# Summation by Parts for Finite Difference Approximations for $d/dx$

Bo Strand

*Department of Scientific Computing, Uppsala University, Uppsala, Sweden*

Finite difference approximations for $d/dx$ which satisfy a summation by parts rule, have been evaluated, using different types of norms, $H$, by implementation in the symbolic language Maple. In the simpler *diagonal norm*, $H = diag(\lambda_0, \lambda_1, ..., \lambda_{2r-1}, 1, ...)$ with all $\lambda_i$ positive, difference operators accurate of order $r \leqslant 4$ at the boundary and accurate of order $2r$ in the interior have been evaluated. However, it was found that the difference operators form multi-parameter families of difference operators when $r \geqslant 3$. In the general *full norm*, $H = diag(\bar{H}, I)$, with $\bar{H} \in \mathfrak{R}^{r+1 \times r+1}$ being SPD, and $I$ the identity matrix, difference operators accurate or order $r = 3, 5$ at the boundary and accurate of order $r + 1$ in the interior have been computed. As in the diagonal norm case we obtain a multiparameter family of operators when $r \geqslant 3$. Finally, a three-parameter family of difference approximations with accuracy three at and near the boundary and with accuracy four in the interior have been cumputed using *restricted full norms*. Here, $H = diag(\bar{H}, I)$, with $\bar{H} \in \mathfrak{R}^{r+2 \times r+2}$ being SPD, and $\bar{H}(:, 1) = \bar{H}(1, :)^\mathsf{T} = \kappa e_1$, where $\kappa$ is a constant and $e_1$ is the vector with the first element being one and the rest zero. Regardless of which norm we use, the parameters can be determined such that the bandwidth of the difference operators are minimized. This is of interest when parallel computers are used, since the bandwidth determines the memory requirement and also the amount of computational work. © 1994 Academic Press, Inc.

## 1. INTRODUCTION

To make the presentation self-contained, we begin by presenting the theory in [1, 4]. This makes it easier for the reader to understand the different properties of the difference approximations presented in Section 3 and beyond. Consider a real system of partial differential equations [1]

$$\frac{\partial u}{\partial t} = P\left(x, t, \frac{\partial}{\partial x}\right) u \tag{1}$$

in a cylindrical domain $\{\Omega \times 0 \leqslant t < \infty\}$ with boundary $\{\partial \Omega \times 0 \leqslant t < \infty\}$. Here, $x = (x^{(1)}, ..., x^{(s)})^\mathsf{T}$ denotes a point in the real Euclidean space $\mathfrak{R}^s$, $u = (u^{(1)}, ..., u^{(n)})^\mathsf{T}$ is a vector function and $P$ is a differential operator with matrix coefficients. For $t = 0$, initial values are given,

$$u(x, 0) = f(x), \qquad x \in \Omega, \tag{2}$$

and on the boundary, homogeneous boundary conditions,

$$B\left(x, \frac{\partial}{\partial x}\right) u = 0, \qquad x \in \partial \Omega, \qquad t \geqslant 0, \tag{3}$$

are prescribed. Here, $B$ denotes a differential operator whose coefficients depend on $x$ but not on $t$. Let

$$(u, v) = \int_\Omega u^* v \, dx, \qquad (u, u) = \|u\|^2 \tag{4}$$

denote the usual $L_2$-scalar product and norm. Assume that the operator $P$ is semibounded, i.e., there is a dense set $S \subset L_2(\Omega)$ of functions $w$ satisfying the boundary conditions (3) such that

$$(w, Pw) + (Pw, w) \leqslant 0, \qquad w \in S. \tag{5}$$

As is well known, (5) implies an energy estimate for those solutions of (1)–(3) which, for every fixed $t$, belong to $S$ because

$$\frac{\partial}{\partial t} \|u\|^2 = \left(u, \frac{\partial u}{\partial t}\right) + \left(\frac{\partial u}{\partial t}, u\right)$$

$$= (u, Pu) + (Pu, u) \leqslant 0. \tag{6}$$

## 2. SEMIBOUNDED DIFFERENCE APPROXIMATIONS FOR $d/dx$

Consider the half-line $\{0 \leqslant x < \infty\}$ and divide it into intervals of length $h > 0$ [4]. Let $x_\nu = \nu h$, $\nu = 0, 1, ...$, denote the grid points and $v_\nu = v(x_\nu)$ be real scalar grid functions with $\sum_{\nu=0}^\infty |v_\nu|^2 h < \infty$. Define a discrete scalar product and norm by

$$(u, v)_h = \langle u', Hv' \rangle_h + \sum_{\nu=r}^\infty u_\nu v_\nu h$$

$$= \sum_{i,j=0}^{r-1} h_{ij} u_i v_j h + \sum_{\nu=r}^\infty u_\nu v_\nu h, \tag{7}$$

$$\|u\|_h^2 = (u, u)_h.$$

581/110/1-4

Here $u^I = (u_0, u_1, ..., u_{r-1})^T$ denotes the vector formed with the first $r$ values of $u$ and $H = H^T > 0$ is a positive definite symmetric $(r \times r)$-matrix. We want to construct difference approximations $Q$ to $d/dx$ such that

$$(u, Qv)_h = -(Qu, v)_h - u_0 v_0, \qquad \forall u, v. \tag{8}$$

This is a discretisation of the integration by parts formula (6). An equivalent formulation is given in the following lemma.

LEMMA 2.1. *The relation (8) is equivalent with*

$$(u, Qu) = -\tfrac{1}{2} u_0^2, \qquad \forall u. \tag{9}$$

*Proof.* $(u, v)_h = (v, u)_h$ shows that (8) implies (9). If (9) holds then

$$(u + v, Q(u + v))_h = -(Q(u + v), u + v)_h - (u_0 + v_0)^2$$

and (8) follows easily using (9) for $u$ and $v$. ∎

Consider $u$ as an infinite column vector $u = (u_0, u_1, ...)^T$ and therefore represent $Q$ as an infinite matrix. Assume that $Q$ has the form

$$hQ = \begin{pmatrix} Q_{11} & Q_{12} \\ -C^T & D \end{pmatrix}. \tag{10}$$

Here

$$Q_{11} = \begin{pmatrix} q_{00} & q_{01} & \cdots & q_{0r-1} \\ \vdots & & & \vdots \\ q_{r-10} & q_{r-11} & \cdots & q_{r-1r-1} \end{pmatrix},$$

$$Q_{12} = \begin{pmatrix} q_{0r} & \cdots & q_{0m} & 0 & \cdots \\ \vdots & & \vdots & \vdots & \\ q_{r-1r} & \cdots & q_{r-1m} & 0 & \cdots \end{pmatrix},$$

$$C = \begin{pmatrix} C_0 & 0 & \cdots \\ C_s & 0 & \cdots \end{pmatrix},$$

where

$$C_s = \begin{pmatrix} \alpha_s & 0 & \cdots & \cdots & 0 \\ \alpha_{s-1} & \alpha_s & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ \alpha_1 & \cdots & \cdots & \alpha_{s-1} & \alpha_s \end{pmatrix},$$

$$D = \begin{pmatrix} 0 & \alpha_1 & \cdots & \alpha_s & 0 & \cdots & & & 0 & \cdots \\ -\alpha_1 & 0 & \alpha_1 & \cdots & \alpha_s & 0 & \cdots & & & \\ \vdots & \ddots & \ddots & \ddots & & \ddots & \ddots & & & \\ -\alpha_s & \cdots & -\alpha_1 & 0 & \alpha_1 & \cdots & \alpha_s & \ddots & & \vdots \\ 0 & \ddots & & \ddots & \ddots & \ddots & & \ddots & \ddots & \vdots \\ \vdots & & -\alpha_s & \cdots & -\alpha_1 & 0 & \alpha_1 & \cdots & \alpha_s & 0 & \cdots \end{pmatrix}, \tag{11}$$

where $D$ is an antisymmetric band matrix and $C_0$ is a $(r - s) \times s$ matrix only with zeros. The part, $h^{-1}(-C^T D)$, represents the operator away from the boundary.

## 2.1. *The Accuracy Conditions*

An operator $Q$ is an approximation of $d/dx$ accurate of order $\tau$ if

$$Qv = \frac{dv}{dx} + \mathcal{O}(h^\tau), \tag{12}$$

where $v$ is a real scalar grid function. It is sufficient to consider polynomials, because using $n + 1$ points, a polynomial of degree $\leqslant n$ can be interpolated. Then (12) gives us

$$Qx^m - \frac{dx^m}{dx} = 0, \qquad m = 0, 1, ..., \tau. \tag{13}$$

Let $h = 1$, and denote by

$$w_j = \begin{pmatrix} e_j \\ f_j \end{pmatrix}, \qquad e_j = (-1)^j \begin{pmatrix} r^j \\ (r-1)^j \\ \vdots \\ 1^j \end{pmatrix}, \tag{14}$$

$$f_j = \begin{pmatrix} 0^j \\ 1^j \\ \vdots \end{pmatrix}, \qquad j = 0, 1, 2, ...$$

the discretisation of $(x - r)^j$, with the conventions $0^0 = 1$, $e_{-1} = 0$.

The following aproximations of increasing order of accuracy $2s$ of the derivative are used:

$$\frac{\partial}{\partial x} \sim D^{[2s]}(h) \equiv \sum_{v=1}^{s} \lambda_v D_0(vh),$$

$$\lambda_v = \frac{-2(-1)^v (s!)^2}{(s+v)! (s-v)!}, \qquad v = 1, 2, .... \tag{15}$$

These formulae use the coordinates of $2s + 1$ symmetric centered points with antisymmetric coefficients and so they obtain the highest accuracy possible with this number of equidistant points. The coefficient of the $v$th right-hand term is

$$\alpha_v = \frac{-2(-1)^v s!}{v(s+v)! (s-v)!}, \tag{16}$$

and the corresponding left-hand term is $\alpha_{-v} = -\alpha_v$. The following lemma characterizes the accuracy of $Q$ in the following lemma.

LEMMA 2.2. *The operator* $h^{-1}(-C^T D)$ *approximates* $d/dx$ *of order* $2s > 0$ *if and only if*

$$\sum_{v=1}^{s} \alpha_v v^{2n+1} = \begin{cases} \frac{1}{2}, & n = 0, \\ 0, & n = 1, 2, ..., s-1. \end{cases} \quad (17)$$

*Proof.* Equation (12) gives us

$$h^{-1}(-C^T D) u$$

$$= \frac{du^{II}}{dx} + \mathcal{O}(h^{2s}), \qquad u^{II} = (u_r, u_{r+1}, ...),$$

$$\Rightarrow h^{-1} \sum_{k=-s}^{s} \alpha_k u(x_{r-1+l+k})$$

$$= \frac{du(x_{r-1+l})}{dx} + \mathcal{O}(h^{2s}), \qquad l = 1, ..., s.$$

Then by using Taylors series for functions of one variable,

$$f(x) = \sum_{j=0}^{\infty} f^{(j)}(a) \frac{(x-a)^j}{j!},$$

one obtains

$$u(x_{r-1+l+k}) = \sum_{j=0}^{\infty} u^{(j)}(x_{r-1+l}) \frac{x_k^j}{j!}.$$

This expression and the fact that $x_v = vh$ give the following relation:

$$h^{-1} \sum_{j=0}^{\infty} u^{(j)}(x_{r-1+l}) \frac{h^j}{j!} \sum_{k=-s}^{s} \alpha_k k^j$$

$$= u'(x_{r-1+l}) + \mathcal{O}(h^{2s}), \qquad l = 1, 2, ..., s.$$

Thus

$$\sum_{k=-s}^{s} \alpha_k = 0, \qquad j = 0,$$

$$\sum_{k=-s}^{s} \alpha_k k = 1, \qquad j = 1,$$

$$\sum_{k=-s}^{s} \alpha_k k^j = 0, \qquad j = 2, ..., s.$$

Since $\alpha_{-v} = -\alpha_v$, these equations can be rewritten as

$$\sum_{k=1}^{s} \alpha_k = 0, \qquad j = 0,$$

$$\sum_{k=1}^{s} \alpha_k k = \frac{1}{2}, \qquad j = 1, \quad (18)$$

$$\sum_{k=1}^{s} \alpha_k \{k^j - (-k)^j\} = 0, \qquad j = 2, ..., s.$$

The third equation is clearly zero for $j$ even and the condition

$$\sum_{k=1}^{s} \alpha_k k^j = 0, \qquad j = 3, 5, ..., 2s-1,$$

must be valid and the proof is completed. ∎

The part, $h^{-1}(Q_{11} Q_{12})$, represents the modification of $Q$ at the boundary points. The accuracy conditions which $Q$ has to fulfill are given in Lemma 2.3.

LEMMA 2.3. *The operator* $h^{-1}(Q_{11} Q_{12})$ *approximates* $d/dx$ *with order of accuracy* $\tau$ *at the points* $x_v$, $v = 0, 1, ..., r-1$, *if and only if*

$$j e_{j-1} = Q_{11} e_j + Q_{12} f_j, \qquad j = 0, 1, ..., \tau. \quad (19)$$

*Proof.* Equations (13) and (14) for the polynomial $(x-r)^j$, with $j = 0, 1, ..., \tau$, give the expression

$$Q w_j = j w_{j-1}, \qquad j = 0, 1, ..., \tau.$$

Thus by (10)

$$Q_{11} e_j + Q_{12} f_j = j e_{j-1}, \qquad j = 0, 1, ..., \tau,$$

$$-C^T e_j + D f_j = j f_{j-1}, \qquad j = 0, 1, ..., \tau,$$

and this concludes the proof. ∎

### 2.2. Necessary and Sufficient Conditions for Q

Necessary and sufficient conditions will be derived such that $Q$ satisfies condition (9).

THEOREM 2.1. *The operator* $Q$ *satisfies the relation* (9) *if and only if it can be written as*

$$hQ = \begin{pmatrix} H^{-1}B & H^{-1}C \\ -C^T & D \end{pmatrix}, \quad (20)$$

*where $B$ is a $(r \times r)$-matrix of the form*

$$B = B_1 + B_2 \quad (21)$$

*with*

$$B_1 = \begin{pmatrix} -\frac{1}{2} & 0 & \cdots & \cdots & 0 \\ 0 & 0 & \cdots & \cdots & 0 \\ \vdots & \vdots & & & \vdots \\ 0 & \cdots & \cdots & \cdots & 0 \end{pmatrix},$$

$$B_2 = \begin{pmatrix} 0 & b_{01} & \cdots & & b_{0r-1} \\ -b_{01} & 0 & b_{12} & \cdots & b_{1r-1} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & b_{r-2r-1} \\ -b_{0r-1} & \cdots & \cdots & -b_{r-2r-1} & 0 \end{pmatrix};$$

*hence $B_2 = -B_2^T$.*

*Proof.* Write (9) in the form

$$-\tfrac{1}{2}u_0^2 = \langle u^{\mathrm{I}}, HQ_{11}u^{\mathrm{I}}\rangle + \langle u^{\mathrm{I}}, HQ_{12}u^{\mathrm{II}}\rangle$$
$$- \langle u^{\mathrm{II}}, C^{\mathrm{T}}u^{\mathrm{I}}\rangle + \langle u^{\mathrm{II}}, Du^{\mathrm{II}}\rangle$$

where $u^{\mathrm{I}} = \langle u_0, ..., u_{r-1}\rangle^{\mathrm{T}}$, $u^{\mathrm{II}} = (u_r, ...)^{\mathrm{T}}$. $D$ is an antisymmetric matrix and therefore $\langle u^{\mathrm{II}}, Du^{\mathrm{II}}\rangle \equiv 0$. Furthermore, if $u^{\mathrm{II}} \equiv 0$ then the above relation becomes

$$-\tfrac{1}{2}u_0^2 = \langle u^{\mathrm{I}}, HQ_{11}u^{\mathrm{I}}\rangle.$$

Therefore $HQ_{11} = B$, must have the form (21). Then the above relation is equivalent with

$$0 = \langle u^{\mathrm{I}}, HQ_{12}u^{\mathrm{II}}\rangle - \langle u^{\mathrm{II}}, C^{\mathrm{T}}u^{\mathrm{I}}\rangle$$
$$= \langle u^{\mathrm{I}}, (HQ_{12} - C)u^{\mathrm{II}}\rangle$$

for all vectors $u^{\mathrm{I}}, u^{\mathrm{II}}$. This is only possible if $Q_{12} = H^{-1}C$ and the theorem is proved. ∎

To obtain the coefficients of $Q$ at the boundary points, the elements of $B_2$ have to be evaluated. This is done by writing (19) as

$$\begin{aligned} B_2 e_j &= g_j, \\ g_j &= jHe_{j-1} - B_1 e_j - Cf_j, \qquad j = 0, 1, ..., \tau. \end{aligned} \tag{22}$$

The equation above, for the antisymmetric matrix $B_2$, can be written as a system of linear equations

$$A_j b = y_j, \qquad j = 0, ..., \tau,$$

$$A_j = \begin{pmatrix}
(1-r)^j & (2-r)^j & \cdots & \cdots & (-1)^j & 0 & 0 \cdots \\
-(-r)^j & 0 & \cdots & \cdots & 0 & (2-r)^j & \cdots \\
0 & -(-r)^j & 0 & \cdots & 0 & -(1-r)^j & 0 \\
\vdots & & \ddots & & & & \ddots \\
\vdots & & & \ddots & & & \\
0 & \cdots & \cdots & 0 & -(-r)^j & 0 & \cdots
\end{pmatrix}$$

$$\begin{pmatrix}
& & 0 & 0 \\
\cdots & (-1)^j & 0 & \cdots & & \\
\cdots & & & & 0 & 0 \\
& & 0 & (-2)^j & (-1)^j & 0 \\
\ddots & & & -(-3)^j & 0 & (-1)^j \\
0 & -(1-r)^j & 0 & \cdots & 0 & -(-3)^j & -(-2)^j
\end{pmatrix} \tag{23}$$

where $b = (b_{01}, b_{02}, ..., b_{0r-1}, b_{12}, ..., b_{r-2r-1})$ contains the $r(r-1)/2$ unknown elements of $B_2$; $A_j$ is the $(r \times r(r-1)/2)$-coefficient matrix, and

$$y_j = g_j - \tfrac{1}{2}(-r)^j, \qquad j = 0, ..., \tau. \tag{24}$$

The term $\tfrac{1}{2}(-r)^j$ stems from $B_1(1, 1)$, which is known $(-\tfrac{1}{2})$. These $\tau + 1$ equations can be written as

$$Gb = z, \tag{25}$$

where $G = [A_0 A_1 \cdots A_\tau]^{\mathrm{T}}$ and $z = [y_0 y_1 \cdots y_\tau]^{\mathrm{T}}$. This is an overdetermined system of $(\tau + 1)r$ equations with $r(r-1)/2$ unknown elements.

By assumption $B_2$ is antisymmetric. Therefore the following compatibility conditions for the system (22) have to be satisfied:

$$\langle e_i, g_j\rangle + \langle e_j, g_i\rangle = \langle e_i, B_2 e_j\rangle + \langle e_j, B_2 e_i\rangle$$
$$= 0, \qquad 0 \leqslant i, \qquad j \leqslant \tau. \tag{26}$$

If these conditions hold, the system (22) can be resolved as it is expressed in Lemma 2.4.

LEMMA 2.4. *Assume that $r \geqslant \tau + 1$ and that the relations (26) hold. Then there is an antisymmetric matrix $B_2$ such that (22) is valid.*

### 2.3. *Derivation of the Norm*

To obtain the elements of the norm matrix $H$, the relations (26) are used. By (22) this can also be written as

$$j\langle e_i, He_{j-1}\rangle + i\langle e_j, He_{i-1}\rangle$$
$$= M_{i,j}, \qquad 0 \leqslant i, \qquad j \leqslant \tau, \tag{27}$$

where

$$M_{i,j} = 2\langle e_i, B_1 e_j\rangle + \langle e_i, Cf_j\rangle + \langle e_j, Cf_i\rangle.$$

LEMMA 2.5. $M_{i,j}$ *can be written as*

$$M_{i,j} = -(-r)^{i+j} + J_{i,i+j},$$
$$J_{i,\sigma} = \sum_{\nu=1}^{s} \alpha_\nu \left( \sum_{\mu=0}^{\nu-1} \mu^{\sigma-i}(\mu - \nu)^i + \mu^i(\mu - \nu)^{\sigma-i} \right), \tag{28}$$
$$\sigma = i + j \geqslant 1.$$

*Proof.* These scalar products can be calculated as

$$\langle e_i, B_1 e_j\rangle = \tfrac{1}{2}(-1)^{i+j+1} r^{i+j}$$
$$\langle e_i, Cf_j\rangle = (-i)^i \sum_{\nu=1}^{s} \alpha_\nu \sum_{\mu=0}^{\nu-1} \mu^j(\nu - \mu)^i.$$

Thus, by using the above expression in (27) and by introducing $\sigma = i + j$, the proof is concluded. ∎

If one introduces the notation $\rho_{i,j} = \langle e_i, He_j\rangle$ then (27) can be written

$$j\rho_{i,j-1} + i\rho_{j,i-1} = M_{i,j}, \qquad 0 \leqslant i, \qquad j \leqslant \tau. \tag{29}$$

Here $\rho_{i,-1} = \rho_{-1,i} = 0$ by the convention for $e_{-1}$ and $\rho_{i,j} = \rho_{j,i}$ by the symmetry of $H$. This system in the $\rho_{i,j}$ is very simple to resolve, and from the solution the elements of $H$ are obtained. The matrix $H$ has to be positive definite in order to be used as norm matrix. Lemma 2.9 below states the equivalence between this condition and the positive definiteness of the matrix defined by the $\rho_{i,j}$. Conditions on $M_{i,j}$ will be derived in Section 2.4 such that the system (29) has a solution (these conditions are resumed in Lemma 2.6) and therefore for the compatibility of (22).

From (29) and (28) with $i = j = 0$ the following condition is obtained:

$$0 = 0\rho_{0,-1} + 0\rho_{0,-1} = M_{0,0} = -1 + 2 \sum_{v=1}^{s} \alpha_v v,$$

i.e.,

$$\sum_{v=1}^{s} \alpha_v v = \tfrac{1}{2}. \tag{30}$$

Therefore, by Lemma 2.2 the approximation at the points has to be at least second-order accurate. But (29) implies $M_{i,j} = j\rho_{i,j-1} + i\rho_{j,i-1} = M_{i,j}$. This condition is satisfied by (28). Therefore (29) has to be considered only for $i \leqslant j$. For $i = 0$ and from (29) the following condition is obtained:

$$\rho_{0,j-1} = \frac{1}{j} M_{0,j}, \qquad j = 1, 2, ..., \tau. \tag{31}$$

If $i > 0$ then $\rho_{i,j-1}$ can be explicitly calculated:

$$\rho_{i,j-1} = \frac{1}{j} M_{i,j} - \frac{i}{j} \rho_{i-1,j}. \tag{32}$$

If $i - 1 > 0$ and $j < \tau$ then (32) can be used to replace $\rho_{i-1,j}$ by $\rho_{i-2,j+1}$ and obtain $\rho_{i,j-1} = (1/j) M_{i,j} - (i/j(j+1)) M_{i-1,j+1} + (i(i-1)/j(j+1)) \rho_{i-2,j+1}$. Therefore by recursion (29) can then be written in the form

$$\rho_{i,j} = \rho_{j,i}, \qquad \rho_{i,-1} = \rho_{-1,i} = 0$$

$$\rho_{i,j-1} = \frac{1}{j} M_{i,j} - \frac{i}{j(j+1)} M_{i-1,j+1} + \cdots$$

$$+ (-1)^{\alpha} \frac{i(i-1)\cdots(i-\alpha+1)}{j(j+1)\cdots(j+\alpha)} M_{i-\alpha,j+\alpha}$$

$$+ (-1)^{\alpha+1} \frac{i(i-1)\cdots(i-\alpha)}{j(j+1)\cdots(j+\alpha)} \rho_{i-\alpha-1,j+\alpha},$$

$$0 < i \leqslant j \leqslant \tau, \tag{33}$$

where $\alpha = \min(i-1, \tau-j)$.

If $i - 1 < \tau - j$, i.e., $i + j \leqslant \tau$, then $\alpha = i - 1$ and the following condition is obtained:

$$\rho_{i-\alpha-1,j+\alpha} = \rho_{0,i+j-1} = \frac{1}{i+j} M_{0,i+j}. \tag{34}$$

Thus (31) and (33) imply that $\rho_{i,j-1}$ is completely determined by the $M_{i,j}$, provided $i + j \leqslant \tau$. If $i + j > \tau$ then $\alpha = \tau - j$ and

$$\rho_{i-\alpha-1,j+\alpha} = \rho_{i+j-1-\tau,\tau} = \rho_{v,\tau}, \qquad v = i + j - 1 - \tau. \tag{35}$$

There are no further relations which $\rho_{v,\tau}$ need to satisfy. If Eq. (33) is used with $i = n, j = n$ and $i = n - 1, j = n + 1$, for $n < \tau$, representations for $\rho_{n,n-1}$ and $\rho_{n-1,n}$ are obtained, but by (28), $\rho_{n,n-1} = \rho_{n-1,n}$. If $i + j = 2n > \tau$, then by (35), these two relations determine $\rho_{v,\tau}$, $v = 2n - 1 - \tau$, and no conditions for the $M_{i,j}$ result. In this case, if $\tau$ is odd, $v$ will be even, but if $\tau$ is even, $v$ will be odd. If Eq. (33) is used with $i + j = 2n + 1 > \tau$, $1 < n < \tau - 1$, and $i + j = 2n - 1 > \tau$, $2 < n < \tau$, then it can be shown that no conditions for the remaining $\rho_{v,\tau}$ result. If $i + j = 2n \leqslant \tau$, then from (33) and (34) the following conditions on $M_{i,j}$ are obtained:

$$\frac{1}{n} M_{n,n} - \frac{n}{n(n+1)} M_{n-1,n+1}$$

$$+ \frac{n(n-1)}{n(n+1)(n+2)} M_{n-2,n+2} - \cdots$$

$$+ (-1)^n \frac{n(n-1)\cdots 1}{n(n+1)\cdots 2n} M_{0,2n}$$

$$= \rho_{n,n-1} = \rho_{n-1,n} = \frac{1}{n+1} M_{n-1,n+1}$$

$$- \frac{(n-1)}{(n+1)(n+2)} M_{n-2,n+2} + \cdots$$

$$+ (-1)^{n-1} \frac{(n-1)\cdots 1}{(n+1)(n+2)\cdots 2n} M_{0,2n}.$$

This relation can also be written as

$$\frac{1}{2n} M_{n,n} - \frac{1}{n+1} M_{n-1,n+1}$$

$$+ \frac{(n-1)}{(n+1)(n+2)} M_{n-2,n+2} - \cdots$$

$$+ (-1)^n \frac{(n-1)\cdots 1}{(n+1)\cdots 2n} M_{0,2n}$$

$$= 0, \qquad n = 1, 2, ...; \qquad \text{with} \quad 2n \leqslant \tau. \tag{36}$$

Thus, the following lemma is given.

**LEMMA 2.6.** *The system* (29) *has a solution if and only if the difference approximation in the is at least second-order accurate and if the $M_{i,j}$ satisfy the relations* (36).

The sought conditions of compatibility for the system (29) will be obtained in terms of accuracy conditions which the chosen scheme has to fulfill. In Section 2.4 it is shown that condition (36) on $M_{i,j}$ implies that

$$\sum_{v=1}^{s} \alpha_v v^{2n+1} = 0, \qquad n = 1, 2, ..., 2n \leqslant \tau.$$

Therefore, from Lemmas 2.6 and 2.2, the following lemma is obtained.

**LEMMA 2.7.** *The system* (29) *has a solution if and only if*

$$\sum_{v=1}^{s} \alpha_v v^{2n+1} = \begin{cases} \frac{1}{2}, & n = 0, \\ 0, & n = 1, 2, ..., 2n \leqslant \tau, \end{cases}$$

*i.e., if and only if the approximation is accurate to order $\tau + 1$ if $\tau$ is odd.*

If the system has a solution when the $\rho_{i,j}$ are determined for $0 \leqslant i, j \leqslant \tau$ if one specifies those $\rho_{v,\tau}$, $v = 0, 1, ..., \tau$, which are not determined by the system. They can be used to define a symmetric matrix

$$R_\tau = \begin{pmatrix} \rho_{0,0} & \rho_{0,1} & \cdots & \rho_{0,\tau} \\ \rho_{0,1} & \rho_{1,1} & \cdots & \rho_{1,\tau} \\ \vdots & & & \vdots \\ \rho_{0,\tau} & \rho_{1,\tau} & \cdots & \rho_{\tau,\tau} \end{pmatrix} = R_\tau^T.$$

The parameters $\rho_{v,\tau}$ should be chosen so as to obtain a positive definite $R_\tau$. The following lemma can be used.

**LEMMA 2.8.** *If*

$$R_{\tau-1} = \begin{pmatrix} \rho_{0,0} & \rho_{0,1} & \cdots & \rho_{0,\tau-1} \\ \rho_{0,1} & \rho_{1,1} & \cdots & \rho_{1,\tau-1} \\ \vdots & & & \vdots \\ \rho_{0,\tau-1} & \rho_{1,\tau-1} & \cdots & \rho_{\tau-1,\tau-1} \end{pmatrix} = R_{\tau-1}^T$$

*is positive definite then one can choose $\rho_{\tau,\tau}$ such that also $R_\tau > 0$ independently of the values of $\rho_{v,\tau}$, $v = 0, 1, ..., \tau - 1$.*

*Proof.* The proof is obtained by developing the determinant of $R_\tau$ by the elements of the last row and balancing with the value of $\rho_{\tau,\tau}$. Note that system (29) does not depend on $\rho_{\tau,\tau}$. ∎

Now a positive definite $(r \times r)$-matrix $H$ has to be determined such that

$$\langle e_i, H e_j \rangle = \rho_{i,j}, \qquad 0 \leqslant i, \quad j \leqslant \tau. \tag{37}$$

In fact the following lemma holds.

**LEMMA 2.9.** *If $r \geqslant \tau + 1$ and the matrix $R_\tau$ is positive definite then there are $H = H^T > 0$ such that* (37) *holds. In particular if $r = \tau + 1$ then $H$ is uniquely defined by*

$$E^T H E = R_\tau, \qquad E = (e_0, ..., e_{r-1}). \tag{38}$$

In this case the norm is going to be referred to as a *full norm* and will have the form

$$H = \begin{pmatrix} h_{0,0} & \cdots & h_{0,\tau} \\ \vdots & & \vdots \\ h_{\tau,0} & \cdots & h_{\tau,\tau} \end{pmatrix}. \tag{39}$$

In stability analysis of first-order systems it is essential that the matrix has the form

$$H = \begin{pmatrix} \lambda_0 & 0 & \cdots & 0 \\ 0 & h_{1,1} & \cdots & h_{1,r-1} \\ \vdots & \vdots & & \vdots \\ 0 & h_{r-1,1} & \cdots & h_{r-1,r-1} \end{pmatrix} = \begin{pmatrix} \lambda_0 & 0 \\ 0 & H_1 \end{pmatrix}. \tag{40}$$

If we set $r = \tau + 2$ the corresponding norm will be referred to as a *restricted full norm*; see also [2]. Then the following lemma holds.

**LEMMA 2.10.** *If $R_\tau$ is positive definite then one can choose $H$ in the form* (40) *such that* (37) *holds. In general one has to take $r \geqslant \tau + 2$.*

*Proof.* Let $e_j = \tilde{e}_j + \tilde{\tilde{e}}_j$, where

$$\tilde{e}_j = (-1)^j \begin{pmatrix} 0 \\ \bar{e}_j \end{pmatrix},$$

$$\tilde{\tilde{e}}_j = (-1)^j \begin{pmatrix} r^j \\ 0 \\ \vdots \\ 0 \end{pmatrix},$$

$$\bar{e}_j = (-1)^j \begin{pmatrix} (r-1)^j \\ \vdots \\ 1^j \end{pmatrix}.$$

Equation (37) is equivalent to

$$\bar{\rho}_{i,j} = \langle \bar{e}_i, H_1 \bar{e}_j \rangle = \rho_{i,j} - \lambda_0 (-1)^{i+j} r^{i+j}. \tag{41}$$

If $R_\tau$ is positive definite then the matrix $\tilde{R}_\tau$ formed by the $\bar{\rho}_{i,j}$ is also positive definite, provided $\lambda_0$ is chosen sufficiently small. ∎

The main result will now be proved.

THEOREM 2.2. *For every $\tau + 1 = 2s$, $s = 1, 2, ...$, there is an $H$ of the form (40) and a scalar product of the form (7) and an approximation $Q$ of $d/dx$ which is accurate to order $\tau$ for $x = x_v$, $v = 0, 1, 2, ..., r - 1$, and accurate to order $\tau + 1$ for $x = x_v$, $v \geqslant r$, such that (8) holds.*

*Proof.* It is sufficient to show that the matrix $R_\tau$ can be made such that $R_\tau > 0$ if $r$ is chosen sufficiently large. The $\rho_{i,j}$ are the solutions of the system (29). It can be split into two parts, $\rho_{i,j} = \tilde{\rho}_{i,j} + \tilde{\tilde{\rho}}_{i,j}$, where $\tilde{\rho}_{i,j}$ is the solution of

$$j\tilde{\rho}_{i,j-1} + i\tilde{\rho}_{j,i-1} = (-1)^{i+j} r^{i+j},$$

$$0 \leqslant i, \quad j \leqslant \tau, \quad i + j > 0 \quad (42)$$

$$\tilde{\rho}_{i,-1} = \tilde{\rho}_{-1,j} = 0, \quad \tilde{\rho}_{i,j} = \tilde{\rho}_{j,i}$$

and

$$j\tilde{\tilde{\rho}}_{i,j-1} + i\tilde{\tilde{\rho}}_{j,i-1} = J_{i,i+j}, \quad 0 \leqslant i, \quad j \leqslant \tau, \quad i + j > 0$$

$$\tilde{\tilde{\rho}}_{i,-1} = \tilde{\tilde{\rho}}_{-1,j} = 0, \quad \tilde{\tilde{\rho}}_{i,j} = \tilde{\tilde{\rho}}_{j,i}. \quad (43)$$

By Lemma 2.11, in Section 2.4, the system (42) has a solution

$$\tilde{\rho}_{i,j} = \frac{(-1)^{i+j} r^{i+j+1}}{i+j+1}, \quad 0 \leqslant i, \quad j \leqslant \tau.$$

The matrix formed by the $\tilde{\rho}_{i,j}$ can be written as

$$\begin{pmatrix} \dfrac{r}{1} & -\dfrac{r^2}{2} & \cdots & (-1)^\tau \dfrac{r^{\tau+1}}{\tau+1} \\ -\dfrac{r^2}{2} & \dfrac{r^3}{3} & \cdots & (-1)^{\tau+1} \dfrac{r^{\tau+2}}{\tau+2} \\ \cdots & \cdots & \cdots & \cdots \\ (-1)^\tau \dfrac{r^{\tau+1}}{\tau+1} & \cdots & \cdots & \dfrac{r^{2\tau+1}}{2\tau+1} \end{pmatrix} = DGD,$$

where

$$D = r^{1/2} \begin{pmatrix} 1 & 0 & \cdots & \cdots & \cdots & 0 \\ 0 & -r & 0 & \cdots & \cdots & 0 \\ 0 & 0 & r^2 & 0 & \cdots & 0 \\ \vdots & & & \ddots & & \vdots \\ \vdots & & & & \ddots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 & (-1)^\tau r^\tau \end{pmatrix}$$

and

$$G = \begin{pmatrix} 1 & \dfrac{1}{2} & \cdots & \cdots & \dfrac{1}{\tau+1} \\ \dfrac{1}{2} & \dfrac{1}{3} & \dfrac{1}{4} & \cdots & \dfrac{1}{\tau+2} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \dfrac{1}{\tau+1} & \dfrac{1}{\tau+2} & \cdots & \cdots & \dfrac{1}{2\tau+1} \end{pmatrix}$$

is the well-known Hilbert matrix which is proved, in Section 2.5, to be positive definite. Also, the eigenvalues of $DGD$ are of the order $r$ and therefore the matrix $DGD$ will dominate the matrix formed by the solution of $\tilde{\tilde{\rho}}_{i,j}$ of the system (43). Thus, for $r$ sufficiently large, the matrix $R_\tau > 0$, which proves the theorem. ∎

EXAMPLE 1. For $\tau = 3$ the approximation must be fourth-order accurate in the interior. Choose $s = 2$, $\alpha_1 = \frac{2}{3}$, $\alpha_2 = -\frac{1}{12}$. A simple calculation shows that the $\rho_{i,j}$ with $0 \leqslant i \leqslant j$ are given by

$$\rho_{0,0} = M_{0,1}, \qquad \rho_{0,1} = \tfrac{1}{2} M_{0,2},$$

$$\rho_{0,2} = \tfrac{1}{3} M_{0,3}, \qquad \rho_{0,3} = M_{1,3} - \tfrac{3}{4} M_{2,2}$$

$$\rho_{1,1} = \tfrac{1}{2} M_{1,2} - \tfrac{1}{6} M_{0,3}, \qquad \rho_{1,2} = \tfrac{1}{4} M_{2,2},$$

$$\rho_{1,3} = \tfrac{1}{2} M_{2,3} - \tfrac{3}{2} \rho_{2,2}, \qquad \rho_{2,3} = \tfrac{1}{6} M_{3,3}$$

No further relations have to be satisfied. Therefore $\rho_{2,2}$, $\rho_{3,3}$ can be chosen arbitrarily and $R_\tau$ can be made such that $R_\tau > 0$, provided

$$\begin{pmatrix} \rho_{0,0} & \rho_{0,1} \\ \rho_{0,1} & \rho_{1,1} \end{pmatrix} = \begin{pmatrix} M_{0,1} & \tfrac{1}{2} M_{0,2} \\ \tfrac{1}{2} M_{0,2} & \tfrac{1}{2} M_{1,2} - \tfrac{1}{6} M_{0,3} \end{pmatrix} > 0.$$

From (28) expressions for the $M_{i,j}$ are obtained:

$$\begin{pmatrix} M_{0,0} & M_{0,1} & M_{0,2} & M_{0,3} \\ & M_{1,1} & M_{1,2} & M_{1,3} \\ & & M_{2,2} & M_{2,3} \\ & & & M_{3,3} \end{pmatrix}$$

$$= \begin{pmatrix} 0 & r - \tfrac{1}{2} & -r^2 + \tfrac{1}{6} & r^3 \\ & -r^2 + \tfrac{1}{6} & r^3 & -r^4 + \tfrac{1}{6} \\ & & -r^4 - \tfrac{1}{6} & r^5 \\ & & & -r^6 + \tfrac{1}{6} \end{pmatrix}.$$

Therefore,

$$\begin{pmatrix} M_{0,1} & \dfrac{1}{2} M_{0,2} \\ \dfrac{1}{2} M_{0,2} & \dfrac{1}{2} M_{1,2} - \dfrac{1}{6} M_{0,3} \end{pmatrix} = \begin{pmatrix} r - \dfrac{1}{2} & \dfrac{-r^2}{2} + \dfrac{1}{12} \\ \dfrac{-r^2}{2} + \dfrac{1}{12} & \dfrac{1}{3} r^3 \end{pmatrix}$$

which is positive definite for $r \geqslant 2$. $H$ is chosen in the form (40) and $r$ is set to $r = \tau + 2$. Then for every choice of $\lambda_0$, $\rho_{2,2}$, $\rho_{3,3}$ $H$ is uniquely determined. (The only restriction on $\lambda_0$, $\rho_{2,2}$, $\rho_{3,3}$ is that $\tilde{R}_\tau$ be positive definite.) After having determined $H$ one can obtain $B_2$ from the system (25).

## 2.4. *Interior Accuracy Conditions*

In this section the sought conditions of compatibility for the system (29) are derived. These conditions are expressed

in terms of accuracy conditions which the chosen scheme in the interior has to fulfill. For the special case when $H$ is a diagonal norm it can be shown that $J_{i,\sigma}$, defined in Eq. (28), depends only on $i + j$, which is done in Section 2.5.

By introducing the diagonal matrix

$$\Lambda = \begin{pmatrix} \lambda_0 & 0 & \cdots & \cdots & 0 \\ 0 & \lambda_1 & 0 & \cdots & 0 \\ \vdots & & \ddots & & \vdots \\ \vdots & & & \ddots & 0 \\ 0 & \cdots & \cdots & 0 & \lambda_{r-1} \end{pmatrix}, \quad (44)$$

and then by expanding Eq. (27), the following expression is obtained

$$j\langle e_i, \Lambda e_{j-1} \rangle = j \sum_{v=0}^{r-1} \lambda_v (r-v)^{i+j-1} (-1)^{i+j-1}$$

and

$$(j+i)(-1)^{i+j-1} \sum_{v=0}^{r-1} \lambda_v (r-v)^{i+j-1} = M_{i,j}. \quad (45)$$

This is the left-hand side of (27), and it is clear that this expression only depends on $\sigma = i + j$. Therefore the same must be true for the right-hand side, i.e., $J_{i,\sigma}$ in Eq. (28) depends only on $\sigma$. Then the following lemma can be given.

LEMMA 2.11.   *If the $M_{i,j}$ depend only on $i + j$, i.e.,*

$$M_{i,j} = R_{i+j}, \qquad \rho_0 = 0,$$

*then the system* (29) *has the solution*

$$\rho_{i,-1} = \rho_{-1,j} = 0,$$

$$\rho_{i,j} = \frac{R_{i+j+1}}{i+j+1}, \qquad 0 \leqslant i, \qquad j \leqslant \tau, \qquad i+j < 2\tau.$$

*Proof.*   Introduce the above expression for $\rho_{i,j}$ into (29). ∎

A consequence of this lemma is that if

$$M_{n,n} = M_{n-1,n+1} = \cdots = M_{0,2n} = R_{2n}$$

then the relation (36) must hold, i.e.,

$$\frac{1}{2n} - \frac{1}{n+1} + \frac{(n-1)}{(n+1)(n+2)} - \cdots$$

$$+ (-1)^n \frac{(n-1)\cdots 1}{(n+1)\cdots 2n}$$

$$= 0, \qquad n = 1, 2, \ldots; \quad \text{with} \quad 2n \leqslant \tau. \quad (46)$$

The $M_{i,j}$ are expressed in terms of $J_{i,\sigma}$ in formula (28). These in turn can be written as functions of sums of type $\sum_{v=1}^{s} \alpha_v v^m$. By introducing (28) into (36) one obtains

$$\frac{1}{2n} J_{n,2n} - \frac{1}{n+1} J_{n-1,2n} + \frac{(n-1)}{(n+1)(n+2)} J_{n-2,2n} - \cdots$$

$$+ (-1)^n \frac{(n-1)\cdots 1}{(n+1)\cdots 2n} J_{0,2n} - (-r)^{2n}$$

$$\times \left[ \frac{1}{2n} - \frac{1}{n+1} + \cdots + (-1)^n \frac{(n-1)\cdots 1}{(n+1)\cdots 2n} \right] = 0. \quad (47)$$

The last term is equal to zero by (46). Now let $N_{i,\sigma}(v)$ be defined by

$$N_{i,\sigma}(v) = \sum_{\mu=0}^{v-1} \mu^{\sigma-i}(\mu-v)^i + \sum_{\mu=1}^{v} (\mu-v)^i \mu^{\sigma-i},$$

$$0 \leqslant i \leqslant \sigma.$$

Then

$$N_{i+1,\sigma}(v) = \sum_{\mu=0}^{v-1} \mu^{\sigma-i-1}(\mu-v)^{i+1} + \sum_{\mu=1}^{v} (\mu-v)^{i+1} \mu^{\sigma-i-1}$$

$$= N_{i,\sigma}(v) - vN_{i,\sigma-1}(v) \quad (48)$$

and

$$J_{i,\sigma} = \sum_{v=1}^{s} \alpha_v N_{i,\sigma} \quad (49)$$

by the relation

$$\sum_{\mu=0}^{v-1} \mu^i(\mu-v)^{\sigma-i} = (-1)^\sigma \sum_{\mu=1}^{v} (\mu-v)^i \mu^{\sigma-i}. \quad (50)$$

The well-known formulas for $\sum \mu^\sigma$ give the expression

$$N_{0,\sigma}(v) = 2 \sum_{\mu=1}^{v} \mu^\sigma - v^\sigma$$

$$= 2 \begin{cases} (\sigma+1)^{-1} v^{\sigma+1} + \dfrac{1}{2} B_2 \dbinom{\sigma}{1} v^{\sigma-1} \\ \qquad + \dfrac{1}{4} B_4 \dbinom{\sigma}{3} v^{\sigma-3} + \cdots + B_\sigma v, \\ \qquad \text{for } \sigma \text{ even,} \\[6pt] (\sigma+1)^{-1} v^{\sigma+1} + \dfrac{1}{2} B_2 \dbinom{\sigma}{1} v^{\sigma-1} \\ \qquad + \dfrac{1}{4} B_4 \dbinom{\sigma}{3} v^{\sigma-3} + \cdots + \dfrac{\sigma}{2} B_{\sigma-1} v^2, \\ \qquad \text{for } \sigma \text{ odd,} \end{cases} \quad (51)$$

where $B_2$, $B_4$, ... are the Bernoulli numbers which can be found in [3]. Therefore, by induction using (48) the $N_{i,\sigma}(v)$ can be written as

$$N_{i,\sigma} = \gamma_1^{(i,\sigma)} v^{\sigma+1} + \gamma_3^{(i,\sigma)} v^{\sigma-1} + \cdots . \qquad (52)$$

Here the last term is proportional to $v$ or $v^2$, depending on whether $\sigma$ is even or odd, respectively. Furthermore, $\gamma_1^{(i,\sigma)} \neq 0$, because in the first approximation

$$N_{i,\sigma} \approx 2 \int_0^v x^{\sigma-i}(x-v)^i \, dx$$

$$= 2(-1)^i v^{\sigma+1} \int_0^1 x^{\sigma-i}(1-x)^i \, dx.$$

By using the derived expression for $N_{i,\sigma}$, and (49) for $\sigma$ even, the relation (47) is reduced to

$$\beta_{n,0} \sum_{v=1}^s \alpha_v v^{2n+1} + \beta_{n,2} \sum_{v=1}^s \alpha_v v^{2n-1} + \cdots + \beta_{n,2n-1} \sum_{v=1}^s \alpha_v v^3,$$

where

$$\beta_{n,0} = \frac{\gamma_1^{(n,2n)}}{2n} - \frac{\gamma_1^{(n-1,2n)}}{n-1} + \cdots \neq 0$$

because the $\gamma_1^{(i,2n)}$ are different from zero, and by (48) and (52) it is clear that they have alternate signs. Then induction gives the sought accuracy conditions for the interior scheme

$$\sum_{v=1}^s \alpha_v v^{2n+1} = 0, \qquad n = 1, 2, ...; \qquad \text{with} \quad 2n \leqslant \tau.$$

### 2.5. Diagonal Norms

When using diagonal norms the following relation, for the norm $A$, is obtained from (27) and (45)

$$(j+i)(-1)^{j+i-1} \sum_{v=0}^{r-1} \lambda_v(r-v)^{j+i-1}$$

$$= -(-r)^{j+i} + J_{i,j+i}, \qquad 0 \leqslant i, j \leqslant \tau$$

$$J_{i,\sigma} = \sum_{v=1}^s \alpha_v \left( \sum_{\mu=0}^{v-1} \mu^{\sigma-i}(\mu-v)^i + \mu^i(\mu-v)^{\sigma-i} \right),$$

$$\sigma = i + j. \qquad (53)$$

The left-hand side of the relation above depends only on $\sigma = i + j$. Therefore the same must be true for the right-hand side, i.e.,

$$J_{i+1,\sigma} = J_{i,\sigma}, \qquad 0 < i+1 \leqslant \sigma, \qquad 0 < i+1 \leqslant \tau. \quad (54)$$

This is easy to see, because using Eqs. (48) and (49) from Section 2.4 for $\sigma$ even, gives the expression

$$J_{i+1,\sigma} = J_{i,\sigma} - \sum_{v=1}^s \alpha_v v N_{i,\sigma-1}(v). \qquad (55)$$

Thus by (52)

$$\sum_{v=1}^s \alpha_v v N_{i,\sigma-1}(v) = \gamma_1^{(i,\sigma-1)} \sum_{v=1}^s \alpha_v v^{\sigma+1} + \cdots, \quad (56)$$

where the last term is proportional to $v^3$ because $\sigma - 1$ is odd. But (56) was proved to be zero in Section 2.4, and therefore the relation (54) is valid for $\sigma$ even; $\sigma = 1$ gives $i = 0$ and the relation (54) is no condition. Thus,

$$J_{0,1} = \sum_{v=1}^s \alpha_v \left( \sum_{\mu=0}^{v-1} \mu + (\mu-v) \right)$$

$$= \sum_{v=1}^s \alpha_v(v^2 - v - v^2) = -\tfrac{1}{2}. \qquad (57)$$

For odd $\sigma > 1$, the relation (50) implies that

$$J_{i,\sigma} = \sum_{v=1}^s \alpha_v \left( \sum_{\mu=0}^{v-1} \mu^{\sigma-i}(\mu-v)^i - \sum_{\mu=1}^v \mu^{\sigma-i}(\mu-v)^i \right)$$

$$= \sum_{v=1}^s \alpha_v(0^{\sigma-i}(-v)^i - v^{\sigma-i}0^i) = 0. \qquad (58)$$

Therefore (54) also holds for odd $\sigma$. Thus Eq. (53) is equivalent to

$$\sum_{v=0}^{r-1} \lambda_v(r-v)^{\sigma-1} = \frac{1}{\sigma}((-r)^\sigma - (-1)^\sigma k_\sigma),$$

$$\sigma = 1, 2, ..., 2\tau, \qquad (59)$$

where by (57) and (51)

$$k_\sigma = \begin{cases} \sum_{v=1}^s \alpha_v N_{0,\sigma}(v) = 2B_\sigma \sum_{v=1}^s \alpha_v v = B_\sigma, \\ \qquad \sigma = 2, 4, ..., 2\tau - 2, \\ 2 \sum_{v=1}^s \alpha_v \sum_{\mu=0}^{v-1} \mu^\tau(\mu-v)^\tau, \qquad \sigma = 2\tau, \\ -\tfrac{1}{2}, \qquad \sigma = 1, \\ 0, \qquad \sigma = 3, 5, ..., 2\tau - 1. \end{cases} \qquad (60)$$

If $r$ is set of $2\tau$, then (59) and (60) define a linear system of equations which has a unique solution.

The following theorem is given.

THEOREM 2.3. *Assume that $r = 2\tau$ and that the system of Eq. (59) has a positive solution. Then there is a scalar product*

(7), *with the diagonal norm $\Lambda$ instead of $H$, and a difference approximation $Q$ which is accurate to order $\tau$ near the boundary and accurate to order $2\tau$ in the interior such that (54) holds.*

## 2.6. *Positive Definiteness of the Hilbert Matrix*

The definiteness of the Hilbert matrix is easy to prove because the elements can be expressed as in integral on $[0, 1]$:

$$g_{i,j} = \int_0^1 t^{i+j-2}\, dt.$$

Let $x = (x_0, ..., x_\tau)^T$, then the definiteness condition $x^T G x > 0$ becomes

$$
\begin{aligned}
x^T G x &= \sum_{i,j=1}^{\tau+1} x_{i-1} g_{ij} x_{j-1} = \int_0^1 \sum_{i,j=1}^{\tau+1} x_{i-1} t^{i+j-2} x_{j-1}\, dt \\
&= \int_0^1 \sum_{j=1}^{\tau+1} \sum_{i=1}^{\tau+1} x_{i-1} t^{i-1} x_{j-1} t^{j-1}\, dt \\
&= \int_0^1 \sum_{i=0}^{\tau} x_i t^i \sum_{j=0}^{\tau} x_j t^j\, dt \\
&= \int_0^1 |p_\tau(t)|^2\, dt \geq 0.
\end{aligned}
$$

It is clear that this implies that $x^T G x > 0$. The relation $x^T G x = 0$ implies that $p_\tau(t) \equiv 0$. However, $\{1, ..., t^\tau\}$ is linearly independent on $[0, 1]$, so $x$ has to satisfy $x_0 = x_1 = \cdots = x_\tau = 0$. Thus it is proved that the Hilbert matrix $G$ is positive definite.

## 3. HIGHER ORDER DIFFERENCE APPROXIMATIONS

In the interior an antisymmetric stencil is used to obtain the highest possible order of accuracy. This is achieved by using $2s + 1$ symmetric centered points. The difference approximation in the interior is accurate to order $2s$ and fulfills the conditions of Lemma 2.2. Then by solving the system of Eqs. (18) we obtain the interior stencil:



Depending on whether we use full norms, restricted full norms, or diagonal norms, the relation between the order of accuracy at the boundary and in the interior will differ for the difference operator $Q$. Properties such as structure and

size of the boundary part of the operators will also differ, so also will the number of arbitrary parameters in the multi-parameter families of the higher order difference approximations.

When using *full norms* (39) we are able to find, if we set $r = \tau + 1$ according to Lemma 2.9, a norm, $H$, defined by (38) in Section 2.3. If we can specify those elements of $\rho_{v,\tau}$, $v = 0, 1, ..., \tau$, not defined by the system (29), such that the matrix $R_\tau$ becomes positive definite, then the norm $H$ is also positive definite. We should note that when $\tau$ is odd then $v$ is even. Therefore there will be two elements to specify when $\tau = 3$, and three elements when $\tau = 5$. Then according to Theorem 2.2, if we set $\tau + 1 = 2s$, the difference operator $Q$ is determined uniquely from the overdetermined system (25) with the order of accuracy $\tau$ at the boundary and $\tau + 1$ in the interior.

However, we want to compute the general form of the difference operator $Q$, i.e., express it in the arbitrary elements $\rho_{v,\tau}$. This is done simply by inserting the arbitrary elements $x_1, x_2, ...$, not determined by the system (29), in $R_\tau$, and then solve for $H$ and $Q$ using a symbolic language such as Maple. We then obtain a multi-parameter family of operators, of which the parameters can be chosen to minimize the bandwidth of the difference operator or in some other way. When the parameters have been chosen the eigenvalues must be evaluated to check the positive definiteness of the corresponding norm.

The structure of the boundary part and the first interior point for the full norm, $H$, and the corresponding difference operator, $Q$, for the case of accuracy three at the boundary, is shown below. The non-zero boundary part of the operator has in general the size $(\tau + 1) \times (3(\tau + 1)/2)$,



In case of *restricted full norms* we set $r = \tau + 2$ according to Lemma 2.10. Then one can choose $H$ in the form (40) such that (37) holds. The structure for the norm and the operator for the case of accuracy three at the boundary is

$$H = \begin{pmatrix} x & & & & & \\ & x & x & x & x & \\ & x & x & x & x & \\ & x & x & x & x & \\ & x & x & x & x & \\ & & & & & 1 \\ & & & & & & \ddots \end{pmatrix},$$

$$Q = \begin{pmatrix} x & x & x & x & x & & & \\ x & x & x & x & x & x & x & \\ x & x & x & x & x & x & x & \\ x & x & x & x & x & x & x & \\ x & x & x & x & x & x & x & \\ & & x & x & 0 & x & x & \\ & & & \ddots & \ddots & \ddots & \ddots & \ddots \end{pmatrix}.$$

The restricted full norms have the advantage over the full norm that they can be used to show stability [4] for first-order systems.

If *diagonal norms* are used, Theorem 2.3 states that if $r = 2\tau$ and the system (60) has a positive solution (which it has when $\tau \leqslant 4$), then there is a difference approximation $Q$ which is accurate to order $\tau$ near the boundary and accurate to order $2\tau$ in the interior. However, $Q$ is not determined uniquely when $\tau \geqslant 3$. The system (25) turns out to give an infinite set of solutions. Therefore to evaluate a specific antisymmetric boundary matrix $B_2$ we instead solve the reduced system,

$$\begin{pmatrix} g_{11} & \cdots & g_{1n-i-1} \\ \vdots & & \vdots \\ g_{m1} & \cdots & g_{mn-i-1} \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ b_{n-i-1} \end{pmatrix}$$

$$= \begin{pmatrix} z_1 \\ \vdots \\ z_{n-i-1} \end{pmatrix} - \sum_{k=0}^{i} b_{n-k} G(:, n-k), \qquad (61)$$

where the $b_n, ..., b_{n-i}$ are chosen arbitrarily and $i$ depends on the order of accuracy, with $i = 1$ when $\tau = 3$ and $i = 3$ when $\tau = 4$. The size of the non-zero boundary part of the operator $Q$ becomes $(2\tau \times 3\tau)$, with the elements $Q(1:\tau, 2\tau:3\tau)$ equal to zero. Below the structure of the norm, $A$, and the corresponding difference approximation, $Q$, which is third-order accurate at and near the boundary and sixth-order accurate in the interior, is shown.

$$A = \begin{pmatrix} x & & & & & & \\ & x & & & & & \\ & & x & & & & \\ & & & x & & & \\ & & & & x & & \\ & & & & & x & \\ & & & & & & 1 \\ & & & & & & & \ddots \end{pmatrix}$$

## 4. RESULTS

The difference approximations, $Q$, of $d/dx$ that have been evaluated, together with norms and interior stencils, are described below. A more detailed description and the exact elements of the difference operators together with norms and interior stencils are given in Appendices A and B.

### 4.1. *Diagonal Norms*

As has been told earlier, difference approximations which are stable in a diagonal norm, have been evaluated with the accuracy of order $\tau = 1, ..., 4$ at the boundary, and $2\tau$ in the interior. They are listed in their complete form in Appendix A.

With $\tau = 1$ and 2 at the boundary, the operators are unique.

When $\tau = 3$ we obtain a one-parameter family of operators. By determining the parameter such that $q_{0,5}$ is zeroed, an operator with minimized bandwidth is computed. The resulting operator has nine non-zero diagonals, four superdiagonals, four subdiagonals, and the main diagonal.

Finally, $\tau = 4$ gives a three-parameter family of difference approximations. Here an operator with minimized bandwidth is computed by solving the linear system of equations $q_{0,6} = q_{0,7} = q_{1,7} = 0$, thus giving the operator a total bandwidth of eleven.

### 4.2. *Full Norms*

Difference approximations which satisfy the summation by parts criterion (9) with respect to full norms, with $\tau = 3, 5$ at the boundary, and $\tau + 1$ in the interior, have been computed, of which the case $\tau = 3$ is listed in Appendix B.

For $\tau = 3$ we have a two-parameter family of norms and difference operators. As earlier we obtain a version with minimized bandwidth of the operator by determining the parameters such that $q_{0,5}$ and $q_{1,5}$ are zeroed. This leads to solving a non-linear system of equations, of which one solution gives a positive definite norm. The corresponding boundary part of the difference operator will have four superdiagonals, three subdiagonals, and the main diagonal. However, if we have a finite computational domain we would like to use the boundary part of the operator at

the upper boundary as well. This can be done simply by reflecting the stencils of the lower boundary. Thus, at the upper boundary we will have four subdiagonals, three superdiagonals, and the main diagonal. Hence, the resulting difference operator would then have totally nine non-zero diagonals.

In the case we have accuracy five at the boundary and six in the interior we obtain a three-parameter family of difference operators. However, the general form of the elements of the operator become so enormous that their practical use must be considered as doubtful. An exact version with optimal bandwidth could not be computed because the obtained non-linear system of equations consists of the arbitrary parameters raised to a power of five.

### 4.3. *Restricted Full Norms*

Difference approximations with $\tau = 3$ at the boundary and $\tau + 1$ in the interior, have been computed and are listed in Appendix B.

This leads to a three-parameter family of norms and operators, and we therefore have one degree of freedom more than for the full norm case, but at the cost of having an extra point to be boundary modified. The parameters are then chosen such that $q_{0,4}$, $q_{1,6}$, and $q_{2,6}$ are zeroed. The resulting difference operator will have a total bandwidth of nine, i.e., the same as we had for the diagonal norms and full norms.

## 5. NUMERICAL EXPERIMENTS

Our theoretical analysis is made for a simple semi-discrete scalar problem. For a study of the numerical behaviour of these difference operators in more general situations we refer to the Ph.D. thesis of Pelle Olsson [2]. In [2] the numerical behaviour of a fourth-order accurate method is compared with a second-order method. The difference operators used are the ones that fulfill the summation by parts rule with respect to a diagonal norm, and which are third-order accurate at and near the boundary, and sixth-order accurate in the interior. The two-dimensional non-linear Euler equations with a forcing function are used to show that the convergence rate to the exact solution is fourth order globally. To measure the efficiency of the fourth-order method versus the second-order method the

simple model equation $u_t + u_x = 0$ is used in [2]. The comparison is done by computing the relative error for the numerical solution and the exact solution and by comparing the consumed CPU time for the two methods. Finally, the two-dimensional non-linear Euler equations over a backward-facing step, are solved on an CM 200, and the consumed CPU time are compared for the fourth- and the second-order methods. The conclusion in [2] is that if very high accuracy is needed, the fourth-order method would be the preferred choice.

## 6. CONCLUSIONS

The results regarding the difference approximations that fulfill the summation by parts criterion (8) and (9), using a *diagonal norm*, turned out to be better and more useful than was expected. This is true when parallel computers are used because, instead of obtaining a unique difference operator we obtained one- and three-parameter sets of solutions for accuracy of orders three and four at the boundary. The parameters were then determined to minimize the bandwidth. For the *full norms* and the *restricted full norms* we obtained multi-parameter families of difference approximations when third-order accuracy or higher was required at the boundary. The parameters should be determined such that the norm (in the full norm and restricted full norm cases) becomes positive definite. If we use parallel computers it is of interest to minimize the bandwidth of the difference operator since it determines the memory requirement and also the amount of computational work. Therefore, if we want third-order accuracy at the boundary all three cases give us a bandwidth of nine, but the diagonal norm gives us an interior stencil which is sixth-order accurate without any extra cost in memory or computation. If we want accuracy five or more at the boundary we have to consider full or restricted full norms. For scalar hyperbolic equations they can both be used to show stability, but if hyperbolic first-order systems are considered, there are no stability results for the full norm case and the restricted full norm should be used, see [4]. However, because the numerator and the denominator of the elements in the exact representation turn out to be very large numbers, it is doubtful if the difference operator satisfies the summation by parts energy norm in finite precision arithmetic.

## APPENDIX A: DIAGONAL NORMS

Here we present the difference operators, which satisfy the summation by parts energy norm (9), with corresponding norms and interior stencils. The boundary part of the operators, $Q = [A^{-1}BA^{-1}C]$, has the size $(2\tau \times 3\tau)$ and are accurate to order $\tau$. The antisymmetric interior stencil uses $2\tau + 1$ symmetric centered points and is accurate to order $2\tau$.

*First-order accuracy at the boundary.*

Interior stencil. $\alpha_1 = \frac{1}{2}$.

Norm. $\lambda_0 = \frac{1}{2}, \lambda_1 = 1$.

Boundary operator. $q_{0,0} = -1, q_{0,1} = 1, q_{0,2} = 0; q_{1,0} = -\frac{1}{2}, q_{1,1} = 0, q_{1,2} = \frac{1}{2}$.

*Second-order accuracy at the boundary.*

Interior stencil. $\alpha_1 = \frac{2}{3}$, $\alpha_2 = -\frac{1}{12}$.

Norm. $\lambda_0 = \frac{17}{48}$, $\lambda_1 = \frac{59}{48}$, $\lambda_2 = \frac{43}{48}$, $\lambda_3 = \frac{49}{48}$.

Boundary operator. $q_{0,0} = -\frac{24}{17}$, $q_{0,1} = \frac{59}{34}$, $q_{0,2} = -\frac{4}{17}$, $q_{0,3} = -\frac{3}{34}$, $q_{0,4} = 0$, $q_{0,5} = 0$; $q_{1,0} = -\frac{1}{2}$, $q_{1,1} = 0$, $q_{1,2} = \frac{1}{2}$, $q_{1,3} = 0$, $q_{1,4} = 0$, $q_{1,5} = 0$; $q_{3,0} = \frac{3}{98}$, $q_{3,1} = 0$, $q_{3,2} = -\frac{59}{98}$, $q_{3,3} = 0$, $q_{3,4} = \frac{32}{49}$, $q_{3,5} = -\frac{4}{49}$.

*Third-order accuracy at the boundary.*

Interior stencil. $\alpha_1 = \frac{3}{4}$, $\alpha_2 = -\frac{3}{20}$, $\alpha_3 = \frac{1}{60}$.

Norm. $\lambda_0 = \frac{13649}{43200}$, $\lambda_1 = \frac{12013}{8640}$, $\lambda_2 = \frac{2711}{4320}$, $\lambda_3 = \frac{5359}{4320}$, $\lambda_4 = \frac{7877}{8640}$, $\lambda_5 = \frac{43801}{43200}$.

Boundary operator. $q_{0,0} = -\frac{21600}{13649}$, $q_{0,1} = 8(16200x_1 - 953)/40947$, $q_{0,2} = (-1036800x_1 + 715489)/81894$, $q_{0,3} = 3(86400x_1 - 62639)/13649$, $q_{0,4} = 5(-207360x_1 + 147127)/81894$, $q_{0,5} = (129600x_1 - 89387)/40947$, $q_{0,6} = 0$, $q_{0,7} = 0$, $q_{0,8} = 0$;

$q_{1,0} = 8(-16200x_1 + 953)/180195$, $q_{1,1} = 0$, $q_{1,2} = (86400x_1 - 57139)/12013$, $q_{1,3} = (-1036800x_1 + 745733)/72078$, $q_{1,4} = 5(25920x_1 - 18343)/12013$, $q_{1,5} = (-345600x_1 + 240569)/120130$, $q_{1,6} = 0$, $q_{1,7} = 0$, $q_{1,8} = 0$;

$q_{2,0} = (1036800x_1 - 715489)/162660$, $q_{2,1} = (-86400x_1 + 57139)/5422$, $q_{2,2} = 0$, $q_{2,3} = (259200x_1 - 176839)/8133$, $q_{2,4} = (-345600x_1 + 242111)/10844$, $q_{2,5} = (259200x_1 - 182261)/27110$, $q_{2,6} = 0$, $q_{2,7} = 0$, $q_{2,8} = 0$;

$q_{3,0} = 3(-86400x_1 + 62639)/53590$, $q_{3,1} = (1036800x_1 - 745733)/64308$, $q_{3,2} = (-259200x_1 + 176839)/16077$, $q_{3,3} = 0$, $q_{3,4} = (259200x_1 - 165041)/32154$, $q_{3,5} = (-1036800x_1 + 710473)/321540$, $q_{3,6} = \frac{72}{5359}$, $q_{3,7} = 0$, $q_{3,8} = 0$;

$q_{4,0} = (207360x_1 - 147127)/47262$, $q_{4,1} = 5(-25920x_1 + 18343)/7877$, $q_{4,2} = (345600x_1 - 242111)/15754$, $q_{4,3} = (-259200x_1 + 165041)/23631$, $q_{4,4} = 0$, $q_{4,5} = 8640x_1/7877$, $q_{4,6} = -\frac{1296}{7877}$, $q_{4,7} = \frac{144}{7877}$, $q_{4,8} = 0$;

$q_{5,0} = (-129600x_1 + 89387)/131403$, $q_{5,1} = (345600x_1 - 240569)/87602$, $q_{5,2} = (-259200x_1 + 182261)/43801$, $q_{5,3} = (1036800x_1 - 710473)/262806$, $q_{5,4} = -43200x_1/43801$, $q_{5,5} = 0$, $q_{5,6} = \frac{32400}{43801}$, $q_{5,7} = -\frac{6480}{43801}$, $q_{5,8} = \frac{720}{43801}$.

*Fourth-order accuracy at the boundary.*

Interior stencil. $\alpha_1 = \frac{4}{5}$, $\alpha_2 = -\frac{1}{5}$, $\alpha_3 = \frac{4}{105}$, $\alpha_4 = -\frac{1}{280}$.

Norm. $\lambda_0 = 1498139/5080320$, $\lambda_1 = 1107307/725760$, $\lambda_2 = 20761/80640$, $\lambda_3 = 1304999/725760$, $\lambda_4 = 299527/725760$, $\lambda_5 = 103097/80640$, $\lambda_6 = 670091/725760$, $\lambda_7 = 5127739/5080320$.

Boundary operator. $q_{0,0} = -2540160/1498139$, $q_{0,1} = 9(2257920x_1 + 11289600x_2 + 22579200x_3 - 15849163)/5992556$, $q_{0,2} = 3(-33868800x_1 - 162570240x_2 - 304819200x_3 + 235236677)/5992556$, $q_{0,3} = (609638400x_1 + 2743372800x_2 + 4572288000x_3 - 3577778591)/17977668$, $q_{0,4} = 3(-16934400x_1 - 67737600x_2 - 84672000x_3 + 67906303)/1498139$, $q_{0,5} = 105(967680x_1 + 2903040x_2 - 305821)/5992556$, $q_{0,6} = 49(-1244160x_1 + 18661400x_3 - 13322233)/17977668$, $q_{0,7} = 3(-6773760x_2 - 33868800x_3 + 24839327)/5992556$, $q_{0,8} = 0$, $q_{0,9} = 0$, $q_{0,10} = 0$, $q_{0,11} = 0$;

$q_{1,0} = 9(-2257920x_1 - 11289600x_2 - 22579200x_3 + 15849163)/31004596$, $q_{1,1} = 0$, $q_{1,2} = 3(7257600x_1 + 33868800x_2 + 60963840x_3 - 47167457)/2214614$, $q_{1,3} = 3(-9676800x_1 - 42336000x_2 - 67737600x_3 + 53224573)/1107307$, $q_{1,4} = 7(55987200x_1 + 217728000x_2 + 261273600x_3 - 211102099)/13287684$, $q_{1,5} = 3(-11612160x_1 - 33868800x_2 + 3884117)/2214614$, $q_{1,6} = 150(24192x_1 - 338688x_3 + 240463)/1107307$, $q_{1,7} = (152409600x_2 + 731566080x_3 - 536324953)/46506894$, $q_{1,8} = 0$, $q_{1,9} = 0$, $q_{1,10} = 0$, $q_{1,11} = 0$;

$q_{2,0} = (33868800x_1 + 162570240x_2 + 304819200x_3 - 235236677)/1743924$, $q_{2,1} = (-7257600x_1 - 33868800x_2 - 60963840x_3 + 47167457)/124566$, $q_{2,2} = 0$, $q_{2,3} = (24192000x_1 + 101606400x_2 + 152409600x_3 - 120219461)/124566$, $q_{2,4} = (-72576000x_1 - 270950400x_2 - 304819200x_3 + 249289259)/249132$, $q_{2,5} = 9(806400x_1 + 2257920x_2 - 290167)/41522$, $q_{2,6} = 6(-134400x_1 + 1693440x_3 - 1191611)/20761$, $q_{2,7} = 5(-2257920x_2 - 10160640x_3 + 7439833)/290654$, $q_{2,8} = 0$, $q_{2,9} = 0$, $q_{2,10} = 0$, $q_{2,11} = 0$;

$q_{3,0} = (-609638400x_1 - 2743372800x_2 - 4572288000x_3 + 3577778591)/109619916$, $q_{3,1} = 3(9676800x_1 + 42336000x_2 + 67737600x_3 - 53224573)/1304999$, $q_{3,2} = 3(-24192000x_1 - 101606400x_2 - 152409600x_3 + 120219461)/2609998$, $q_{3,3} = 0$, $q_{3,4} = 9(16128000x_1 + 56448000x_2 + 56448000x_3 - 47206049)/5219996$, $q_{3,5} = 3(-19353600x_1 - 50803200x_2 + 7628371)/2609998$, $q_{3,6} = 2(10886400x_1 - 114307200x_3 + 79048289)/3914997$, $q_{3,7} = 75(1354752x_2 + 5419008x_3 - 3952831)/18269986$, $q_{3,8} = 0$, $q_{3,9} = 0$, $q_{3,10} = 0$, $q_{3,11} = 0$;

$q_{4,0} = 3(16934400x_1 + 67737600x_2 + 84672000x_3 - 67906303)/2096689$, $q_{4,1} = 7(-55987200x_1 - 217728000x_2 - 261273600x_3 + 211102099)/3594324$, $q_{4,2} = 3(72576000x_1 + 270950400x_2 + 304819200x_3 - 249289259)/1198108$, $q_{4,3} = 9(-16128000x_1 - 56448000x_2 - 56448000x_3 + 47206049)/1198108$, $q_{4,4} = 0$, $q_{4,5} = 105(414720x_1 + 967680x_2 - 165527)/1198108$, $q_{4,6} = 15(-967680x_1 + 6773760x_3 - 4472029)/1198108$, $q_{4,7} = (-304819200x_2 - 914457600x_3 + 657798011)/25160268$, $q_{4,8} = -2592/299527$, $q_{4,9} = 0$, $q_{4,10} = 0$, $q_{4,11} = 0$;

$q_{5,0} = 5(-967680x_1 - 2903040x_2 + 305821)/1237164$, $q_{5,1} = (11612160x_1 + 33868800x_2 - 3884117)/618582$, $q_{5,2} = 9(-806400x_1 - 2257920x_2 + 290167)/206194$, $q_{5,3} = (19353600x_1 + 50803200x_2 - 7628371)/618582$, $q_{5,4} = 35(-414720x_1 - 967680x_2 + 165527)/1237164$, $q_{5,5} = 0$, $q_{5,6} = 80640x_1/103097$, $q_{5,7} = 80640x_2/103097$, $q_{5,8} = 3072/103097$, $q_{5,9} = -288/103097$, $q_{5,10} = 0$, $q_{5,11} = 0$;

$q_{6,0} = 7(1244160x_1 - 18662400x_3 + 13322233)/8041092$, $q_{6,1} = 150(-24192x_1 + 338688x_3 - 240463)/670091$, $q_{6,2} = 54(134400x_1 - 1693440x_3 + 1191611)/670091$, $q_{6,3} = 2(-10886400x_1 + 114307200x_3 - 79048289)/2010273$, $q_{6,4} = 15(967680x_1 - 6773760x_3 + 4472029)/2680364$, $q_{6,5} = -725760x_1/670091$, $q_{6,6} = 0$, $q_{6,7} = 725760x_3/670091$, $q_{6,8} = -145152/670091$, $q_{6,9} = 27648/670091$, $q_{6,10} = -2592/670091$, $q_{6,11} = 0$;

$q_{7,0} = 3(6773760x_2 + 33868800x_3 - 24839327)/20510956$, $q_{7,1} = (-152409600x_2 - 731566080x_3 + 536324953)/30766434$, $q_{7,2} = 45(2257920x_2 + 10160640x_3 - 7439833)/10255478$, $q_{7,3} = 75(-1354752x_2 - 5419008x_3 + 3952831)/10255478$, $q_{7,4} = (304819200x_2 + 914457600x_3 - 657798011)/61532868$, $q_{7,5} = -5080320x_2/5127739$, $q_{7,6} = -5080320x_3/5127739$, $q_{7,7} = 0$, $q_{7,8} = 4064256/5127739$, $q_{7,9} = -1016064/5127739$, $q_{7,10} = 193536/5127739$, $q_{7,11} = -18144/5127739$.

## A.1. *Minimum Bandwidth Operators*

Since the third-order accurate difference operator is really a one-parameter family of operators, and the fourth-order accurate difference opertor is a three-parameter family, the parameters can be used to minimize the bandwidth.

The minimization of the bandwidth is of interest when parallel computers are used. In the case of a third-order accurate operator we therefore choose $x_1$ such that $q_{0,5}$ is zeroed. The resulting operator will then have four non-zero super-diagonals, four subdiagonals, and the main diagonal. Thus, a total of nine non-zero diagonals are needed to store. For the case of accuracy four at the boundary, $x_1, x_2, x_3$ are determined such that $q_{0,6}$, $q_{0,7}$, $q_{1,7}$ are zeroed. Thus, resulting in an operator with a total of eleven non-zero diagonals. The operators with norms and interior stencils are given below.

*Third-order accuracy at the boundary, minimized bandwidth, $q_{0,5} = 0$.*

Interior stencil. $\alpha_1 = \frac{3}{4}$, $\alpha_2 = -\frac{3}{20}$, $\alpha_3 = \frac{1}{60}$.

Norm. $\lambda_0 = \frac{13649}{43200}$, $\lambda_1 = \frac{12013}{8640}$, $\lambda_2 = \frac{2711}{4320}$, $\lambda_3 = \frac{5359}{4320}$, $\lambda_4 = \frac{7877}{8640}$, $\lambda_5 = \frac{43801}{43200}$.

Boundary operator, minimized bandwidth. $q_{0,0} = -21600/13649$, $q_{0,1} = 81763/40947$, $q_{0,2} = 131/27298$, $q_{0,3} = -9143/13649$, $q_{0,4} = 20539/81894$, $q_{0,5} = 0$, $q_{0,6} = 0$, $q_{0,7} = 0$, $q_{0,8} = 0$;

$q_{1,0} = -81763/180195$, $q_{1,1} = 0$, $q_{1,2} = 7357/36039$, $q_{1,3} = 30637/72078$, $q_{1,4} = -2328/12013$, $q_{1,5} = 6611/360390$, $q_{1,6} = 0$, $q_{1,7} = 0$, $q_{1,8} = 0$;

$q_{2,0} = -131/54220$, $q_{2,1} = -7357/16266$, $q_{2,2} = 0$, $q_{2,3} = 645/2711$, $q_{2,4} = 11237/32532$, $q_{2,5} = -3487/27110$, $q_{2,6} = 0$, $q_{2,7} = 0$, $q_{2,8} = 0$;

$q_{3,0} = 9143/53590$, $q_{3,1} = -30637/64308$, $q_{3,2} = -645/5359$, $q_{3,3} = 0$, $q_{3,4} = 13733/32154$, $q_{3,5} = -\frac{67}{4660}$, $q_{3,6} = \frac{72}{5359}$, $q_{3,7} = 0$, $q_{3,8} = 0$;

$q_{4,0} = -20539/236310$, $q_{4,1} = 2328/7877$, $q_{4,2} = -11237/47262$, $q_{4,3} = -13733/23631$, $q_{4,4} = 0$, $q_{4,5} = 89387/118155$, $q_{4,6} = -\frac{1296}{7877}$, $q_{4,7} = \frac{144}{7877}$, $q_{4,8} = 0$;

$q_{5,0} = 0$, $q_{5,1} = -6611/262806$, $q_{5,2} = 3487/43801$, $q_{5,3} = 1541/87602$, $q_{5,4} = -89387/131403$, $q_{5,5} = 0$, $q_{5,6} = 32400/43801$, $q_{5,7} = -6480/43801$, $q_{5,8} = 720/43801$.

*Fourth-order accuracy at the boundary, minimized bandwidth, $q_{0,6} = q_{0,7} = q_{1,7} = 0$.*

Interior stencil. $\alpha_1 = \frac{4}{5}$, $\alpha_2 = -\frac{1}{5}$, $\alpha_3 = \frac{4}{105}$, $\alpha_4 = -\frac{1}{280}$.

Norm. $\lambda_0 = 1498139/5080320$, $\lambda_1 = 1107307/725760$, $\lambda_2 = 20761/80640$, $\lambda_3 = 1304999/725760$, $\lambda_4 = 299527/725760$, $\lambda_5 = 103097/80640$, $\lambda_6 = 670091/725760$, $\lambda_7 = 5127739/5080320$.

Boundary operator, minimized bandwidth. $q_{0,0} = -2540160/1498139$, $q_{0,1} = 37052897/17977668$, $q_{0,2} = 7891273/8988834$, $q_{0,3} = -7624221/2996278$, $q_{0,4} = 15181679/8988834$, $q_{0,5} = -6971555/17977668$, $q_{0,6} = 0$, $q_{0,7} = 0$, $q_{0,8} = 0$, $q_{0,9} = 0$, $q_{0,10} = 0$, $q_{0,11} = 0$;

$q_{1,0} = -5293271/13287684$, $q_{1,1} = 0$, $q_{1,2} = -2931787/6643842$, $q_{1,3} = 12616429/6643842$, $q_{1,4} = -6970881/4429228$, $q_{1,5} = 4026475/6643842$, $q_{1,6} = -101360/1107307$, $q_{1,7} = 0$, $q_{1,8} = 0$, $q_{1,9} = 0$, $q_{1,10} = 0$, $q_{1,11} = 0$;

$q_{2,0} = -607021/603666$, $q_{2,1} = 2931787/1121094$, $q_{2,2} = 0$, $q_{2,3} = -23428253/2242188$, $q_{2,4} = 9447614/560547$, $q_{2,5} = -217571/20761$, $q_{2,6} = 2847947/1121094$, $q_{2,7} = -1185475/15695316$, $q_{2,8} = 0$, $q_{2,9} = 0$, $q_{2,10} = 0$, $q_{2,11} = 0$;

$q_{3,0} = 7624221/18269986$, $q_{3,1} = -12616429/7829994$, $q_{3,2} = 23428253/15659988$, $q_{3,3} = 0$, $q_{3,4} = -2184329/1304999$, $q_{3,5} = 7700062/3914997$, $q_{3,6} = -3371361/5219996$, $q_{3,7} = 2783695/54809958$, $q_{3,8} = 0$, $q_{3,9} = 0$, $q_{3,10} = 0$, $q_{3,11} = 0$;

$q_{4,0} = -15181679/12580134$, $q_{4,1} = 6970881/1198108$, $q_{4,2} = -9447614/898581$, $q_{4,3} = 2184329/299527$, $q_{4,4} = 0$, $q_{4,5} = -10921405/3594324$, $q_{4,6} = 3604685/1797162$, $q_{4,7} = -1462269/4193378$, $q_{4,8} = -2592/299527$, $q_{4,9} = 0$, $q_{4,10} = 0$, $q_{4,11} = 0$;

$q_{5,0} = 6971555/77941332$, $q_{5,1} = -4026475/5567238$, $q_{5,2} = 217571/103097$, $q_{5,3} = -7700062/2783619$, $q_{5,4} = 10921405/11134476$, $q_{5,5} = 0$, $q_{5,6} = 1714837/5567238$, $q_{5,7} = -1022551/38970666$, $q_{5,8} = 3072/103097$, $q_{5,9} = -288/103097$, $q_{5,10} = 0$, $q_{5,11} = 0$;

$q_{6,0} = 0$, $q_{6,1} = 101360/670091$, $q_{6,2} = -2847947/4020546$, $q_{6,3} = 3371361/2680364$, $q_{6,4} = -3604685/4020546$, $q_{6,5} = -1714837/4020546$, $q_{6,6} = 0$, $q_{6,7} = 6445687/8041092$, $q_{6,8} = -145152/670091$, $q_{6,9} = 27648/670091$, $q_{6,10} = -2592/670091$, $q_{6,11} = 0$;

$q_{7,0} = 0$, $q_{7,1} = 0$, $q_{7,2} = 1185475/61532868$, $q_{7,3} = -2783695/30766434$, $q_{7,4} = 1462269/10255478$, $q_{7,5} = 1022551/30766434$, $q_{7,6} = -45119809/61532868$, $q_{7,7} = 0$, $q_{7,8} = 4064256/5127739$, $q_{7,9} = -1016064/5127739$, $q_{7,10} = 193536/5127739$, $q_{7,11} = -18144/5127739$.

## APPENDIX B: FULL NORMS AND RESTRICTED FULL NORMS

Difference approximations with corresponding full and restricted full norms and interior stencils are presented in their general form. Also, difference operators with minimized bandwidth are given.

### 8.1. *Full Norms*

The boundary part of the operators, $Q = [H^{-1}BH^{-1}C]$, has the size $((\tau + 1) \times 3(\tau + 1)/2)$ and are accurate to order $\tau$. The antisymmetric interior stencil uses $\tau + 2$ symmetric centered points and is accurate to order $\tau + 1$.

*Third-order accuracy at the boundary.*

Interior stencil. $\alpha_1 = \frac{2}{3}$, $\alpha_2 = -\frac{1}{2}$.

Norm. $h_{0,0} = -\frac{4421}{54} + \frac{1}{12}x_1 + \frac{1}{36}x_2$, $h_{0,1} = \frac{235099}{864} - \frac{3}{8}x_1 - \frac{1}{12}x_2$, $h_{0,2} = -\frac{53225}{216} + \frac{1}{4}x_1 + \frac{1}{12}x_2$, $h_{0,3} = \frac{16265}{288} + \frac{1}{24}x_1 - \frac{1}{36}x_2$;

$h_{1,0} = \frac{235099}{864} - \frac{3}{8}x_1 - \frac{1}{12}x_2$, $h_{1,1} = -\frac{8485}{9} + \frac{7}{4}x_1 + \frac{1}{4}x_2$, $h_{1,2} = \frac{88159}{96} - \frac{13}{8}x_1 - \frac{1}{4}x_2$, $h_{1,3} = -\frac{53191}{216} + \frac{1}{4}x_1 + \frac{1}{12}x_2$;

$h_{2,0} = -\frac{53225}{216} + \frac{1}{4}x_1 + \frac{1}{12}x_2$, $h_{2,1} = \frac{88159}{96} - \frac{13}{8}x_1 - \frac{1}{4}x_2$, $h_{2,2} = -\frac{33953}{36} + \frac{7}{4}x_1 + \frac{1}{4}x_2$, $h_{2,3} = \frac{234971}{864} - \frac{3}{8}x_1 - \frac{1}{12}x_2$;

$h_{3,0} = \frac{16265}{288} + \frac{1}{24}x_1 - \frac{1}{36}x_2$, $h_{3,1} = -\frac{53191}{216} + \frac{1}{4}x_1 + \frac{1}{12}x_2$, $h_{3,2} = \frac{234971}{864} - \frac{3}{8}x_1 - \frac{1}{12}x_2$, $h_{3,3} = -\frac{8759}{108} + \frac{1}{12}x_1 + \frac{1}{36}x_2$.

Boundary operator.

$fq_{0,0} = -(33048723219840x_1 + 85705756992x_2 - 422039808x_1x_2 - 154648248192x_1^2 + 248583168x_1^3 - 2415569939181083)/6$

$fq_{0,1} = (-127236096x_1x_2 - 631668895465685 + 8529226885512x_1 + 25889227344x_2 - 39533821632x_1^2 + 63265536x_1^3)$

$fq_{0,2} = -3(-57376512x_1x_2 - 175395768352673 + 2221076753280x_1 + 11730429216x_2 - 9796298112x_1^2 + 15303168x_1^3)/2$

$fq_{0,3} = -(111082752x_1x_2 + 38379992551957 + 169611442992x_1 - 22853387520x_2 - 3113385984x_1^2 + 6718464x_1^3)/3$

$fq_{0,4} = 8(1114560x_1x_2 - 3126314790419 - 230699286x_2 + 53751114648x_1 + 489888x_1^3 - 288112464x_1^2)$

$fq_{0,5} = -6(176256x_1x_2 - 630683707861 - 36417960x_2 + 10493831412x_1 + 93312x_1^3 - 55295136x_1^2)$

$fq_{1,0} = (-2467535470632x_1 - 7560127440x_2 + 37366272x_1x_2 + 11487655872x_1^2 - 18475776x_1^3 + 182013256354915)/3$

$fq_{1,1} = (31746816x_1x_2 + 349348112239985 - 4934674023552x_1 - 6484065984x_2 + 23561995392x_1^2 - 38071296x_1^3)/2$

$fq_{1,2} = (-32866560x_1x_2 - 330480117255149 + 4654174009968x_1 + 6702105024x_2 - 22182937344x_1^2 + 35831808x_1^3)$

$fq_{1,3} = (24862464x_1x_2 + 744757008535873 - 10734014433024x_1 - 5213420064x_2 + 51861845376x_1^2 - 83980800x_1^3)/6$

$fq_{1,4} = (352512x_1x_2 - 33880491134138 - 61652880x_2 + 495784036872x_1 + 3919104x_1^3 - 2416267584x_1^2)$

$fq_{1,5} = (41472x_1x_2 + 4889298728344 - 9473328x_2 - 71106806016x_1 - 559872x_1^3 + 345347712x_1^2)$

$fq_{2,0} = (4468035121344x_1 + 11761519680x_2 - 57708288x_1x_2 - 20898684288x_1^2 + 33592320x_1^3 - 326921290536385)/6$

$fq_{2,1} = (51549696x_1x_2 + 298516556982667 - 4087820882232x_1 - 10494673200x_2 + 19145771712x_1^2 - 30792960x_1^3)$

$fq_{2,2} = (-88729344x_1x_2 - 473456622343649 + 6436561653504x_1 + 18126860832x_2 - 29998843776x_1^2 + 48148992x_1^3)/2$

$fq_{2,3} = (28512000x_1x_2 + 79236815169347 - 985971485232x_1 - 5935521024x_2 + 4306473216x_1^2 - 6718464x_1^3)/3$

$fq_{2,4} = (-8211456x_1x_2 - 39288284063904 + 1682660304x_2 + 528786861984x_1 + 3919104x_1^3 - 2447947008x_1^2)$

$fq_{2,5} = (1140480x_1x_2 + 5574648286010 - 233163792x_2 - 75262164984x_1 - 559872x_1^2 + 349220160x_1^3)$

$fq_{3,0} = (-2373828501672x_1 - 2530216656x_2 + 12358656x_1x_2 + 11400595776x_1^2 - 18475776x_1^3 + 166688181811999)/3$

$fq_{3,1} = (-40414464x_1x_2 - 506993048493859 + 7205150657088x_1 + 8262613440x_2 - 34559125632x_1^2 + 55987200x_1^3)/2$

$fq_{3,2} = (51570432x_1x_2 + 533353205953455 - 7531408278288x_1 - 10511667648x_2 + 35980360704x_1^2 - 58226688x_1^3)$

$fq_{3,3} = (-125680896x_1x_2 - 1783398165357899 + 25434754033536x_1 + 25733909088x_2 - 122263572096x_1^2 + 198194688x_1^3)/6$

$$fq_{3,4} = (-16775424x_1x_2 - 44437957089478 + 3387345264x_2 + 560447103240x_1 + 3919104x_1^3 - 2477013696x_1^2)$$
$$fq_{3,5} = (2239488x_1x_2 + 6251575671936 - 452563632x_2 - 79463658528x_1 - 559872x_1^3 + 353279232x_1^2),$$

where the denominator is given by $f = (-35811072x_1x_2 - 2249831450017201 + 7261488864x_2 + 3101605909056x_1 + 23514624x_1^3 - 14592286080x_1^2)$.

## 8.2. Restricted Full Norms

The boundary part of the difference operators, $Q = [H^{-1}BH^{-1}C]$, has the size $((\tau+2) \times [(\tau+2) + (\tau+1)/2])$ and are accurate to order $\tau$. The antisymmetric interior stencil uses $\tau + 2$ symmetric centered points and is accurate to order $\tau + 1$. The case of third-order accuracy at the boundary is given below.

*Third-order accuracy at the boundary.*

Interior stencil. $\alpha_1 = \frac{2}{3}$, $\alpha_2 = -\frac{1}{12}$.

Norm. $h_{0,0} = x_1$, $h_{0,1} = 0$, $h_{0,2} = 0$, $h_{0,3} = 0$, $h_{0,4} = 0$;

$h_{1,0} = 0$, $h_{1,1} = -\frac{3203}{9} - 16x_1 + \frac{1}{12}x_2 + \frac{1}{36}x_3$, $h_{1,2} = 1000129/864 + 24x_1 - \frac{3}{8}x_2 - \frac{1}{12}x_3$, $h_{1,3} = -233699/216 - 16x_1 + \frac{1}{4}x_2 + \frac{1}{12}x_3$, $h_{1,4} = 244393/864 + 4x_1 + \frac{1}{24}x_2 - \frac{1}{36}x_3$;

$h_{2,0} = 0$, $h_{2,1} = 1000129/864 + 24x_1 - \frac{3}{8}x_2 - \frac{1}{12}x_3$, $h_{2,2} = -34867/9 - 36x_1 + \frac{7}{4}x_2 + \frac{1}{4}x_3$, $h_{2,3} = 364837/96 + 24x_1 - \frac{13}{8}x_2 - \frac{1}{4}x_3$, $h_{2,4} = -234337/216 - 6x_1 + \frac{1}{4}x_2 + \frac{1}{12}x_3$;

$h_{3,0} = 0$, $h_{3,1} = -233699/216 - 16x_1 + \frac{1}{4}x_2 + \frac{1}{12}x_3$, $h_{3,2} = 364837/96 + 24x_1 - \frac{13}{8}x_2 - \frac{1}{4}x_3$, $h_{3,3} = -139673/36 - 16x_1 + \frac{7}{4}x_2 + \frac{1}{4}x_3$, $h_{3,4} = 1005377/864 + 4x_1 - \frac{3}{8}x_2 - \frac{1}{12}x_3$;

$h_{4,0} = 0$, $h_{4,1} = 244393/864 + 4x_1 + \frac{1}{24}x_2 - \frac{1}{36}x_3$, $h_{4,2} = -234337/216 - 6x_1 + \frac{1}{4}x_2 + \frac{1}{12}x_3$, $h_{4,3} = 1005377/864 + 4x_1 - \frac{3}{8}x_2 - \frac{1}{12}x_3$, $h_{4,4} = -38977/108 - x_1 + \frac{1}{12}x_2 + \frac{1}{36}x_3$.

Boundary operator. $q_{0,0} = -\frac{1}{2}/x_1$, $q_{0,1} = -\frac{1}{3}(-6 + 13x_1)/x_1$, $q_{0,2} = \frac{1}{2}(-6 + 19x_1)/x_1$, $q_{0,3} = -(-2 + 7x_1)/x_1$, $q_{0,4} = \frac{1}{6}(-3 + 11x_1)/x_1$, $q_{0,5} = 0$, $q_{0,6} = 0$;

$$fq_{1,0} = \left\{ \begin{array}{l} 12(1389312x_1x_2x_3 + 23561475245750x_1 + 382288508604x_2 + 1529781768x_3 - \\ 60141543984x_1x_2 - 863009856x_1x_3 - 546271776x_2^2 - 2477952x_2x_3 + \\ 35754048x_1x_2^2 + 279936x_3^3 + 248632x_1^2x_2x_3 + \\ 373248x_1^2x_2^2 - 2468423808x_1^3x_2 + 697073039520x_1^2 - \\ 93675178378625 - 92472192x_1^2x_3) \end{array} \right.$$

$$fq_{1,1} = \left\{ \begin{array}{l} -(-565843968x_1x_2x_3 - 20937429673405920x_1 - 288832087875456x_2 - \\ 245865521664x_3 + 90463714786944x_1x_2 + 334373780160x_1x_3 + 454278305664x_2^2 + \\ 400142592x_2x_3 - 1367823283320x_1x_2^2 + 73903104x_1x_3^2 - 241864704x_2^3 + \\ 71663616x_1^2x_2x_3 + 107495424x_1^2x_2^2 - 710906056704x_1^3x_2 + \\ 200757035381760x_1^2 - 26631991296x_1^2x_3 + 62069541085794875)/6 \end{array} \right.$$

$$fq_{1,2} = \left\{ \begin{array}{l} 3(-51591168x_1x_2x_3 - 1912449225706160x_1 - 25382367160632x_2 - 27354017520x_3 + \\ 82567759034888x_1x_2 + 303629281920x_1x_3 + 39682872000x_2^2 + 44499456x_2x_3 - \\ 12460386816x_1x_2^2 + 6718464x_1x_3^2 - 21088512x_2^3 + 5971968x_1^2x_2x_3 + \\ 8957952x_1^2x_2^2 - 59242171392x_1^3x_2 + 16729752948480x_1^2 - \\ 2219332608x_1^2x_3 + 5506314313768875) \end{array} \right.$$

$$fq_{1,3} = \left\{ \begin{array}{l} -(-88584192x_1x_2x_3 - 4692072170000640x_1 - 47292926246784x_2 - 63583974432x_3 + \\ 10338701506491875 + 22212732746880x_1x_2 + 48561073344x_1x_3 + 73659466368x_2^2 + \\ 103783680x_2x_3 - 359031838464x_1x_2^2 + 20155392x_1x_3^2 - 39191040x_2^3 + \\ 23887872x_1^2x_2x_3 + 35831808x_1^2x_2^2 - 236968685568x_1^3x_2 + \\ 66919011793920x_1^2 - 8877330432x_1^2x_3)/2 \end{array} \right.$$

$$fq_{1,4} = \left\{ \begin{array}{l} (17169408x_1x_2x_3 - 875467742864040x_1 + 10802001206304x_2 - 36381362976x_3 + \\ 5784648613056x_1x_2 - 14483945664x_1x_3 - 18615401856x_2^2 + 59450112x_2x_3 - \\ 11106180864x_1x_2^2 + 6718464x_1x_3^2 + 10077696x_2^3 + 8957952x_1^2x_2x_3 + \\ 13436928x_1^2x_2^2 - 888632576088x_1^3x_2 + 25004629422720x_1^2 - \\ 1948629566030375 - 3328998912x_1^2x_3)/3 \end{array} \right.$$

$$fq_{1,5} = \left\{ \begin{array}{l} -24(5217696x_1x_2^2 - 9941107104x_1x_2 + 300672x_1x_2x_3 + 4288182479270x_1 - \\ 198051264x_1x_3 + 163296x_2^3 - 304131024x_2^2 + 233280x_2x_3 + 186213363048x_2 - \\ 142938162x_3 - 37481321648875) \end{array} \right.$$

$$fq_{1,6} = \left\{ \begin{array}{l} 2(8211456x_1x_2^2 - 15889893120x_1x_2 + 497664x_1x_2x_3 - 324207360x_1x_3 + \\ 6870618696000x_1 + 279936x_2^3 - 522290592x_2^2 + 321474293532x_2 + \\ 300672x_2x_3 - 65287532057125 - 184148856x_3) \end{array} \right.$$

$$fq_{2,0} = \left\{ \begin{array}{l} 48(67392x_1x_2x_3 + 3051548096729x_1 + 90042984288x_2 + 193899366x_3 - \\ 8380331496x_1x_2 - 56223936x_1x_3 - 134894160x_2^2 - 305856x_2x_3 + \\ 5999184x_1x_2^2 - 20700319137235 + 69984x_2^3 + 124416x_1^2x_2x_3 + \\ 746496x_1^2x_2^2 - 2111349888x_1^3x_2 + 815128616928x_1^2 - 58672512x_1^2x_3) \end{array} \right.$$

$$fq_{2,1} = \left\{ \begin{array}{l} -(-313352832x_1x_2x_3 - 546425633238912x_1 + 20797530562008x_2 + 57732323184x_3 + \\ 4487651144064x_1x_2 + 9351671616x_1x_3 - 30003618240x_2^2 - 88169472x_2x_3 - \\ 9688025088x_1x_2^2 + 6718464x_1x_3^3 + 15116544x_2^3 - 4998147438898793 + \\ 71663616x_1^2x_2x_3 + 429981696x_1^2x_2^2 - 1216137535488x_1^3x_2 + \\ 469514083350528x_1^2 - 33795366912x_1^2x_3)/3 \end{array} \right.$$

$$fq_{2,2} = \left\{ \begin{array}{l} 3(61378560x_1x_2x_3 + 1665736866755776x_1 + 35723653878528x_2 + 64812507264x_3 - \\ 5414441860224x_1x_2 - 40742371008x_1x_3 - 54250217856x_2^2 - 103037184x_2x_3 + \\ 5866463232x_1x_2^2 - 2239488x_1x_3^3 + 28366848x_2^3 - 8065911877683451 + \\ 23887872x_1^2x_2x_3 + 143327232x_1^2x_2^2 - 405379178496x_1^3x_2 + \\ 156504694450176x_1^2 - 11265122304x_1^2x_3)/2 \end{array} \right.$$

$$fq_{2,3} = \left\{ \begin{array}{l} -(149299200x_1x_2x_3 + 3710751586940736x_2 + 69818790100464x_2 + 116786952576x_3 - \\ 12766585571712x_1x_2 - 94861606464x_1x_3 - 106676020224x_2^2 - 186727680x_2x_3 + \\ 15170291712x_1x_2^2 - 6718464x_1x_3^3 + 55987200x_2^3 - 15636088717639621 + \\ 23887872x_1^2x_2x_3 + 143327232x_1^2x_2^2 - 405379178496x_1^3x_2 + \\ 156504694450176x_1^2 - 11265122304x_1^2x_3) \end{array} \right.$$

$$fq_{2,4} = \left\{ \begin{array}{l} (335923200x_1x_2x_3 + 620778328107062 4x_1 + 137954728520064x_2 + 207352150560x_3 - \\ 18509803454592x_1x_2 - 215085409344x_1x_3 - 211342504320x_2^2 - 328935168x_2x_3 + \\ 18069308928x_1x_2^2 - 6718464x_1x_2^3 + 110854656x_2^3 - 30749882535748609 + \\ 35831808x_1^2x_2x_3 + 214990848x_1^2x_2^2 - 608068767744x_1^2x_2 + \\ 234757041675264x_1^2 - 16897683456x_1^2x_3)/6 \end{array} \right.$$

$$fq_{2,5} = \left\{ \begin{array}{l} -2(91072512x_1x_2^2 - 195196760064x_1x_2x_3 + 7216128x_1x_2x_3 + 87271245005664x_1 - \\ 4586212224x_1x_2 + 1959552x_2^3 - 3730777056x_2^2 - 4655232x_2x_3 + \\ 2419457482692x_2 - 533783285704223 + 2958900408x_3) \end{array} \right.$$

$$fq_{2,6} = \left\{ \begin{array}{l} 24(995328x_1x_2^2 - 2175393024x_1x_2 + 82944x_1x_2x_3 - 52372224x_1x_3 + \\ 976751816000x_1 + 23328x_2^3 - 44365968x_2^2 + 28759018176x_2 - 57024x_2x_3 - \\ 6343369900435 + 36009954x_3) \end{array} \right.$$

$$fq_{3,0} = \left\{ \begin{array}{l} 12(-850176x_1x_2x_3 - 8327527149946x_1 + 330550181532x_2 - 5439096x_3 + \\ 19074174192x_1x_2 + 480530496x_1x_3 - 527096160x_2^2 + 31104x_2x_3 - 7915968x_1x_2^2 + \\ 279936x_2^3 + 746496x_1^2x_2x_3 + 10077696x_1^2x_2^2 - 20469013632x_1^2x_2 + \\ 8408569556448x_1^2 - 69197124803235 - 426653568x_1^2x_3) \end{array} \right.$$

$$fq_{3,1} = \left\{ \begin{array}{l} -(-97044480x_1x_2x_3 + 802170337139424x_1 + 14567309911699 2x_4 + 100322381952x_3 - \\ 5837919765120x_1x_2 + 50217900480x_1x_3 - 228511539072x_2^2 - 155748096x_2x_3 + \\ 11785305600x_1x_2^2 - 6718464x_1x_2^3 + 120932352x_2^3 + 214990848x_1^2x_2x_3 + \\ 2902376448x_1^2x_2^2 - 5895075926016x_1^2x_2 + 2421668032257024x_1^2 - \\ 31341180273033409 - 122876227584x_1^2x_3)/6 \end{array} \right.$$

$$fq_{3,2} = \left\{ \begin{array}{l} 3(23390208x_1x_2x_3 + 805097807026512x_1 + 23535330061320x_2 + 30016883088x_3 - \\ 3179822330880x_1x_2 - 14929175232x_1x_3 - 36421265088x_2^2 - 48024576x_2x_3 + \\ 4435305984x_1x_2^2 - 2239488x_1x_2^3 + 19222272x_2^3 + 17915904x_1^2x_2x_3 + \\ 241864704x_1^2x_2^2 - 491256327168x_1^2x_2 + 201805669354752x_1^2 - \\ 10239685632x_1^2x_3 - 5175484656324797) \end{array} \right.$$

$$fq_{3,3} = \left\{ \begin{array}{l} -3(49434624x_1x_2x_3 + 1104761834978176x_1 + 33721099789824x_2 + 52828971552x_3 - \\ 3849768369792x_1x_2 - 31177517760x_1x_3 - 51767185536x_2^2 - 84540672x_2x_3 + \\ 4739254272x_1x_2^2 - 2239488x_1x_2^3 + 27247104x_2^3 + 23687872x_1^2x_2x_3 + \\ 322486272x_1^2x_2^2 - 655008436224x_1^2x_2 + 269074225806336x_1^2 - \\ 13652914176x_1^2x_3 - 7506161222686091)/2 \end{array} \right.$$

$$fq_{3,4} = \left\{ \begin{array}{l} (50015232x_1x_2x_3 - 954978330702312x_1 + 13955741779104x_2 + 56850305760x_3 + \\ 6303459071808x_1x_2 - 32638479552x_1x_3 - 19804569984x_2^2 - 89849088x_2x_3 - \\ 11800422144x_1x_2^2 + 6718464x_1x_2^3 + 10077696x_2^3 + 26873856x_1^2x_2x_3 + \\ 362797056x_1^2x_2^2 - 736884490752x_1^2x_2 + 302708504032128x_1^2 - \\ 15359528448x_1^2x_3 - 3454300789008023)/3 \end{array} \right.$$

$$fq_{3,5} = \left\{ \begin{array}{l} -16(14912230420107x_1 - 33371320248x_1x_2 + 14731632x_1x_2^2 - 828598896x_1x_3 + \\ 1353024x_1x_2x_3 + 940161357x_3 - 74653422358664 + 244944x_2^2 - 1513728x_2x_3 - \\ 470490768x_2^2 + 317926705698x_2) \end{array} \right.$$

$$fq_{3,6} = \left\{ \begin{array}{l} 2(15676416x_1x_2^2 - 36078213888x_1x_2 + 1492992x_1x_2x_3 - 912777984x_1x_3 + \\ 16198255166400x_1 + 279936x_2^3 - 536707296x_2^2 + 361410270012x_2 - \\ 1669248x_2x_3 - 84486787162055 + 1035477000x_3) \end{array} \right.$$

$$fq_{4,0} = \left\{ \begin{array}{l} 48(-492480x_1x_2x_3 - 10390548152065x_1 + 69878810916x_2 - 203331978x_3 + \\ 27393015024x_1x_2 + 317817216x_1x_3 - 124384896x_2^2 + 321408x_2x_3 - \\ 17515440x_1x_2^2 + 69984x_2^3 + 248832x_1^2x_2x_3 + 5971968x_1^2x_2^2 - \\ 10269338112x_1^2x_2 + 4214870894400x_1^2 - 167090688x_1^2x_3 - \\ 12144968733650) \end{array} \right.$$

$$fq_{4,1} = \left\{ \begin{array}{l} -(-334430208x_1x_2x_3 - 8101800797826240x_1 + 9933326262360x_2 - 145646157456x_3 + \\ 24624376265088x_1x_2 + 214458725952x_1x_3 - 24181229376x_2^2 + 230957568x_2x_3 - \\ 23012978688x_1x_2^2 + 6718464x_1x_2^3 + 15116544x_2^3 + 143327232x_1^2x_2x_3 + \\ 3439853568x_1^2x_2^2 - 5915138752512x_1^2x_2 + 2427765635174400x_1^2 - \\ 96244236288x_1^2x_3 - 438610159220645)/3 \end{array} \right.$$

$$fq_{4,2} = \left\{ \begin{array}{l} 3(-146976768x_1x_2x_3 - 4130297290252480x_1 - 17167695218496x_2 - 69044258304x_3 + \\ 14149731860352x_1x_2 + 93377918016x_1x_3 + 23950235520x_2^2 + 109963008x_2x_3 - \\ 16308449280x_1x_2^2 + 6718464x_1x_2^3 - 11943936x_2^3 + 47775744x_1^2x_2x_3 + \\ 1146617856x_1^2x_2^2 - 1971712917504x_1^2x_2 + 809255211724800x_1^2 - \\ 32081412096x_1^2x_3 + 4293793001873985)/2 \end{array} \right.$$

$$fq_{4,3} = \left\{ \begin{array}{l} -(-269733888x_1x_2x_3 - 8599520736070080x_1 - 81467392389264x_2 - 144422472960x_3 + \\ 32400894519936x_1x_2 + 168566526144x_1x_3 + 123996904704x_2^2 + 231641856x_2x_3 - \\ 42423367680x_1x_2^2 + 20155392x_1x_2^3 - 64945152x_2^3 + 47775744x_1^2x_2x_3 + \\ 1146617856x_1^2x_2^2 - 1971712917504x_1^2x_2 + 809255211724800x_1^2 - \\ 32081412096x_1^2x_3 + 18332637080024935) \end{array} \right.$$

$$fq_{4,4} = \left\{ \begin{array}{l} (-673339392x_1x_2x_3 - 25986204959404320x_1 - 306855074673792x_2 - 353717343072x_3 + \\ 104545776200832x_1x_2 + 420457459392x_1x_3 + 476625596544x_2^2 + 566611200x_2x_3 - \\ 146927208960x_1x_2^2 + 73903104x_1x_2^3 - 251942400x_2^3 + \\ 71663616x_1^2x_2x_3 + 1719926784x_1^2x_2^2 - 2957569376256x_1^2x_2 + \\ 1213882817587200x_1^2 - 48122118144x_1^2x_3 + 67107296980771115)/6 \end{array} \right.$$

$$fq_{4,5} = \left\{ \begin{array}{l} -6(58475520x_1x_2^2 - 124624942848x_1x_2 + 48107052x_1x_2x_3 - 2834784000x_1x_3 + \\ 53132220804800x_1 + 653184x_2^3 - 1272036960x_2^2 + 896526448236x_2 - \\ 6521472x_2x_3 - 221957121908115 + 3990637800x_3) \end{array} \right.$$

$$fq_{4,6} = \left\{ \begin{array}{l} 16(2743670587200x_1 - 6405516288x_1x_2 + 2985984x_1x_2^2 - 147142656x_1x_3 + \\ 248832x_1x_2x_3 + 203240475x_3 - 11724346138850 + 34992x_2^3 - 331776x_2x_3 - \\ 68012784x_2^2 + 47671712730x_2) \end{array} \right.$$

where $f = (6718464x_1x_2^3 - 13403335680x_1x_2^2 + 9859389455232x_1x_2 - 89579520x_1x_2x_3 - 2546440283980800x_1 + 54170731584x_1x_3 - 30233088x_2^3 + 57302681472x_2^2 + 99512064x_2x_3 - 37301376049344x_2 + 8301045567048625 - 61336247328x_3)$.

## 8.3. Minimum Bandwidth Operators

The parameters can be used to minimize the bandwidth of the difference operators. In the case of full norms we have two parameters, and we can choose $q_{0,5} = q_{1,5} = 0$ to minimize the bandwidth. This gives a non-linear system of equations with two solutions, one corresponding to an indefinite norm. The other solution gives the positive definite norm

$$fh_{0,0} = -(-263779327 + 582347\sqrt{7}\sqrt{24943})/20736$$

$$fh_{0,1} = (-85294475 + 209367\sqrt{7}\sqrt{24943})/2304$$

$$fh_{0,2} = -7(-39830729 + 96637\sqrt{7}\sqrt{24943})/6912$$

$$fh_{0,3} = 7(-49913119 + 118955\sqrt{7}\sqrt{24943})/20736$$

$$fh_{1,0} = (-85294475 + 209367\sqrt{7}\sqrt{24943})/2304$$

$$fh_{1,1} = -(-68373221 + 156329\sqrt{7}\sqrt{24943})/768$$

$$fh_{1,2} = (-73504403 + 177871\sqrt{7}\sqrt{24943})/768$$

$$fh_{1,3} = -(-91953365 + 219321\sqrt{7}\sqrt{24943})/2304$$

$$fh_{2,0} = -7(-39830729 + 96637\sqrt{7}\sqrt{24943})/6912$$

$$fh_{2,1} = (-73504403 + 177871\sqrt{7}\sqrt{24943})/768$$

$$fh_{2,2} = -(-207379375 + 483131\sqrt{7}\sqrt{24943})/2304$$

$$fh_{2,3} = (-253102241 + 610693\sqrt{7}\sqrt{24943})/6912$$

$$fh_{3,0} = 7(-49913119 + 118955\sqrt{7}\sqrt{24943})/20736$$

$$fh_{3,1} = -(-91953365 + 219321\sqrt{7}\sqrt{24943})/2304$$

$$fh_{3,2} = (-253102241 + 610693\sqrt{7}\sqrt{24943})/6912$$

$$fh_{3,3} = -(-220497151 + 311435\sqrt{7}\sqrt{24943})/20736,$$

where $f = -2716 + 17\sqrt{7}\sqrt{24943}$. In decimal form

$$H = \begin{pmatrix} 0.2247105 & 0.2166550 & -0.1267943 & -0.01596010 \\ 0.2166550 & 0.9053611 & 0.2432041 & 0.03061306 \\ -0.1267943 & 0.2432041 & 0.5442500 & 0.06850688 \\ -0.01596010 & 0.03061306 & 0.06850688 & 0.9932290 \end{pmatrix}.$$

The eigenvalues of $H$ are

$$1.0837, \quad 0.9690, \quad 0.5382, \quad 0.0766.$$

To this positive definite matrix the corresponding boundary operator with minimized bandwidth is

$$f_1 q_{0,0} = -(-44534043 + 105337\sqrt{7}\sqrt{24943})/6$$

$$f_1 q_{0,1} = (-12112595 + 28577\sqrt{7}\sqrt{24943})$$

$$f_1 q_{0,2} = -3(-4000241 + 9355\sqrt{7}\sqrt{24943})/2$$

$$f_1 q_{0,3} = (-3888369 + 8843\sqrt{7}\sqrt{24943})/3$$

$$f_1 q_{0,4} = 32(-437 + 2\sqrt{7}\sqrt{24943})$$

$$q_{0,5} = 0;$$

$$f_2 q_{1,0} = -(-2322245967111 + 5556025261\sqrt{7}\sqrt{24943})/3$$

$$f_2 q_{1,1} = -7(-752572077947 + 1798259497\sqrt{7}\sqrt{24943})/6$$

$$f_2 q_{1,2} = (-1858955051919 + 4443227269\sqrt{7}\sqrt{24943})$$

$$f_2 q_{1,3} = -5(-186474644307 + 443526257\sqrt{7}\sqrt{24943})/6$$

$$f_2 q_{1,4} = -8(-19303788133 + 46366583\sqrt{7}\sqrt{24943})/3$$

$$q_{1,5} = 0;$$

$$f_2 q_{2,0} = (-45425453907169 + 108756928139\sqrt{7}\sqrt{24943})/90$$

$$f_2 q_{2,1} = -(-125808978873263 + 301133010757\sqrt{7}\sqrt{24943})/45$$

$$f_2 q_{2,2} = (-21614189741851 + 51772725129\sqrt{7}\sqrt{24943})/10$$

$$f_2 q_{2,3} = -7(-1134888541649 + 2748754651\sqrt{7}\sqrt{24943})/45$$

$$f_2 q_{2,4} = 8(-2027181750139 + 4857666461\sqrt{7}\sqrt{24943})/45$$

$$f_3 q_{2,5} = -2(-98399 + 247\sqrt{7}\sqrt{24943})/15;$$

$$f_3 q_{3,0} = -(-931273 + 2263\sqrt{7}\sqrt{24943})/15$$

$$f_3 q_{3,1} = (-9235079 + 22021\sqrt{7}\sqrt{24943})/30$$

$$f_3 q_{3,2} = -129(-30989 + 71\sqrt{7}\sqrt{24943})/5$$

$$f_3 q_{3,3} = (-9653609 + 23411\sqrt{7}\sqrt{24943})/30$$

$$f_3 q_{3,4} = 2(-2002441 + 4049\sqrt{7}\sqrt{24943})/15$$

$$f_3 q_{3,5} = -2(-87535 + 179\sqrt{7}\sqrt{24943})/5,$$

where

$$f_1 = -4056177 + 9611\sqrt{7}\sqrt{24943}$$

$$f_2 = -2167815662047 + 5185092597\sqrt{7}\sqrt{24943}$$

$$f_3 = -567063 + 1213\sqrt{7}\sqrt{24943}.$$

Since $q_{0,4} \neq 0$, the difference stencil corresponding to the first point will need four neighbors to the right. The rest of the points will need three or less neighbors to the left and/or to the right. When using parallel computers we have to store the main diagonal and four superdiagonals. Thus we have the same number of diagonals as for the operators corresponding to the diagonal norms.

In the case of restricted full norms we choose the parameters such that $q_{0,4}$, $q_{1,6}$, $q_{2,6}$ are zeroed. As for the full norm case this leads to a non-linear system of equations with two solutions, with one root corresponding to a positive definite norm and one corresponding to an indefinite norm. The first solution then gives the following positive definite norm

$$h_{0,0} = \tfrac{3}{11}, \quad h_{0,1} = 0, \quad h_{0,2} = 0, \quad h_{0,3} = 0, \quad h_{0,4} = 0;$$

$$h_{1,0} = 0$$

$$f h_{1,1} = (299913292801 + 56278767\sqrt{26116897})/228096$$

$$f h_{1,2} = -(64756272879 + 310129\sqrt{26116897})/76032$$

$$f h_{1,3} = -(-50615837729 + 5284177\sqrt{26116897})/76032$$

$$f h_{1,4} = (-5026701941 + 948741\sqrt{26116897})/20736;$$

$$h_{2,0} = 0$$

$$f h_{2,1} = -(64756272879 + 310129\sqrt{26116897})/76032$$

$$f h_{2,2} = -7(-6989673895 + 13527\sqrt{26116897})/25344$$

$$f h_{2,3} = 49(-657605303 + 100423\sqrt{26116897})/25344$$

$$f h_{2,4} = -49(-75022899 + 14467\sqrt{26116897})/6912;$$

$$h_{3,0} = 0$$

$$fh_{3,1} = -(-50615837729 + 5284177\sqrt{26116897})/76032$$

$$fh_{3,2} = 49(-657605303 + 100423\sqrt{26116897})/25344$$

$$fh_{3,3} = -(-45333081425 + 982369\sqrt{26116897})/25344$$

$$fh_{3,4} = (-3355209517 + 597005\sqrt{26116897})/6912;$$

$$h_{4,0} = 0$$

$$fh_{4,1} = (-5026701941 + 948741\sqrt{26116897})/20736$$

$$fh_{4,2} = -49(-75022899 + 14467\sqrt{26116897})/6912$$

$$fh_{4,3} = (-3355209517 + 597005\sqrt{26116897})/6912$$

$$fh_{4,4} = 5(35213725709 + 5139171\sqrt{26116897})/228096,$$

where $f = 591223 + 146\sqrt{26116897}$. In decimal form

$$H = \begin{pmatrix} 0.2727273 & 0 & 0 & 0 & 0 \\ 0 & 1.926028 & -0.6524409 & 0.2322075 & -0.006425591 \\ 0 & -0.6524409 & 1.429281 & -0.2087529 & 0.005776559 \\ 0 & 0.2322075 & -0.2087529 & 1.189382 & -0.03291229 \\ 0 & -0.006425591 & 0.005776559 & -0.03291229 & 1.007677 \end{pmatrix}.$$

The eigenvalues of $H$ are

$$0.2727, \quad 2.4520, \quad 1.1308, \quad 1.0019, \quad 0.9676.$$

To this positive definite matrix the corresponding boundary operator with minimized bandwidth is

$$q_{0,0} = -\tfrac{11}{6}, \quad q_{0,1} = 3, \quad q_{0,2} = -\tfrac{3}{2}, \quad q_{0,3} = \tfrac{1}{3}$$

$$q_{0,4} = 0, \quad q_{0,5} = 0, \quad q_{0,6} = 0;$$

$$f_1 q_{1,0} = -24(-779042810827742869 + 104535124033147\sqrt{26116897})$$

$$f_1 q_{1,1} = -(-176530817412806109689 + 29768274816875927\sqrt{26116897})/6$$

$$f_1 q_{1,2} = 343(-1710791116122226871 + 27975630462649\sqrt{26116897})$$

$$f_1 q_{1,3} = -3(-7475554291248533227 + 1648464218793925\sqrt{26116897})/2$$

$$f_1 q_{1,4} = (-2383792768180030915 + 1179620587812973\sqrt{26116897})/3$$

$$f_1 q_{1,5} = -1232(-115724529581315 + 37280576429\sqrt{26116897})$$

$$q_{1,6} = 0;$$

$$f_2 q_{2,0} = -12(-380966843 + 86315\sqrt{26116897})$$

$$f_2 q_{2,1} = (5024933015 + 2010631\sqrt{26116897})/3$$

$$f_2 q_{2,2} = -231(-431968921 + 86711\sqrt{26116897})/2$$

$$f_2 q_{2,3} = (-65931742559 + 12256337\sqrt{26116897})$$

$$f_2 q_{2,4} = -(-50597298167 + 9716873\sqrt{26116897})/6$$

$$f_2 q_{2,5} = -88(-15453061 + 2911\sqrt{26116897})$$

$$q_{2,6} = 0;$$

$$f_1 q_{3,0} = 48(-56020909845192541 + 9790180507043 \sqrt{26116897})$$

$$f_1 q_{3,1} = (-99182490049237586011 + 1463702013196501 \sqrt{26116897})/6$$

$$f_1 q_{3,2} = -13(-4130451756851441723 + 664278707201077 \sqrt{26116897})$$

$$f_1 q_{3,3} = 3(-26937108467782666617 + 5169063172799767 \sqrt{26116897})/2$$

$$f_1 q_{3,4} = -(6548308508012371315 + 3968886380989379 \sqrt{26116897})/3$$

$$f_1 q_{3,5} = 88(-91337851897923397 + 19696768305507 \sqrt{26116897})$$

$$f_3 q_{3,6} = 242(-120683 + 15 \sqrt{26116897});$$

$$f_3 q_{4,0} = 264(-120683 + 15 \sqrt{26116897})$$

$$f_3 q_{4,1} = (-43118111 + 23357 \sqrt{26116897})/3$$

$$f_3 q_{4,2} = -47(-28770085 + 2259 \sqrt{26116897})/2$$

$$f_3 q_{4,3} = -3(1003619433 + 11777 \sqrt{26116897})$$

$$f_3 q_{4,4} = -11(-384168269 + 65747 \sqrt{26116897})/6$$

$$f_3 q_{4,5} = 22(87290207 + 10221 \sqrt{26116897})$$

$$f_3 q_{4,6} = -66(3692405 + 419 \sqrt{26116897}),$$

where

$$f_1 = -56764003702447356523 + 8154993476273221 \sqrt{26116897}$$

$$f_2 = -55804550303 + 9650225 \sqrt{26116897}$$

$$f_3 = 3262210757 + 271861 \sqrt{26116897}.$$

## REFERENCES

1. H.-O. Kreiss and G. Scherer, "Finite Element and Finite Difference Methods for Hyperbolic Partial Differential Equations," in *Mathematical Aspects of Finite Elements in Partial Differential Equations* (Academic Press, New York/London, 1974).
2. P. Olsson, Ph.D. thesis, Department of Scientific Computing, Uppsala University, 1992.
3. I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products* (Academic Press, New York/London, 1965), pp. 7, 1076.
4. G. Scherer, Ph.D. thesis, Department of Scientific Computing, Uppsala University, 1977.