the matrix

$$(4.1b) \quad R(u) = (r^1(u), \dots, r^m(u))$$

is invertible. The rows $l^1(u), \dots, l^m(u)$ of $R^{-1}(u)$ constitute the biorthonormal system of left-eigenvectors; thus

$$(4.1c) \quad l^i r^j = \delta_{ij},$$

and

$$(4.1d) \quad R^{-1}AR = \Lambda, \quad \Lambda_{ij} = a^i \delta_{ij}.$$

Next, we define characteristic variables w with respect to the state u by

$$(4.2a) \quad w^k = l^k(u)u, \quad 1 \leq k \leq m,$$

or

$$(4.2b) \quad u = \sum_{k=1}^m w^k r^k(u).$$

In the constant coefficient case, (1.1) decouples into m scalar equations for the characteristic variables

$$(4.3) \quad w_t^k + a^k w_x^k = 0, \quad a^k = \text{constant},$$

and therefore $\text{TV}(w^k)$ is diminishing in time for all $1 \leq k \leq m$.

To retain this property, we extend our scheme to systems of conservation laws by applying (3.1) to each of the *scalar* locally defined characteristic variables, as follows:

Let $v_{j+1/2} = V(v_j, v_{j+1})$ be some symmetric average of v_j and v_{j+1} , and denote by $a_{j+1/2}^k$, $r_{j+1/2}^k$ and $l_{j+1/2}^k$ the respective quantities of $A(v_{j+1/2})$. Let $\alpha_{j+1/2}^k$ be the component of $(v_{j+1} - v_j)$ in the k th characteristic direction, i.e.,

$$(4.4a) \quad \alpha_{j+1/2}^k = l_{j+1/2}^k \Delta_{j+1/2} v$$

or

$$(4.4b) \quad \Delta_{j+1/2} v = \sum_{k=1}^m \alpha_{j+1/2}^k r_{j+1/2}^k,$$

and define

$$(4.5a) \quad v_j^{n+1} = v_j^n - \lambda(\bar{f}_{j+1/2} - \bar{f}_{j-1/2}),$$

$$(4.5b) \quad \lambda \bar{f}_{j+1/2} = \frac{\lambda}{2}(f_j + f_{j+1}) + \sum_{k=1}^m r_{j+1/2}^k \left\{ \frac{1}{2}(g_j^k + g_{j+1}^k) - \sum_{l=-K+1}^{K-1} C_l(v^k + \gamma^k)_{j+l+1/2} \alpha_{j+l+1/2}^k \right\},$$

where

$$(4.5c) \quad g_j^k = s \cdot \max\left[0, \min\left(\sigma_{j+1/2}^k |\alpha_{j+1/2}^k|, s \cdot \sigma_{j-1/2}^k \alpha_{j-1/2}^k\right)\right],$$

$$s = \text{sgn}(\alpha_{j+1/2}^k),$$

$$(4.5d) \quad \gamma_{j+1/2}^k = \Delta_{j+1/2} g^k / \alpha_{j+1/2}^k,$$

and

$$(4.5e) \quad v_{j+1/2}^k = \lambda a_{j+1/2}^k;$$

here $\alpha_{j+1/2}^k$ is (4.4), $\sigma_{j+1/2}^k = \sigma(v_{j+1/2}^k)$ (2.14c) and $C_l(x)$ are (2.12).

It is easy to see that in the constant coefficient case indeed

$$(4.6) \quad \text{TV}(w^k) = \sum_{j=-\infty}^{\infty} |\alpha_{j+1/2}^k|$$

is diminishing in time for all $1 \leq k \leq m$.

This property does not depend on the particular forms of averaging $v_{j+1/2} = V(v_j, v_{j+1})$. However, if we want the scheme (4.5) for $m = 1$ to be identical to (3.1), we have to choose (4.5e) so that $v_{j+1/2}$ is the mean value CFL number (2.9c). This can be accomplished by taking the eigenvalues $a_{j+1/2}^k$ and the eigenvectors $l_{j+1/2}^k$ and $r_{j+1/2}^k$ in (4.5) to be those of $A(v_j, v_{j+1/2})$, where $A(u, v)$ is Roe's mean-value Jacobian (see [9]). This matrix satisfies

$$(4.7a) \quad (i) \quad f(v) - f(u) = A(u, v)(v - u),$$

$$(4.7b) \quad (ii) \quad A(u, u) = A(u),$$

$$(iii) \quad A(u, v) \text{ has real eigenvalues and a complete set of eigenvectors.}$$

(See [2] and [3] for more details.)

We note that in increasing K , the only extra computational work is in the calculation of

$$(4.8) \quad \sum_{l=-K+1}^{K-1} C_l(v^k + \gamma^k)_{j+l+1/2} \alpha_{j+l+1/2}^k$$

in (4.5b). Since (4.8) is a scalar quantity, and $C_l(x)$ are rather easy-to-compute polynomials of degree K (see table), one can expect to gain in computational efficiency by taking $K > 1$.

The enlargement of the numerical domain of dependence of the scheme with a larger K requires a special treatment of boundaries. In the Appendix we show that this can be accomplished in a rather simple manner.

5. Numerical Examples. In this section we present some numerical examples in order to demonstrate the performance of the scheme (4.5) for various choices of K .

We consider a Riemann problem for the Euler equations of gas dynamics

$$(5.1a) \quad u_t + f(u)_x = 0, \quad u(x, 0) = \begin{cases} u_L, & x < 0, \\ u_R, & x > 0, \end{cases}$$

$$(5.1b) \quad u = (\rho, m, E)^T, \quad f(u) = (m, m^2/\rho + P, m(E + P)/\rho)^T,$$

$$(5.1c) \quad P = (\gamma - 1)(E - \frac{1}{2}m^2/\rho).$$

Here ρ, m, E and P are the density, momentum, total energy, and pressure, respectively; we take $\gamma = 1.4$.

As in [2] we choose

$$(5.2) \quad u_L = (0.445, 0.311, 8.928)^T, \quad u_R = (0.5, 0, 1.4275)^T.$$

In Figures 1 and 2, we present calculations with 140 cells and $h = 0.1$. The solid line shows the exact solution; the circles are the discrete values of the numerical solution. In all these calculations we have used Roe's linearization (see [2] for more details).

In Figure 1, we present calculations performed by the scheme (4.5) with various choices of K ; these are shown at about the same physical time. Figure 1(a) shows the results for $K = 1$ with a CFL number 0.9 after 111 time steps. Figure 1(b) shows $K = 2$ with a CFL number 1.8 after 55 time steps. Figure 1(c) shows $K = 4$ with a CFL number 3.6 after 27 time steps. In Figure 1(d), we show $K = 6$ with a CFL number of 5.4 after 18 time steps.

We observe that the shock in Figure 1(d) has propagated about 52 cells in 18 time steps. Unlike other large time-step methods, the deterioration in quality with increasing CFL number exhibits itself in the form of excessive smearing rather than the creation of spurious oscillations. It is interesting to note that it is the rarefaction wave that produces such oscillations.

In Figure 2 we repeat the calculation in Figure 1, but now we add an artificial compression/rarefaction term to the scheme (4.5) by modifying g_j^k in (4.5c) to be

$$(5.3a) \quad (1 + \mu_j^k \theta_j^k) \cdot g_j^k,$$

where θ_j^k is an automatic switch of the form

$$(5.3b) \quad \theta_j^k = |\alpha_{j+1/2}^k - \alpha_{j-1/2}^k| / (|\alpha_{j+1/2}^k| + |\alpha_{j-1/2}^k|).$$

$\alpha_{j+1/2}^k$ is (4.4a). θ is $O(h)$ in regions of smoothness and always $0 \leq \theta \leq 1$. Taking $\mu > 0$ in (5.3a) corresponds to applying artificial compression, while $\mu < 0$ has the effect of artificial rarefaction. Here we take $\mu_j^k \equiv 1$ in the linearly-degenerate characteristic field and

$$(5.3c) \quad \mu_j^k = -\text{sgn}(a_{j+1/2}^k - a_{j-1/2}^k)$$

in the genuinely nonlinear characteristic fields (see [2]).

Comparing Figure 1 to Figure 2 we see that the addition of artificial rarefaction acts to suppress the spurious oscillations at the expense of excessive rounding of the corners of the rarefaction wave. The addition of artificial compression improves the resolution of both the shock and the contact discontinuity.

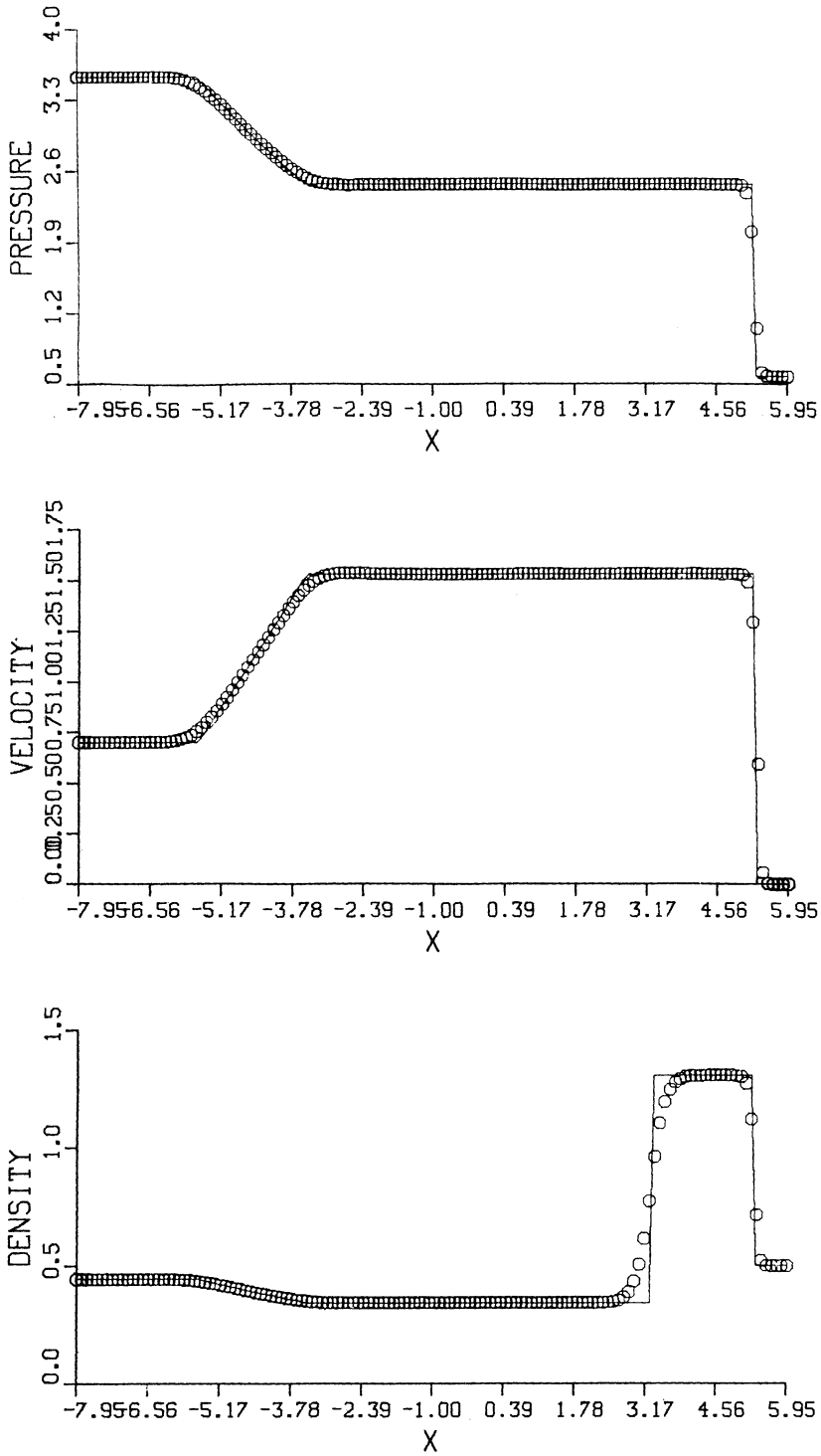


FIGURE 1(a)
K = 1 with CFL 0.9 after 111 time steps

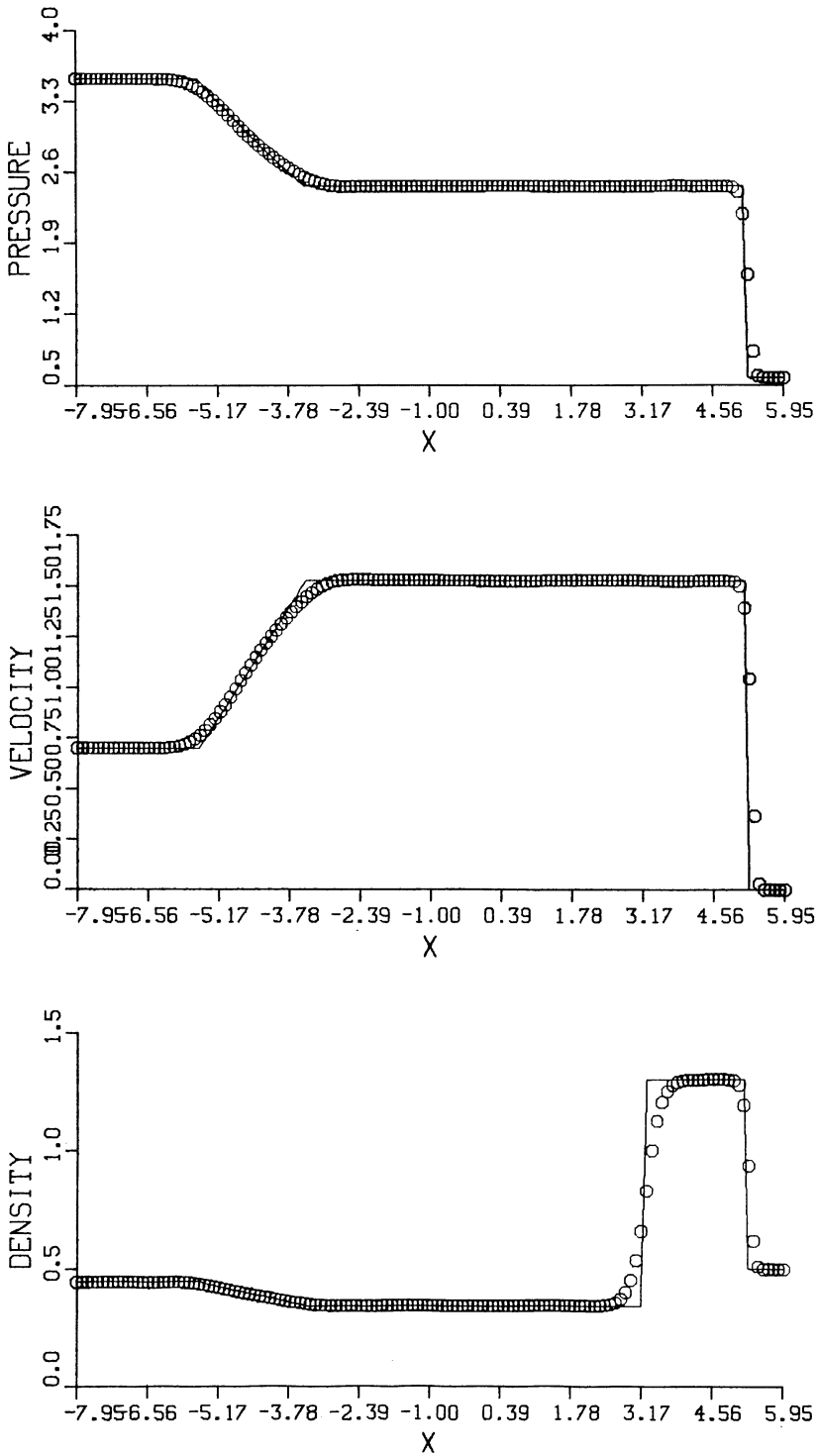


FIGURE 1(b)
K = 2 with CFL 1.8 after 55 time steps

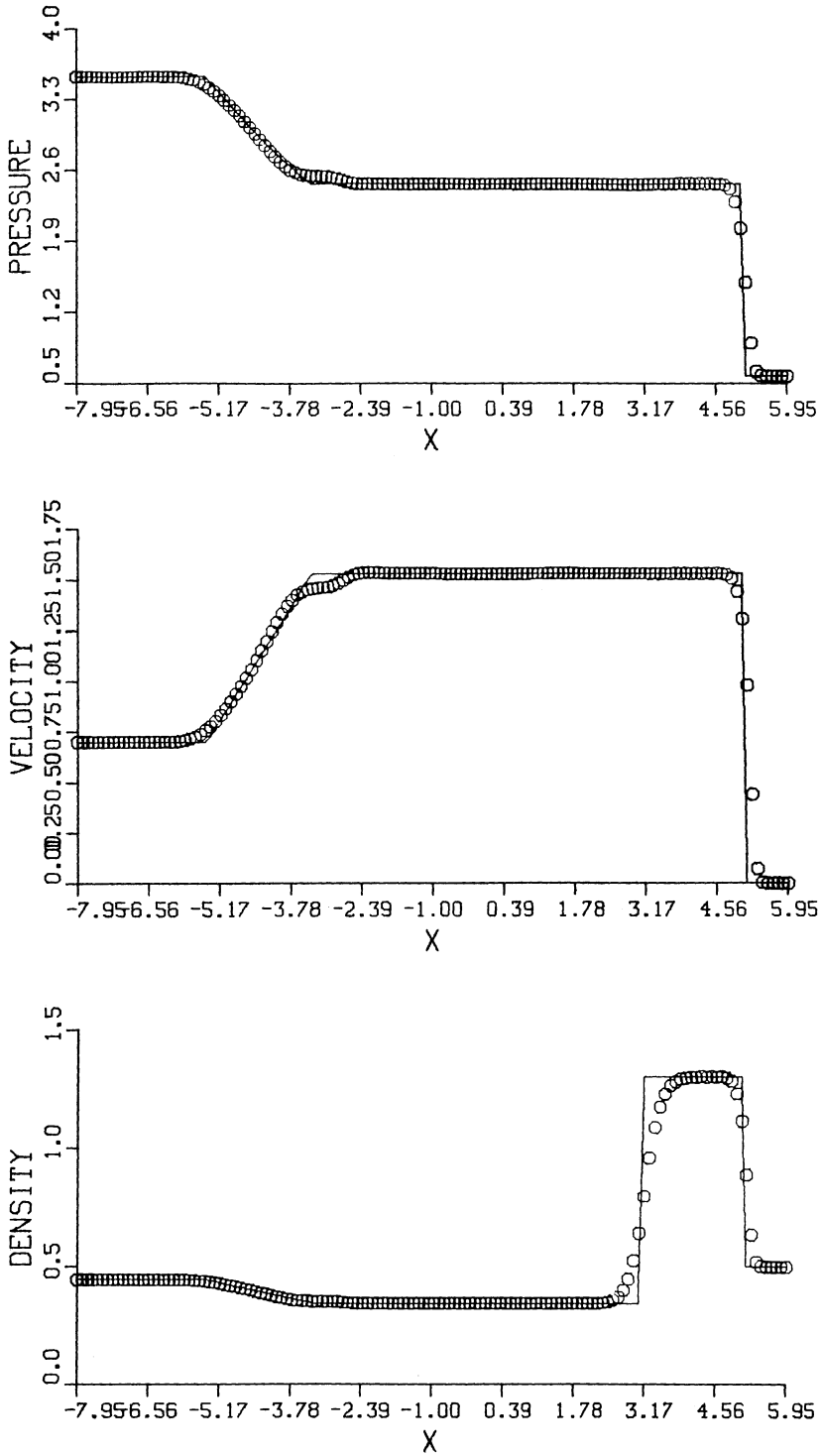


FIGURE 1(c)

$K = 4$ with CFL 3.6 after 27 time steps

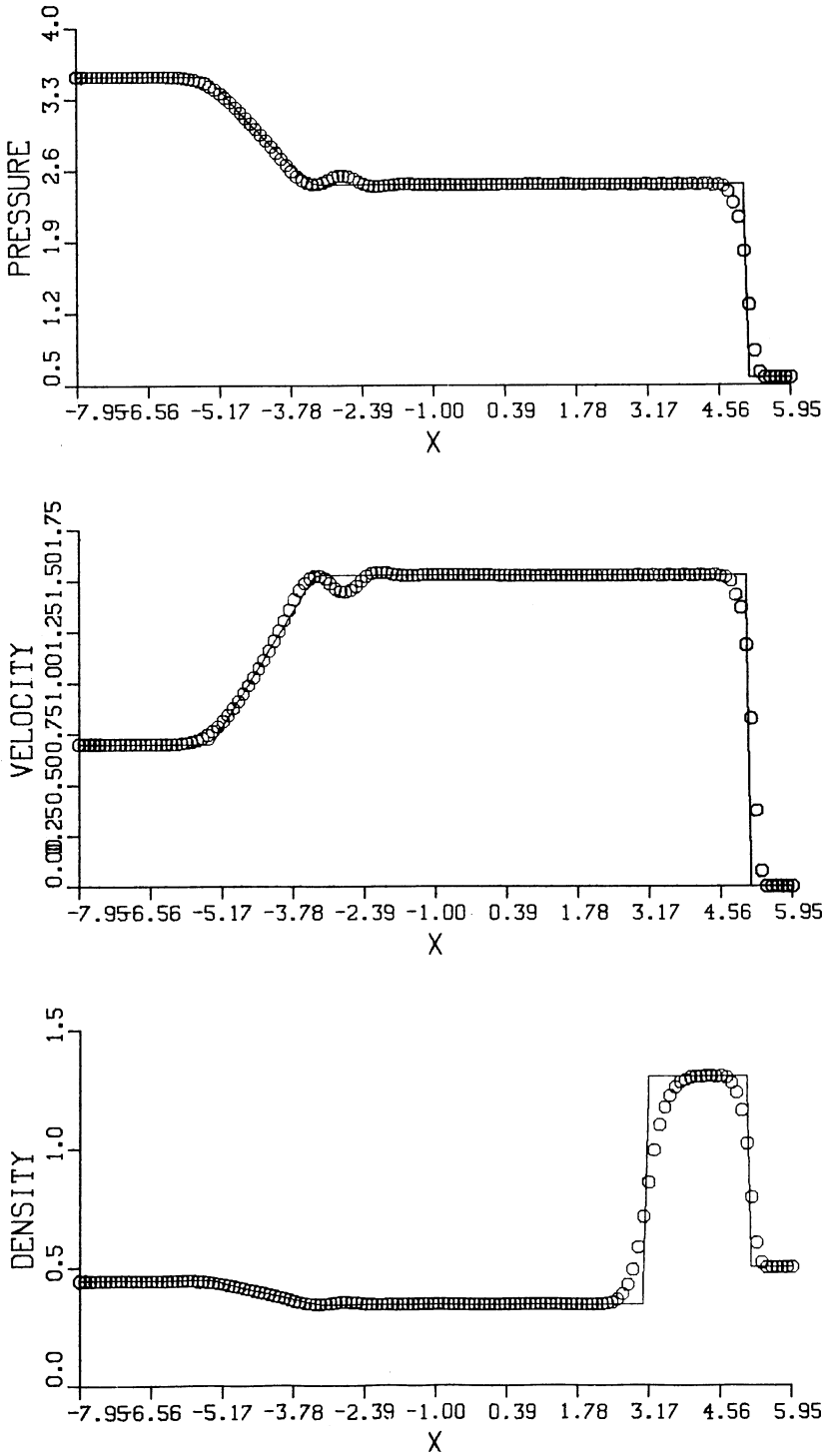


FIGURE 1(d)
K = 6 with CFL 5.4 after 18 time steps

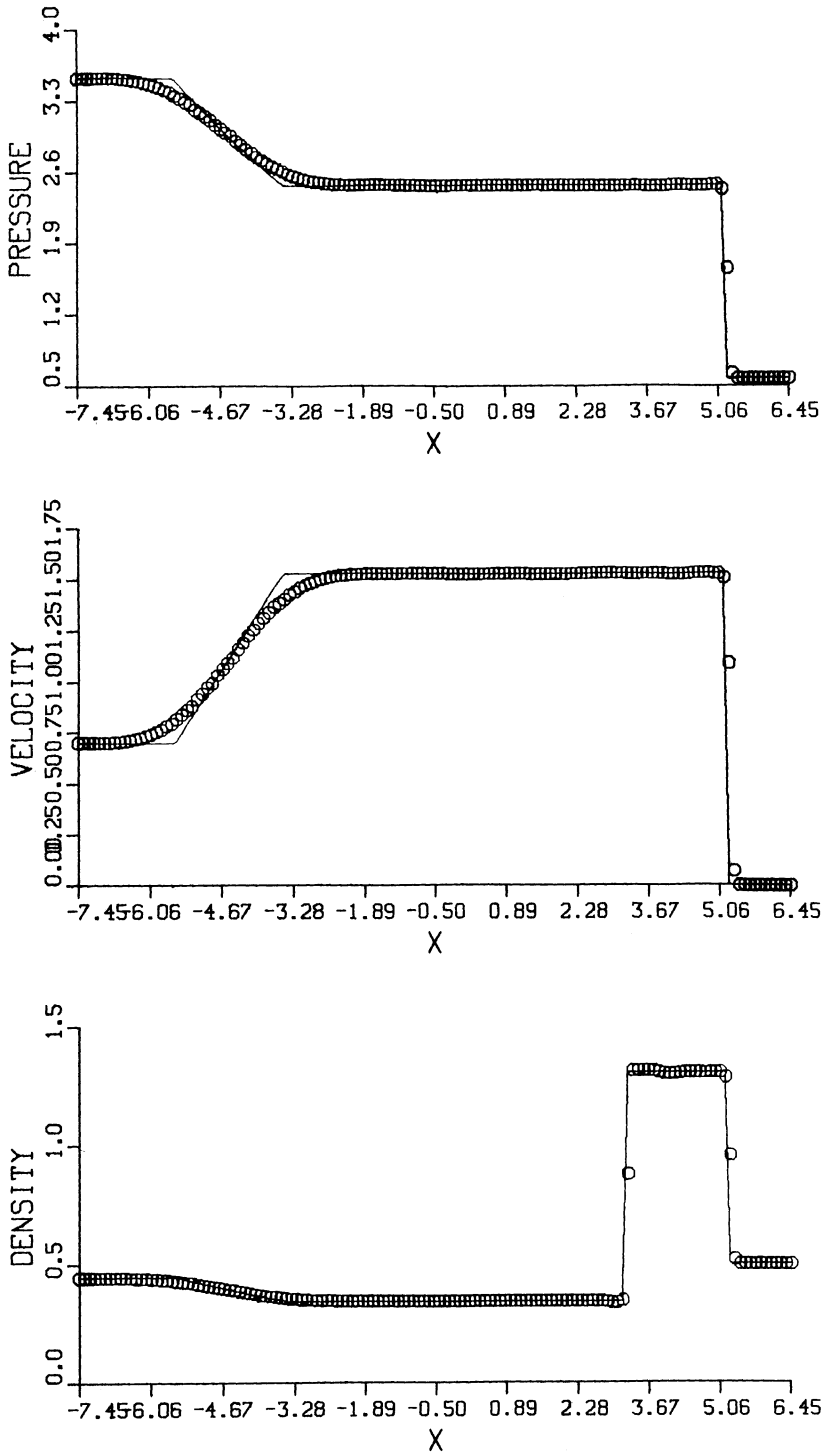


FIGURE 2(a)

$K = 1$, $CFL = 0.9$, 111 time steps with artificial compression/rarefaction

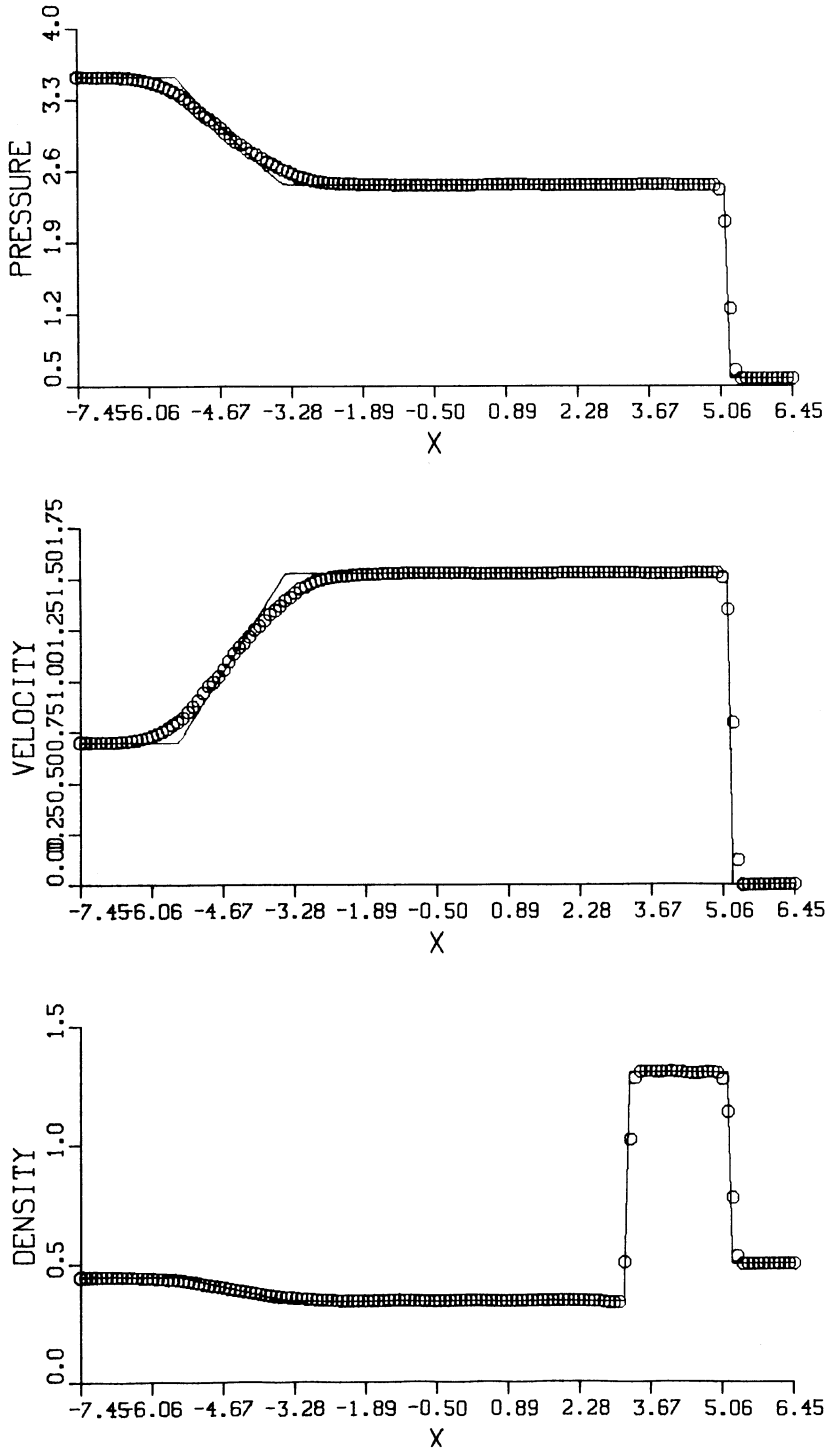


FIGURE 2(b)

$K = 2$, $CFL = 1.8$, 55 time steps with artificial compression/rarefaction

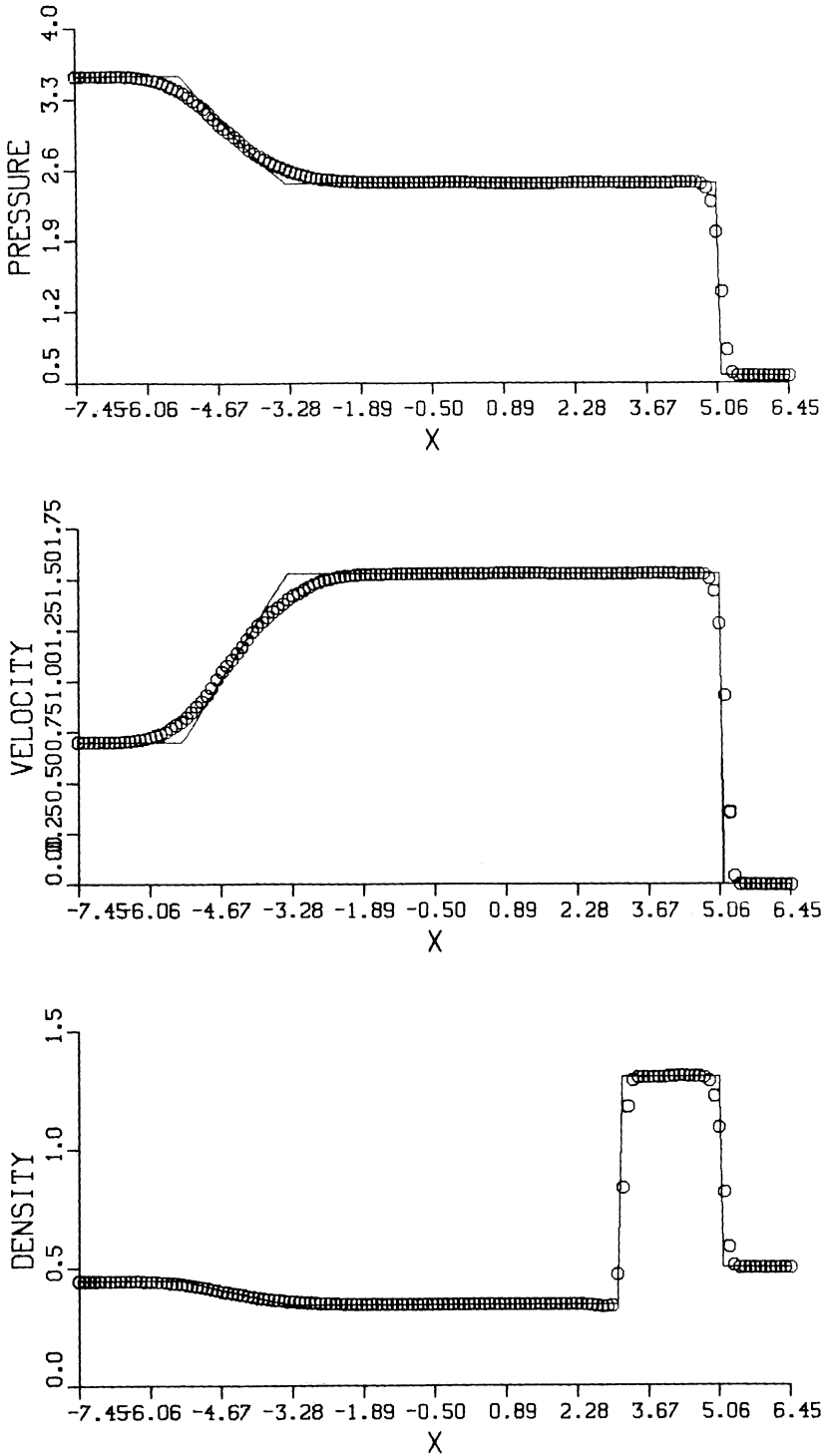


FIGURE 2(c)

$K = 4$, $CFL = 3.6$, 27 time steps with artificial compression/rarefaction

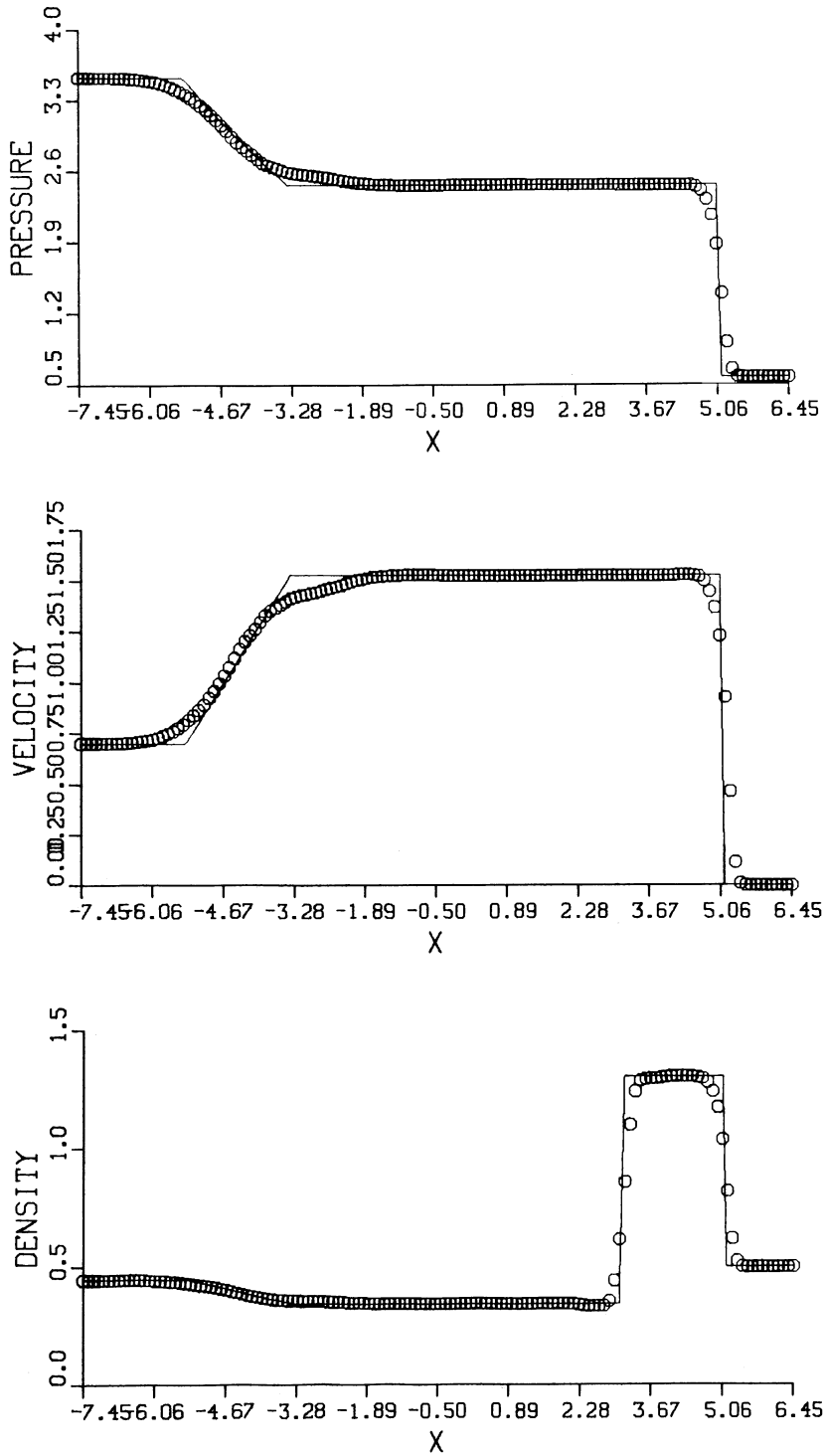


FIGURE 2(d)

$K = 6$, $CFL = 5.4$, 18 time steps with artificial compression/rarefaction

APPENDIX A

Treatment of Boundaries. In this section, we discuss numerical approximation for the initial boundary-value problem

$$(A.1) \quad u_t + f(u)_x = 0, \quad 0 \leq x < \infty, t \geq 0.$$

Let w^k be the characteristic variable (4.2). The initial boundary-value problem is well-posed if

$$(A.2) \quad w^k(0, t) = \begin{cases} \text{specified} & \text{if } a^k(u(0, t)) > 0, \\ \text{unspecified} & \text{if } a^k(u(0, t)) \leq 0. \end{cases}$$

We assume that we know v_j^n for $j \geq 1$. We set v_0^n to be the following: Let $w_1^k = l^k(v_1^n)v_1^n$ be the characteristic variable of v_1^n , and define

$$(A.3a) \quad v_0^n = v_1^n - \sum_{k=1}^m \alpha_1^k r^k(v_1^n),$$

where

$$(A.3b) \quad \alpha_1^k = \begin{cases} w_1^k - w^k(0, t_n) & \text{if } a^k(v_1) \geq 0, \\ 0 & \text{if } a^k(v_1) < 0. \end{cases}$$

Next we define in (4.5)

$$(A.3c) \quad g_0^k = \sigma(v_{1/2}^k) \alpha_{1/2}^k, \quad 1 \leq k \leq m.$$

Note that the definition (A.3c) does not alter the properties (3.2).

With this setup, we can use the scheme (4.5) to compute v_j^{n+1} , $j \geq K$. Clearly, for $K = 1$ the boundary treatment is complete, and we can compute the numerical solution for all $j \geq 0$, $n \geq 0$. Hence, for $K > 1$ we can complete the calculation of v_j^{n+1} for $1 \leq k \leq K - 1$ by successively applying (4.5) with $K = 1$ and (A.3) to obtain

$$v_{J_l}^{n+(l+1)/K}, \quad 0 \leq l \leq K - 1, \quad 1 \leq J_l \leq 3(K - 1) - 2l.$$

To maintain conservation form of the combined scheme, we use a technique suggested by Osher [7]. Let $\tilde{f}_{j+1/2,K}^n$ denote the numerical flux (4.5b) with an integer K computed at v^n , and evaluate the combined scheme as follows: Use (A.3) and

$$(A.4a) \quad v_{J_l}^{n+(l+1)/K} = v_{J_l}^{n+l/K} - \frac{\lambda}{K} (\tilde{f}_{J_l+1/2,1}^{n+l/K} - \tilde{f}_{J_l-1/2,1}^{n+l/K})$$

for $0 \leq l \leq K - 1, 1 \leq J_l \leq 3(K - 1) - 2l,$

to evaluate v_j^{n+1} for $1 \leq j \leq K - 1$. To calculate v_j^{n+1} for $K + 1 \leq j$ we use

$$(A.4b) \quad v_j^{n+1} = v_j^n - \lambda (\tilde{f}_{j+1/2,K}^n - \tilde{f}_{j-1/2,K}^n).$$

The advance to the level $n + 1$ is completed by

$$(A.4c) \quad v_K^{n+1} = v_K^n = \lambda (\hat{f}_{K+1/2,k}^n - \hat{f}_{K-1/2}^n)$$

where

$$(A.4d) \quad f_{K-1/2} = \frac{1}{K} \sum_{l=0}^{K-1} \tilde{f}_{L-1/2,1}^{n+l/K}.$$

It is the particular definition of v_K^{n+1} in (A.4c, d) that enforces the conservation form of the combined scheme.

The combined scheme (A.4) has an $O(h^3)$ truncation error for all $j \geq 1$. Using the first-order accurate 3-point scheme (2.9) in (A.4a) reduces its domain of dependence to $1 \leq J_i \leq 2(K - 1 - i)$. This is certainly adequate if $a^k(v_1^n) < 0$ for all $1 \leq k \leq m$.

We remark that if we take $\varepsilon = 0$ in (2.12a), i.e., $Q(v) \equiv |v|$, then (4.5) becomes an upstream-differencing scheme with respect to the characteristic field $a + \gamma/\lambda$. Relations (3.2c) and (3.4c) imply that

$$(A.5a) \quad |\gamma| < |v| \quad \text{for } |v| \leq K$$

and, therefore,

$$(A.5b) \quad \text{sgn}(v + \gamma) = \text{sgn}(v).$$

This shows that (4.5) is an upstream-differencing scheme with respect to the original characteristic field $a(u)$.

Next, let us consider the scalar case and assume that $a_{j+k+1/2}^n < 0$ for $-K + 1 \leq k \leq 0$. It follows from (2.9), (2.7c, d) and (2.8a, b) that the numerical flux (3.1b) can be expressed as

$$(A.5c) \quad \lambda \bar{f}_{j+1/2} = \lambda f_{j+1} + g_{j+1} - \sum_{k=1}^{K-1} C_k(v_{j+k+1/2} + \gamma_{j+k+1/2}) \Delta_{j+k+1/2} v^n.$$

Thus, if $a_{j+1/2} < 0$ for $0 \leq j \leq K - 1$, then $\lambda \bar{f}_{j+1/2}$ in (A.5c) is defined for $j \geq 1$, and therefore v_j^{n+1} can be calculated by (3.1) for all $j \geq 0$. This also shows that if $a_{j+1/2} < 0$ near $j = 0$, then the combined scheme (A.4) does not use the extrapolated value v_j^0 .

TABLE

The coefficients $c_k(x)$ (2.7d) for $2 \leq K \leq 5$

K	c_1	c_2	c_3	c_4
2	x^2			
3	$x^2(3 - x)$	x^3		
4	$x^2(6 - 4x + x^2)$	$2x^3(2 - x)$	x^4	
5	$x^2(10 - 10x + 5x^2 - x^3)$	$x^3(10 - 10x + 3x^2)$	$x^4(5 - 3x)$	x^5

Department of Mathematics
Tel-Aviv University
Tel-Aviv, Israel

1. Y. BRENNIER, private communication.
2. A. HARTEN, "High resolution schemes for hyperbolic conservation laws," *J. Comput. Phys.*, v. 49, 1983, pp. 357-393.
3. A. HARTEN, "On a class of high resolution total-variation stable finite-difference schemes," *SIAM J. Numer. Anal.*, v. 21, 1984, pp. 1-23.
4. A. HARTEN, P. D. LAX & B. VAN LEER, "On upstream differencing and Godunov-type schemes for hyperbolic conservation laws," *SIAM Rev.*, v. 25, 1983, pp. 35-61.
5. R. J. LEVEQUE, "Large time-step shock capturing techniques for scalar conservation laws," *SIAM J. Numer. Anal.*, v. 19, 1982, pp. 1091-1109.

6. R. J. LEVEQUE, *Towards a Large Time-Step Algorithm for Systems of Conservation Laws: Preliminary Results Ignoring Interactions*, Numerical Analysis Report 7/82, University of Reading, 1982.
7. S. OSHER & R. SANDERS, "Numerical approximations to nonlinear conservation laws with locally varying time and space grids," *Math. Comp.*, v. 41, 1983, pp. 321–336.
8. P. L. ROE, Proc. Seventh Internat. Conf. on Numerical Methods in Fluid Dynamics, Stanford/NASA Ames, June 1980, Springer-Verlag.
9. P. L. ROE, "Approximate Riemann solvers, parameter vectors, and difference schemes," *J. Comput. Phys.*, v. 43, 1981, pp. 357–372.
10. H. C. YEE, R. F. WARMING & A. HARTEN, *Application of TVD Schemes for the Euler Equations of Gas Dynamics*, Proceedings of the AMS–SIAM Summer Seminar on Large Scale Computations in Fluid Mechanics, Lectures in Appl. Math., Vol. 22, 1985.