

# Uniform Expansions for a Class of Finite Difference Schemes for Elliptic Boundary Value Problems

By Harry Munz

**Abstract.** For a class of finite difference schemes for the Dirichlet problem on a bounded region  $\Omega \subset \mathbf{R}^n$ , the existence of uniform expansions of the approximate solution for meshlength  $h \rightarrow 0$  is shown. The results also improve error bounds which Pereyra, Proskurowski, and Widlund obtained with respect to certain discrete  $L_2$ -norms.

**1. Introduction.** In [7], Pereyra, Proskurowski, and Widlund discuss a class of finite difference schemes, due to H.-O. Kreiss, for the Poisson-equation  $\Delta u = f$  in an arbitrary bounded region  $\Omega \subset \mathbf{R}^n$  with Dirichlet boundary conditions.

On a uniform mesh, they replace the second-order partial derivatives by the standard three-point finite difference approximations. In points near the boundary  $\Gamma$  of  $\Omega$ , it may happen that one or more of the points needed in this approximation lie outside  $\Omega$ . For those points, provisional values are calculated by one-dimensional polynomial extrapolation of fixed degree  $k$  along the corresponding mesh line, thereby using the boundary value at the intersection of the mesh line and  $\Gamma$ .

For  $k \leq 6$ , Pereyra et al. [7] could show the stability and the convergence of these schemes and the existence of asymptotic expansions for  $h \rightarrow 0$  of the finite difference solution  $U$  with respect to certain discrete  $L_2$ -norms, which allow the use of Richardson extrapolation or deferred correction methods. For a sufficiently smooth boundary  $\Gamma$ , their expansions have the form

$$(1) \quad U = \Delta_h(u + h^2 e^{(1)} + h^4 e^{(2)}) + O(h^{k-0.5}),$$

where  $\Delta_h$  is the restriction operator to meshpoints in  $\Omega$ , and  $e^{(1)}, e^{(2)}$  are certain continuous functions on  $\bar{\Omega}$ , independent of  $h$ .

Pereyra et al. [7] conjecture the existence of similar expansions in the discrete maximum norm.

In this paper, the schemes of Pereyra et al. are applied to general linear second-order elliptic equations without mixed derivatives with Dirichlet boundary conditions.

For  $k \leq 2$  and a sufficiently fine mesh, the finite difference operators obtained are of inverse monotone type. Therefore, the classical convergence proof works [2]. For  $k \leq 4$ , an idea of Bramble and Hubbard [1] can be used to show the convergence of the schemes for the generalized problem and the existence of asymptotic expansions of the approximate solution with respect to both the discrete maximum norm and the discrete  $L_2$ -norms of Pereyra et al. [7]; the difference operator is modified near the boundary such that it becomes inverse monotone,

---

Received December 5, 1979; revised March 4, 1980.  
1980 *Mathematics Subject Classification*. Primary 65B05, 65N15.

and the points where a modification is necessary are discussed separately. The expansions obtained have the form

$$(2) \quad U = \Delta_h(u + h^2 e^{(1)} + h^4 e^{(2)}) + O(h^{k+1})$$

in both norms.

Finally, we report on numerical tests in which we exploited the asymptotic expansions by a modified deferred correction method. Unfortunately, the theorems of Pereyra [6], on the gain in accuracy obtained by using deferred correction methods, do not fit the present case.

**2. The Difference Operator.** For  $h > 0$ , let  $\mathbf{R}_h^n$  denote a uniform mesh of mesh size  $h$  on the  $\mathbf{R}^n$ . We assume that there are enough meshpoints on each mesh line in  $\Omega$  so that the extrapolation operations described in Section 1 and below are possible. Let  $\Omega_h := \mathbf{R}_h^n \cap \Omega$ . We denote by  $R_h$  the set of regular meshpoints, i.e. of those points  $x \in \Omega_h$  which have all their closest neighbors  $x \pm he_i$ ,  $i = 1, \dots, n$ , in  $\Omega$  ( $e_i$  is the unit vector parallel to the  $i$ th coordinate axis), and define  $\Gamma_h^x := \Omega_h \setminus R_h$ . The set of all  $x_i^\Gamma$  (see Figure 1), i.e. of all intersections of mesh lines meeting  $\Omega$  with the boundary  $\Gamma$  of  $\Omega$ , is called  $\Gamma_h$  and  $\bar{\Omega}_h := \Omega_h \cup \Gamma_h$ . Finally, we assume  $R_h$  to be meshwise connected, i.e. for every pair of points in  $R_h$  there is a path which consists of mesh lines connecting the two points.

We will consider a class of finite difference approximations to the linear second-order elliptic equation,

$$(3) \quad Lu := - \sum_{i=1}^n a_i u_{2x_i} + 2 \sum_{i=1}^n b_i u_{x_i} + c = f \quad \text{in } \Omega,$$

with Dirichlet boundary conditions

$$(4) \quad u|_\Gamma = g \quad \text{on } \Gamma.$$

Here  $a_i$ ,  $b_i$  ( $i = 1, \dots, n$ ), and  $c$  are continuous real valued functions on  $\bar{\Omega}$  which satisfy the following conditions:

- (i)  $\exists \bar{\alpha}, \bar{\alpha} > 0$ :  $\bar{\alpha} \leq a_i(x) \leq \bar{\alpha}$  for all  $x \in \bar{\Omega}$ ,  $i \in \{1, \dots, n\}$ ,
- (ii)  $\exists \bar{\beta} > 1$ :  $|b_i(x)| < \bar{\beta}$  for all  $x \in \Omega$ ,  $i \in \{1, \dots, n\}$ ,
- (iii)  $\exists \bar{\gamma} > 0$ :  $0 \leq c(x) \leq \bar{\gamma}$  for all  $x \in \bar{\Omega}$ .

We assume, that problem (3), (4) has a solution  $u$ , which is unique by the maximum principle.

We are now in a position to define the finite difference operator  $L_{h,k}$  which is used in the approximation of (3) and (4). We will use a notation that differs from that of Pereyra et al. [7].  $L_{h,k}$  is a linear operator on the finite dimensional vector space  $F(\bar{\Omega}_h)$  of all real valued functions on  $\bar{\Omega}_h$ . Let  $W \in F(\bar{\Omega}_h)$ .

For  $x \in R_h$ , the operator  $L_{h,k}$  is obtained by replacing  $u_{2x_i}$  by

$$(5a) \quad (D_{2x_i} W)(x) := h^{-2}(W(x - he_i) - 2W(x) + W(x + he_i))$$

and  $u_{x_i}$  by

$$(5b) \quad (D_{x_i} W)(x) := (2h)^{-1}(W(x + he_i) - W(x - he_i)),$$

which gives

$$(6) \quad \begin{aligned} (L_{h,k} W)(x) = & - \sum_{i=1}^n a_i(x)(D_{2x_i} W)(x) \\ & + 2 \sum_{i=1}^n b_i(x)(D_{x_i} W)(x) + c(x)W(x). \end{aligned}$$

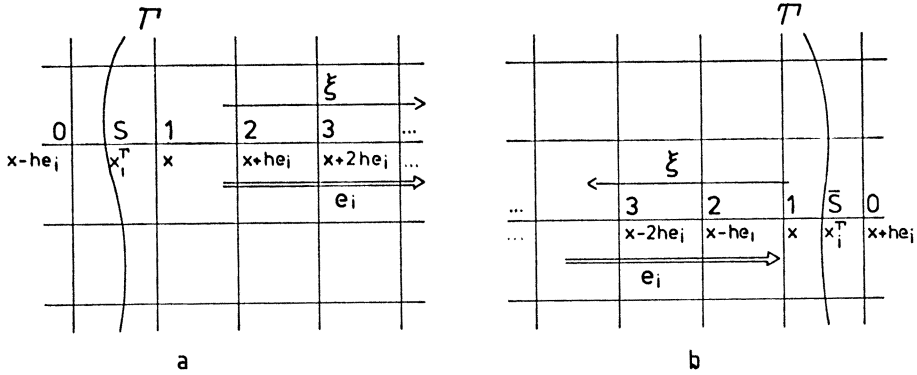


FIGURE 1

For  $x \in \Gamma_h^x$  and  $x - he_i \notin \Omega$ , say, an auxiliary one-dimensional coordinate system  $\xi$  is introduced according to Figure 1a. Let  $Q(\xi) := \sum_{j=0}^k c_j l_j(\xi)$  be a polynomial satisfying  $Q(s) = W(x_i^\Gamma)$ , where  $c_j := W(x + (j - 1)he_i)$ ,  $j \geq 1$ , and  $l_j(\xi)$  is that polynomial of degree  $k$  which satisfies  $l_j(j') = \delta_{j,j'}$ ,  $j, j' \in \{0, \dots, k\}$ . Then the interpolated value of  $W$  in  $x - he_i$  is given by

$$(7a) \quad W(x - he_i) := c_0 = \frac{1}{\alpha_0^k} W(x_i^\Gamma) - \sum_{j=1}^k W(x + (j - 1)he_i) \frac{\alpha_j^k}{\alpha_0^k},$$

where

$$(7b) \quad \alpha_j^k := \alpha_j^k(s) := \prod_{l=0; l \neq j}^k \frac{s - l}{j - l}.$$

As  $x \in \Gamma_h^x$  implies  $x \in \Omega$ , we have  $0 \leq s < 1$ . Therefore,  $c_0$  is well defined for all possible values of  $s$ . Using this provisional value for  $W(x - he_i)$  in (5) gives

$$(8) \quad \begin{aligned} (D_{2x_i}^x W)(x) := & \frac{1}{h^2} \left( \left( -2 - \frac{\alpha_1^k}{\alpha_0^k} \right) W(x) + \left( 1 - \frac{\alpha_2^k}{\alpha_0^k} \right) W(x + he_i) \right. \\ & \left. - \frac{\alpha_3^k}{\alpha_0^k} W(x + 2he_i) - \dots - \frac{\alpha_k^k}{\alpha_0^k} W(x + (k - 1)he_i) \right) \\ & + \frac{1}{\alpha_0^k h^2} W(x_i^\Gamma), \end{aligned}$$

whereas (6) yields

$$\begin{aligned}
 (D_x^x W)(x) &:= \frac{1}{2h} \left( \left( \frac{\alpha_1^k}{\alpha_0^k} \right) W(x) + \left( 1 + \frac{\alpha_2^k}{\alpha_0^k} \right) W(x + he_i) \right. \\
 (9) \quad &\quad \left. + \frac{\alpha_3^k}{\alpha_0^k} W(x + 2he_i) + \dots + \frac{\alpha_k^k}{\alpha_0^k} W(x + (k - 1)he_i) \right) \\
 &\quad - \frac{1}{2h\alpha_0^k} W(x_i^\Gamma).
 \end{aligned}$$

For  $x + he_i \notin \Omega_h$ , the auxiliary coordinate system is defined according to Figure 1b. In that case Eq. (8) remains unchanged whereas in (9) signs are reversed.

It is easily seen that for  $x \pm he_i = x_i^\Gamma$  one gets  $\alpha_j^k = \delta_{0,j}$ ,  $j = 0, \dots, k$ , which means  $c_0 = W(x_i^\Gamma)$ , as expected.

Proper denumeration of points of  $\bar{\Omega}_h$  allows us to write each  $W \in F(\bar{\Omega}_h)$  as  $W = (W^R, W^x, W^d)^T$  where  $W^R$ ,  $W^x$ , and  $W^d$  are functions on  $R_h$ ,  $\Gamma_h^x$ , and  $\Gamma_h$ , respectively.

The finite difference approximation problem to (1) is now given by

$$(10) \quad L_{h,k} U = \begin{bmatrix} \Delta_h f \\ \Delta_h^d g \end{bmatrix},$$

where  $\Delta_h$ ,  $\Delta_h^d$  are the restriction operators to  $\Omega_h$ ,  $\Gamma_h$ , respectively, and  $L_{h,k} W$  decomposes into

$$(11) \quad L_{h,k} W = \begin{bmatrix} L_h^{11} & L_h^{12} & 0 \\ L_{h,k}^{21} & L_{h,k}^{22} & L_{h,k}^{23} \\ 0 & 0 & I \end{bmatrix} \begin{bmatrix} W^R \\ W^x \\ W^d \end{bmatrix},$$

where  $I$  is an identity matrix of appropriate size.

The first line of  $L_{h,k}$ , which is independent of  $k$ , gives  $(L_{h,k} W)(x)$  for  $x \in R_h$ , the second for  $x \in \Gamma_h^x$ , and the third incorporates the boundary conditions.

**3. Two Lemmas.** In order to avoid unnecessary notational complexity, we will restrict ourselves to the cases  $n = 2$  and  $k = 4$ . The generalizations to  $n > 2$  and  $k \leq 3$  are straightforward. However, to show the dependence of the results on  $k$ , we will continue to write  $k$  instead of its specific value 4.

For the sequel, we assume the following conditions (cond) to be satisfied:

(i) Let  $\Lambda_k := (h_n)_{n \in \mathbb{N}}$  be a monotonically decreasing sequence of positive real numbers, satisfying  $\lim_{n \rightarrow \infty} h_n = 0$ , such that there is a mesh  $\mathbf{R}_{h_n}^2$  for each  $h_n$ , which satisfies the condition of Section 2.

(ii) The unique solution  $u$  of (3), (4) and the functions  $e^{(1)}$ ,  $e^{(2)}$ , defined below, are smooth enough, so that all derivatives needed exist and are uniformly bounded. (In a tedious definition, exact differentiability conditions for these functions, which depend on  $k$ , have been given in [5].)

On  $\Omega$  we define

$$(12) \quad D_1(L, u) := \sum_{i=1}^2 \left( -\frac{2a_i}{4!} u_{4x_i} + \frac{2b_i}{3!} u_{3x_i} \right),$$

$$(13) \quad D_2(L, u, e^{(1)}) := \sum_{i=1}^2 \left( -2a_i \left( \frac{u_{6x_i}}{6!} + \frac{e_{4x_i}^{(1)}}{4!} \right) + 2b_i \left( \frac{u_{5x_i}}{5!} + \frac{e_{3x_i}^{(1)}}{3!} \right) \right),$$

where  $e^{(1)}$  is the solution of

$$(14) \quad \begin{aligned} Le^{(1)} &= D_1(L, u) \quad \text{in } \Omega, \\ e^{(1)}|_{\Gamma} &= 0. \end{aligned}$$

$e^{(2)}$  will then denote the solution of

$$(15) \quad \begin{aligned} Le^{(2)} &= D_2(L, u, e^{(1)}) \quad \text{in } \Omega, \\ e^{(2)}|_{\Gamma} &= 0. \end{aligned}$$

Finally, let  $V, E(h, k) \in F(\bar{\Omega}_h)$  be defined by

$$(16) \quad V := \bar{\Delta}_h(u + h^2e^{(1)} + h^4e^{(2)}),$$

$$(17) \quad E(h, k)(x) := \begin{cases} |(L_{h,k}V)(x) - (\Delta_h f)(x)|, & x \in \Omega_h, \\ 0, & x \in \Gamma_h, \end{cases}$$

where  $\bar{\Delta}_h$  is the restriction operator to  $\bar{\Omega}_h$ .

*Remark.* Theorem 1 will show that for  $k = 2, 3$  it is sufficient to define  $V := \bar{\Delta}_h(u + h^2e^{(1)})$ , whereas for  $k = 1$ ,  $V = \bar{\Delta}_h u$  will do. This reduces the smoothness conditions on  $u$  (and  $e^{(1)}$ ) for  $k \leq 3$ . Details have been worked out in [5].

The following two lemmas are generalizations of lemmas used by Pereyra et al. [7, Section 4]. The proofs of these lemmas have been worked out in [5]. It does not seem appropriate to present all the arguments here.

LEMMA 1. *Let condition (cond) be satisfied. Then there exists a constant  $C_1 > 0$  such that*

$$E(h, k)(x) \leq C_1 h^{k+1}$$

for all  $h \in \Lambda_k$  and  $x \in R_h$ .

*Proof.* The proof is straightforward using Taylor expansion of  $u, e^{(1)}$ , and  $e^{(2)}$  about  $x$  in  $(L_{h,k}V)(x)$  and the uniform boundedness of the partial derivatives of  $u, e^{(1)}$ , and  $e^{(2)}$  occurring in the last terms of the Taylor expansion.

LEMMA 2. *Let condition (cond) be satisfied. There exists a constant  $C_{2,k}$ , depending on  $k$ , so that*

$$E(h, k)(x) \leq C_{2,k} h^{k-1}$$

for all  $h \in \Lambda_k$  and  $x \in \Gamma_h^x$ .

*Proof.* To prove this lemma, we have to handle the case where provisional values for  $V$  outside  $\Omega$  have to be calculated. We assume the existence of Taylor expansions of  $u, e^{(1)}$ , and  $e^{(2)}$  about  $x_i^\Gamma$  (Figure 1) of sufficient order, thus continuing these functions sufficiently smoothly along the corresponding mesh line. The interpolation values of  $u, e^{(1)}$  and  $e^{(2)}$  in  $x + he_i$  and  $x - he_i$ , respectively, can now be interpreted as interpolation values of these continued functions and therefore can be replaced by the values of the continued functions with the appropriate interpolation error terms added. These additional error terms are responsible for

the factor  $h^{k-1}$  in Lemma 2 instead of  $h^{k+1}$  as in Lemma 1. As the continued functions are smooth enough, Taylor expansion about  $x$  is possible, so that the remainder of the proof is similar to the proof of Lemma 1.

**4. Asymptotic Expansions.** For  $h \in \Lambda_k$ , we define the operator  $\check{L}_{h,k}$  on  $F(\bar{\Omega}_h)$  by

$$(18) \quad \check{L}_{h,k} := \begin{bmatrix} L_h^{11} & L_h^{12} & 0 \\ 0 & I^{22} & 0 \\ 0 & 0 & I \end{bmatrix},$$

where  $I^{22}$  is an identity matrix of appropriate size and  $I$  is the corresponding matrix of (11).

We have the following discrete maximum principle.

**LEMMA 3.** *Assume condition (cond) to be satisfied. There exists an  $n_k \in \mathbf{N}$  such that for  $h \in \Lambda'_h := \{h_n \in \Lambda_k | n \geq n_k\}$  and  $W_h \in F(\bar{\Omega}_h)$ ,  $(\check{L}_{h,k} W_h)(x) \leq 0$  for all  $x \in R_h$  implies*

$$(19) \quad W_h(x) \leq \max\left(0, \max_{y \in \Gamma_h^x \cup \Gamma_h} (W(y))\right) \text{ for all } x \in \bar{\Omega}_h.$$

Furthermore,  $\check{L}_{h,k}$  is monotone (Young [9, p. 44]). Specifically,  $\check{L}_{h,k}$  is nonsingular and its inverse is given by

$$(20) \quad \check{L}_{h,k}^{-1} = \begin{bmatrix} (L_h^{11})^{-1} & -(L_h^{11})^{-1} L_h^{12} & 0 \\ 0 & I^{22} & 0 \\ 0 & 0 & I \end{bmatrix}.$$

*Proof.* Choose  $n_k \in \mathbf{N}$  such that  $h < \bar{\alpha}/\bar{\beta}$  for all  $h \in \Lambda'_k$ . Observing the special structure of  $L_{h,k}$  and the meshwise connectedness of  $R_h$ , the lemma now follows immediately from Theorem 1 and Theorem 3 of Ciarlet [3].

For  $h \in \Lambda_k$ ,  $W \in F(\bar{\Omega}_h)$ , and an operator  $L_h$  on  $F(\bar{\Omega}_h)$ , let  $l_h(x, y)$  denote that element of the matrix representation of  $L_h$  which is multiplied with  $W(y)$  in the computation of  $(L_h W)(x)$ . The operator norm  $|\cdot|_\infty$  generated by the maximum norm  $\|\cdot\|_\infty$  on  $F(\bar{\Omega}_h)$  is then given by

$$(22) \quad |L_h|_\infty := \max_{x \in \bar{\Omega}_h} \left( \sum_{y \in \bar{\Omega}_h} |l_h(x, y)| \right).$$

**LEMMA 4.** *Let condition (cond) be satisfied. Then there are constants  $x_1^0, x_2^0, \eta \in \mathbf{R}$ , and  $n_k \in \mathbf{N}$  such that*

$$(23) \quad |\check{L}_{h,k}^{-1}|_\infty \leq C_m := \max_{(x_1, x_2) \in \Gamma} (\exp(\eta[(x_1 - x_1^0) + (x_2 - x_2^0)]))$$

for all  $h \in \Lambda'_k := \{h_n \in \Lambda_k | n \geq n_k\}$ .

*Proof.* The technique used in this proof goes back to L. Bers; cf. [1]. Let  $\eta := 2 \cdot \bar{\beta}^2/\bar{\alpha} + \bar{\gamma} + 1$  and choose  $x_1^0, x_2^0$  such that  $x_1 - x_1^0 \geq 0, x_2 - x_2^0 \geq 0$  for all  $x = (x_1, x_2) \in \bar{\Omega}$ , which is always possible as  $\Omega$  is bounded.  $v: \mathbf{R}^2 \rightarrow \mathbf{R}$  can now be defined by

$$(24) \quad \mathbf{R}^2 \ni x = (x_1, x_2) \mapsto v(x_1, x_2) := \exp(\eta[(x_1 - x_1^0) + (x_2 - x_2^0)]).$$

For  $h \in \Lambda_k$ , let  $W_h := \bar{\Delta}_h v \in F(\bar{\Omega}_h)$ . Obviously,  $W_h(x) > 0$  for all  $x \in \bar{\Omega}_h$ . For  $x \in \mathbf{R}^2$ , Taylor expansion gives, for  $i = 1, 2$ ,

$$\begin{aligned}
 (25) \quad 2v(x) - v(x - he_i) - v(x + he_i) &= -h^2 v_{2x_i}(x) - \frac{h^4}{24} (v_{4x_i}(\bar{x}) + v_{4x_i}(\bar{\bar{x}})) \\
 &= -h^2 v_{2x_i}(x) - \frac{h^4}{12} v_{4x_i}(x + \delta_i he_i),
 \end{aligned}$$

where  $\bar{x}, \bar{\bar{x}} \in (x - he_i, x + he_i)$ , which is the line segment connecting  $x - he_i$  and  $x + he_i$  and  $-1 \leq \delta_i \leq 1$ . Similarly, we have, for  $i = 1, 2$ ,

$$(26) \quad v(x + he_i) - v(x - he_i) = 2hv_{x_i}(x) + \frac{h^3}{3} v_{3x_i}(x + \delta_{2+i} he_i),$$

where  $-1 \leq \delta_{2+i} \leq 1$ .

Using (24)–(26) in the definition of  $\check{L}_{h,k} W_h$ , gives, for  $x \in R_h$ ,

$$\begin{aligned}
 (27) \quad (\check{L}_{h,k} W_h)(x) &= \exp(\eta[(x_1 - x_1^0) + (x_2 - x_2^0)]) \\
 &\quad \cdot \left[ (-a_1 \eta^2 - a_2 \eta^2 + 2b_1 \eta + 2b_2 \eta + c) \right. \\
 &\quad \quad \left. + h^2 \left( -\frac{a_1}{12} \eta^4 \exp(\eta \delta_1 h) - \frac{a_2}{12} \eta^4 \exp(\eta \delta_2 h) \right. \right. \\
 &\quad \quad \quad \left. \left. + \frac{b_1}{3} \eta^3 \exp(\eta \delta_3 h) + \frac{b_2}{3} \eta^3 \exp(\eta \delta_4 h) \right) \right] \\
 &\leq \exp(\eta[(x_1 - x_1^0) + (x_2 - x_2^0)]) \\
 &\quad \cdot \left[ (-2\bar{\alpha} \eta^2 + 4\bar{\beta} \eta + \bar{\gamma}) \right. \\
 &\quad \quad \left. + h^2 \left( -\frac{\bar{\alpha}}{6} \eta^4 \exp(-\eta h) + \frac{2}{3} \bar{\beta} \eta^3 \exp(\eta h) \right) \right],
 \end{aligned}$$

where we used  $a_i := a_i(x)$ ,  $b_i := b_i(x)$  ( $i = 1, 2$ ),  $c = c(x)$ , and  $\eta > 0$ .

As a consequence of our definition of  $\eta$ , we have

$$\begin{aligned}
 (28) \quad &-2\bar{\alpha} \eta^2 + 4\bar{\beta} \eta + \bar{\gamma} \\
 &= -8 \frac{\bar{\beta}^4}{\bar{\alpha}} - 8\bar{\beta}^2(\bar{\gamma} + 1) - 2\bar{\alpha}(\bar{\gamma} + 1)^2 + \frac{8\bar{\beta}^3}{\bar{\alpha}} + 4\bar{\beta}(\bar{\gamma} + 1) + \bar{\gamma} \\
 &\leq -4\bar{\beta}^2(\bar{\gamma} + 1) - 2\bar{\alpha}(\bar{\gamma} + 1)^2 + \bar{\gamma} \\
 &\leq -4\bar{\gamma} - 4 + \bar{\gamma} \leq -4,
 \end{aligned}$$

which makes the first term in (27) less than  $-4$ .

We can choose  $n_k \in \mathbb{N}$  such that for all  $h \in \Lambda'_k$  the statements of Lemma 3 and the following inequalities are true:

$$(29) \quad (i) \quad h < \frac{1}{\eta}, \quad (ii) \quad h^2 \left( -\frac{\bar{\alpha}}{6} \eta^4 \frac{1}{e} + \frac{8\bar{\beta}}{3} \eta^3 \cdot e \right) < 3.$$

This choice of  $n_k$  makes the second summand of (27) less than 3.

Hence, we have

$$(30) \quad (\check{L}_{h,k} W_h)(x) < -1$$

for all  $h \in \Lambda'_k$  and  $x \in R_h$ .

For  $h \in \Lambda'_k$ , the operator  $\check{L}_{h,k}^{-1}$  exists by Lemma 3, and we can define  $Q_{h,k} := \check{L}_{h,k}^{-1}E \in F(\bar{\Omega}_h)$ ,  $E(x) := 1$  for all  $x \in \bar{\Omega}_h$ , which is equivalent to  $Q_{h,k}(x) = \sum_{y \in \bar{\Omega}_h} \check{l}_{h,k}^{-1}(x, y) = \sum_{y \in \bar{\Omega}_h} |l_{h,k}^{-1}(x, y)|$  for all  $x \in \bar{\Omega}_h$ , by the monotonicity of  $\check{L}_{h,k}$ . Here,  $l_{h,k}^{-1}(x, y)$  is the appropriate element of the matrix representation of  $\check{L}_{h,k}^{-1}$  (cf. (22)).

This implies  $Q_{h,k}(x) = 1$  for all  $x \in \Gamma_h^x \cup \Gamma_h$ . For  $x \in R_h$ , we have

$$(\check{L}_{h,k}(W_h + Q_{h,k}))(x) < -1 + 1 = 0,$$

by (30) and the definition of  $Q_{h,k}$ . Use of the discrete maximum principle of Lemma 3 for  $\check{L}_{h,k}$  gives for  $x \in R_h$ ,

$$\begin{aligned} (31) \quad Q_{h,k}(x) + 1 &\leq Q_{h,k}(x) + \min_{x \in \bar{\Omega}_h} W_h(x) \leq (W_h + Q_{h,k})(x) \\ &\leq \max\left(0, \max_{\tilde{x} \in \Gamma_h^x \cup \Gamma_h} ((W_h + Q_{h,k})(\tilde{x}))\right) \\ &\leq \max_{\tilde{x} \in \Gamma_h^x \cup \Gamma_h} (W_h(\tilde{x})) + 1 \\ &\leq \max_{(x_1, x_2) \in \Gamma} (\exp(\eta[(x_1 - x_1^0) + (x_2 - x_2^0)])) + 1. \end{aligned}$$

As the last expression of (31) is always greater than 2,

$$\begin{aligned} (32) \quad |\check{L}_{h,k}^{-1}|_\infty &:= \max_{x \in \bar{\Omega}_h} \left( \sum_{y \in \bar{\Omega}_h} |l_{h,k}^{-1}(x, y)| \right) \\ &\leq \max_{(x_1, x_2) \in \Gamma} (\exp(\eta[(x_1 - x_1^0) + (x_2 - x_2^0)])), \end{aligned}$$

which proves the lemma.

LEMMA 5. *There is a constant  $d_k, 0 \leq d_k < 1$ , such that*

$$(33) \quad d_k \left[ 2 + (1 - \varepsilon) \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right] \geq (1 + \varepsilon) \cdot \left[ \sum_{i=2}^k \frac{|\alpha_i^k(s)|}{|\alpha_0^k(s)|} + 1 \right]$$

for all  $s \in [0, 1)$  and for all  $0 \leq \varepsilon \leq 1/10$ .

*Proof.* By definition of  $\alpha_j^k$ , we have

$$(34) \quad \alpha_0^k(s) > 0 \quad \text{for all } s \in [0, 1),$$

$$(35) \quad \text{sgn } \alpha_j^k(s) = (-1)^{j+1} \quad \text{for all } j \in \{1, \dots, k\}, s \in (0, 1),$$

$$(36) \quad \alpha_j^k(0) = 0 \quad \text{for all } j \in \{1, \dots, k\}.$$

As both sides of (33) depend monotonously on  $\varepsilon$ , we can assume  $\varepsilon = 1/10$ . When we set  $d_k := 19/20$  and use (34)–(36), then (33) is equivalent to

$$24(160\alpha_0^4(s) + 171\alpha_1^4(s) + 220\alpha_2^4(s) - 220\alpha_3^4(s) + 220\alpha_4^4(s)) \geq 0.$$

As

$$\begin{aligned} (37) \quad &24(160\alpha_0^4(s) + 171\alpha_1^4(s) + 220\alpha_2^4(s) - 220\alpha_3^4(s) + 220\alpha_4^4(s)) \\ &= 4(474s^4 - 3371s^3 + 6909s^2 - 3946s + 960) \\ &\geq 4(-3371s^3 + 6909s^2 - 3946s + 960) =: p(s), \end{aligned}$$

it suffices to show  $p(s) \geq 0$  for all  $s \in [0, 1]$ .



Applying the well-known division algorithm for polynomials to

$$f_0(s) := p(s), \quad f_1(s) := p'(s) = -10113s^2 + 138185s - 3946,$$

one gets

$$f_2(s) = -\frac{5218922}{10113}s - \frac{620842}{10113}, \quad f_3(s) = \zeta = \text{const.}$$

Hence, the number  $Z(x)$  of sign changes of the sequence  $(f_0(x), f_1(x), f_2(x), \text{ and } f_3(x))$  is

$$Z(1) = Z(0) = \begin{cases} 2 & \text{if } \text{sgn } \zeta = 1, \\ 1 & \text{for all other cases.} \end{cases}$$

By a well-known theorem of Sturm, this implies that there are no zeros of  $p(s) = f_0(s)$  in  $[0, 1)$ . Hence,  $p(1) > 0$  gives  $p(s) > 0$  for all  $s \in [0, 1]$ .

*Remark.* The cases  $k \leq 3$  can be proved by elementary calculation. Any choice of  $d_k$  which satisfies  $1 > d_k > 0.59$  will do.

For  $k \geq 5$ , an  $\bar{s} \in [0, 1)$  can be found such that

$$2 + \frac{\alpha_1^k(\bar{s})}{\alpha_0^k(\bar{s})} \leq \sum_{l=2}^k \frac{|\alpha_l^k(\bar{s})|}{|\alpha_0^k(\bar{s})|} + 1.$$

Therefore, Lemma 5 is not valid in that case. This is the reason why we can show the validity of our expansions only for the cases  $k \leq 4$ .

LEMMA 6. *Let condition (cond) be satisfied. Then there are constants  $n_k \in \mathbb{N}$  and  $\delta_k, 0 < \delta_k < 1$ , such that for all  $h \in \Lambda'_k := \{h_n \in \Lambda_k | n > n_k\}$  and all  $x \in \Gamma_h^x$*

$$(38) \quad \delta_k |l_{h,k}(x, x)| \geq \sum_{y \in \Omega_h; y \neq x} |l_{h,k}(x, y)|.$$

Furthermore

$$(39) \quad |l_{h,k}(x, x)| \geq 4\bar{\alpha}h^{-2} \quad \text{for all } x \in \Gamma_h^x.$$

*Proof.* Choose  $n_k$  such that  $h\bar{\beta}/\bar{\alpha} \leq 10^{-1}$  for all  $h \in \Lambda'_k$  and define

$$(40) \quad \delta_k := (\bar{\bar{\alpha}} + \bar{\alpha}d_k) / (\bar{\bar{\alpha}} + \bar{\alpha}),$$

where  $d_k$  is the constant of Lemma 5. By

$$(41) \quad 1 - \delta_k = \bar{\alpha}(1 - d_k) / (\bar{\bar{\alpha}} + \bar{\alpha}) > 0$$

and

$$(42) \quad \delta_k - d_k = \bar{\bar{\alpha}}(1 - d_k) / (\bar{\bar{\alpha}} + \bar{\alpha}) > 0,$$

we have  $0 \leq d_k < \delta_k < 1$ .

For  $x \in \Gamma_h^x$ , we distinguish the following two cases.

(i) One of the four closest neighbors  $x \pm he_i$  ( $i = 1, 2$ ) of  $x$  is not in  $\Omega_h$ , i.e. interpolation is necessary for exactly one coordinate direction.

(ii) Two of the closest neighbors of  $x$  are not in  $\Omega_h$ .

(i) Let  $x + he_1 \notin \Omega_h$ . By (34)–(36) and our choice of  $n_k$ , we have

$$\begin{aligned}
 h^2 |l_{h,k}(x, x)| &= 2a_1(x) + 2a_2(x) + |a_1(x) - hb_1(x)| \cdot \left| \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right| + c(x) \\
 &\geq a_1(x) \left[ 2 + \left| 1 - h \frac{b_1(x)}{a_1(x)} \right| \cdot \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right] + 2a_2(x) \\
 (43) \quad &\geq a_1(x) \left[ 2 + \left| 1 - h \frac{\bar{\beta}}{\bar{\alpha}} \right| \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right] + 2a_2(x) \\
 &\geq a_1(x) \left[ 2 + \frac{9}{10} \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right] + 2a_2(x) \geq 4\bar{\alpha},
 \end{aligned}$$

which proves the second part of the lemma. Further,

$$\begin{aligned}
 h^2 \cdot \sum_{y \in \Omega_h; y \neq x} |l_{h,k}(x, y)| \\
 (44) \quad &= a_1(x) \left[ \sum_{l=2}^k \left( \left| 1 - h \frac{b_1(x)}{a_1(x)} \right| \left| \frac{\alpha_l^k(s)}{\alpha_0^k(s)} \right| \right) + \left| 1 + h \frac{b_1(x)}{a_1(x)} \right| \right] + 2a_2(x) \\
 &\leq a_1(x) \left[ \sum_{l=2}^k \frac{11}{10} \left| \frac{\alpha_l^k(s)}{\alpha_0^k(s)} \right| + \frac{11}{10} \right] + 2a_2(x).
 \end{aligned}$$

(41)–(44) and Lemma 5 now give

$$\begin{aligned}
 h^2 \delta_k \cdot |l_{h,k}(x, x)| &\geq \delta_k \left[ a_1(x) \left( 2 + \frac{9}{10} \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right) + 2a_2(x) \right] \\
 &\geq 2\delta_k a_2(x) + 2(\delta_k - d_k) a_1(x) + d_k \left( a_1(x) \left( 2 + \frac{9}{10} \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right) \right) \\
 &= 2a_2(x) - 2(1 - \delta_k) a_2(x) + 2(\delta_k - d_k) a_1(x) \\
 &\quad + d_k \left( a_1(x) \left( 2 + \frac{9}{10} \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right) \right) \\
 (45) \quad &\geq 2a_2(x) - 2(1 - \delta_k) \bar{\alpha} + 2(\delta_k - d_k) \bar{\alpha} + d_k \left( a_1(x) \left( 2 + \frac{9}{10} \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right) \right) \\
 &= 2a_2(x) + d_k a_1(x) \left( 2 + \frac{9}{10} \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right) \\
 &\geq 2a_2(x) + \frac{11}{10} a_1(x) \left[ \sum_{l=2}^k \left| \frac{\alpha_l^k(s)}{\alpha_0^k(s)} \right| + 1 \right] \\
 &> h^2 \sum_{y \in \Omega_h; y \neq x} |l_{h,k}(x, y)|,
 \end{aligned}$$

which proves the lemma. The other cases of (i) are proved similarly.

(ii) Let  $x + he_1, x + he_2 \notin \Omega_h$ . By (34)–(36) and our choice of  $n_k$ , we have

$$\begin{aligned}
 h^2 |l_{h,k}(x, x)| &= 2a_1(x) + |a_1(x) - hb_1(x)| \left| \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right| \\
 &\quad + 2a_2(x) + |a_2(x) - hb_2(x)| \left| \frac{\alpha_1^k(t)}{\alpha_0^k(t)} \right| + c(x) \\
 (46) \quad &\geq a_1(x) \left[ 2 + \left| 1 - h \frac{b_1(x)}{a_1(x)} \right| \left| \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right| \right] + a_2(x) \left[ 2 + \left| 1 - h \frac{b_2(x)}{a_2(x)} \right| \left| \frac{\alpha_1^k(t)}{\alpha_0^k(t)} \right| \right] \\
 &\geq a_1(x) \left[ 2 + \left| 1 - h \frac{\bar{\beta}}{\bar{\alpha}} \right| \left| \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right| \right] + a_2(x) \left[ 2 + \left| 1 - h \frac{\bar{\beta}}{\bar{\alpha}} \right| \left| \frac{\alpha_1^k(t)}{\alpha_0^k(t)} \right| \right] \\
 &\geq a_1(x) \left[ 2 + \frac{9}{10} \frac{\alpha_1^k(s)}{\alpha_0^k(s)} \right] + a_2(x) \left[ 2 + \frac{9}{10} \frac{\alpha_1^k(t)}{\alpha_0^k(t)} \right] \geq 4\bar{\alpha},
 \end{aligned}$$

which proves the second part of the lemma. Further, by our choice of  $n_k$ ,

$$\begin{aligned}
 h^2 \cdot \sum_{y \in \Omega_h; y \neq x} |l_{h,k}(x, y)| &= \sum_{l=2}^k |a_1(x) - hb_1(x)| \left| \frac{\alpha_l^k(s)}{\alpha_0^k(s)} \right| + |a_1(x) + hb_1(x)| \\
 (47) \quad &\quad + \sum_{l=2}^k |a_2(x) - hb_2(x)| \left| \frac{\alpha_l^k(t)}{\alpha_0^k(t)} \right| + |a_2(x) + hb_2(x)| \\
 &\leq a_1(x) \left[ \sum_{l=2}^k \frac{11}{10} \cdot \left| \frac{\alpha_l^k(s)}{\alpha_0^k(s)} \right| + \frac{11}{10} \right] + a_2(x) \left[ \sum_{l=2}^k \frac{11}{10} \cdot \left| \frac{\alpha_l^k(t)}{\alpha_0^k(t)} \right| + \frac{11}{10} \right],
 \end{aligned}$$

which together with (46),  $\delta_k \geq d_k$  and Lemma 5 again proves the lemma.

The other cases of (ii) are treated in a similar way.

*Remark.* This lemma holds for arbitrary space dimension  $n$ . The arguments of the proof have to be modified to account for the space dimension.

LEMMA 7. Let condition (cond) be satisfied. Then there is an  $n_k \in \mathbb{N}$  such that, for all  $h \in \Lambda'_k := \{h_n \in \Lambda_k \mid n \geq n_k\}$ ,

(a)  $|(L_h^{11})^{-1}L_h^{12}|_\infty = 1$  (cf. Eq. (22)).

(b)  $L_{h,k}$  is nonsingular.

*Proof.* Let  $n_k$  be the maximum of the corresponding constants of Lemma 3 and Lemma 6.

(i) For  $x \in R_h$ , the choice of  $n_k$  ensures

$$(48) \quad \sum_{y \in \bar{\Omega}_h} l_{h,k}(x, y) = \sum_{y \in \Omega_h} l_{h,k}(x, y) = 0.$$

(ii) Let  $W \in F(\bar{\Omega}_h)$ ,  $W(x) = 1$  for all  $x \in \bar{\Omega}_h$ . Lemma 3 ensures the existence of  $(L_h^{11})^{-1}$ . Hence, using (48), we have

$$\begin{aligned}
 W|_{R_h} &= (L_h^{11})^{-1} L_h^{11}(W|_{R_h}) \\
 (49) \quad &= (L_h^{11})^{-1} \cdot (L_h^{11}(W|_{R_h}) + L_h^{12}(W|_{\Gamma_h^x})) - (L_h^{11})^{-1} L_h^{12}(W|_{\Gamma_h^x}) \\
 &= - (L_h^{11})^{-1} L_h^{12}(W|_{\Gamma_h^x}).
 \end{aligned}$$

By the monotonicity of  $\check{L}_{h,k}$  (Lemma 3), this implies (a).

(iii) Let  $W \in F(\bar{\Omega}_h)$  satisfy  $L_{h,k}W = 0$ , and assume the existence of  $x \in \bar{\Omega}_h$  such that  $|W(x)| > 0$ . The definition of  $L_{h,k}$  immediately gives  $W|_{\Gamma_h} = 0$ .

Hence,  $L_{h,k}W = 0$  is equivalent to

$$(50a) \quad L_h^{11}W|_{R_h} + L_h^{12}W|_{\Gamma_h^x} = 0$$

and

$$(50b) \quad L_{h,k}^{21}W|_{R_h} + L_{h,k}^{22}W|_{\Gamma_h^x} = 0.$$

(50b) is equivalent to

$$- \sum_{y \in \Omega_h; y \neq x} l_{h,k}(x, y)W(y) = l_{h,k}(x, x)W(x)$$

for  $x \in \Gamma_h^x$ .

Lemma 6 now gives

$$\begin{aligned}
 |l_{h,k}(x, x)W(x)| &\leq \sum_{y \in \Omega_h; y \neq x} |l_{h,k}(x, y)W(y)| \\
 (51) \quad &\leq \max_{y \in \Omega_h; y \neq x} |W(y)| \cdot \sum_{y \in \Omega_h; y \neq x} |l_{h,k}(x, y)| \\
 &\leq \max_{y \in \Omega_h; y \neq x} |W(y)| \cdot \delta_k |l_{h,k}(x, x)|.
 \end{aligned}$$

If there were  $\bar{x} \in \Gamma_h^x$  with  $W(\bar{x}) = \max_{y \in \Omega_h} |W(y)| > 0$ , we would have

$$(l_{h,k}(x, x) \neq 0), \quad 0 < |W(\bar{x})| \leq \delta_k \max_{y \in \Omega_h; y \neq x} |W(y)| < \max_{y \in \Omega_h; y \neq x} |W(y)|,$$

which contradicts the maximality of  $W(\bar{x})$ . Hence, the maximum of  $|W(x)|$  is assumed in a point  $\bar{x} \in R_h$ . Using (ii) and (50a), this gives

$$|W(\bar{x})| = \|W|_{R_h}\|_\infty \leq |(L_h^{11})^{-1} L_h^{12}|_\infty \cdot \|W|_{\Gamma_h^x}\|_\infty = \|W|_{\Gamma_h^x}\|_\infty.$$

This implies that there must be an  $x \in \Gamma_h^x$  where  $|W(x)|$  assumes its maximum, too. We just saw that this is possible. Therefore there cannot exist an  $x \in \bar{\Omega}_h$  with  $|W(x)| > 0$ , which proves part (b) of the lemma. We are now in a position to state and prove the main theorem of this paper.

**THEOREM 1.** *Let  $\Omega \subset \mathbb{R}^2$  be a bounded region and  $k \in \{1, \dots, 4\}$ . Assume, further, condition (cond) to be satisfied. Then there is an  $n_k \in \mathbb{N}$  such that, for all  $h \in \Lambda'_k := \{h_n \in \Lambda_k | n \geq n_k\}$ , there exists a unique solution  $U_{h,k}$  of (10). Furthermore,  $n_k$  can be chosen such that there is a constant  $C_k > 0$  giving*

$$(52) \quad \|V - U_{h,k}\|_\infty \leq C_k h^{k+1}$$

for all  $h \in \Lambda'_k$ .

*Proof* (cf. [1]). Let  $n_k$  be the maximum of the corresponding constants of the previous lemmas. Now Lemma 7 immediately gives the existence and the uniqueness of the solution  $U_{k,h}$  of (10).

For  $h \in \Lambda'_k$ , define

$$(53) \quad R_{h,k} := V - U_{h,k}.$$

$V|_{\Gamma_h} = \Delta_h^d g = U_{h,k}|_{\Gamma_h}$  gives  $R_{h,k}|_{\Gamma_h} = 0$ . By Lemma 3,  $\check{L}_{h,k}^{-1}$  exists for all  $h \in \Lambda'_k$ . Hence, for all  $h \in \Lambda'_k$ , we have, using (20),

$$(54) \quad R_{h,k} = (\check{L}_{h,k}^{-1}) \begin{bmatrix} (L_{h,k} R_{h,k})|_{R_h} \\ R_{h,k}|_{\Gamma_h^x} \\ 0 \end{bmatrix} = (\check{L}_{h,k}^{-1}) \begin{bmatrix} (L_{h,k} R_{h,k})|_{R_h} \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -(L_h^{11})^{-1} L_h^{12} R_{h,k}|_{\Gamma_h^x} \\ R_{h,k}|_{\Gamma_h^x} \\ 0 \end{bmatrix},$$

which implies

$$(55) \quad \begin{aligned} \|R_{h,k}\|_\infty &\leq |\check{L}_{h,k}^{-1}|_\infty \cdot \max_{x \in R_h} |(L_{h,k} R_{h,k})(x)| \\ &+ \max \left[ \|-(L_h^{11})^{-1} L_h^{12}\|_\infty \max_{x \in \Gamma_h^x} |R_{h,k}(x)|, \max_{x \in \Gamma_h^x} |R_{h,k}(x)| \right] \\ &\leq C_m \max_{x \in R_h} |(L_{h,k} R_{h,k})(x)| + \max_{x \in \Gamma_h^x} |R_{h,k}(x)|. \end{aligned}$$

By Lemma 1, there exists a constant  $\tilde{C}_1$  such that, for all  $h \in \Lambda'_k$ ,

$$(56) \quad \max_{x \in R_h} |(L_{h,k} R_{h,k})(x)| \leq \tilde{C}_1 h^{k+1}.$$

For  $h \in \Lambda'_k$  and  $x \in \Gamma_h^x$ , observation of  $R_{h,k}|_{\Gamma_h} = 0$  gives

$$(L_{h,k} R_{h,k})(x) = \sum_{y \in \Omega_h} l_{h,k}(x, y) R_{h,k}(y),$$

which implies

$$(57) \quad |R_{h,k}(x)| \leq \sum_{y \in \Omega_h; y \neq x} \frac{|l_{h,k}(x, y)|}{|l_{h,k}(x, x)|} \cdot \max_{y \in \Omega_h; y \neq x} |R_{h,k}(y)| + \frac{|(L_{h,k} R_{h,k})(x)|}{|l_{h,k}(x, x)|}.$$

Lemma 2 and the second part of Lemma 6 ensure

$$(58) \quad \frac{|(L_{h,k} R_{h,k})(x)|}{|l_{h,k}(x, x)|} \leq \frac{1}{4} \frac{\bar{C}_k}{\bar{\alpha}} h^{k+1},$$

whereas the first part of Lemma 6 gives

$$(59) \quad \sum_{y \in \Omega_h; y \neq x} \frac{|l_{h,k}(x, y)|}{|l_{h,k}(x, x)|} \leq \delta_k.$$

Now we use (58) and (59) in (57) and take the maximum of  $|R_{h,k}(x)|$  for  $x \in \Gamma_h^x$ . Then,

$$(60) \quad \max_{x \in \Gamma_h^x} |R_{h,k}(x)| \leq \frac{1}{4} \frac{C_{2,k}}{\bar{\alpha}} h^{k+1} + \delta_k \|R_{h,k}\|_\infty.$$

Using (56) and (60) in (55) yields

$$(61) \quad \|R_{h,k}\|_\infty \leq C_m \tilde{C}_1 h^{k+1} + \frac{1}{4} \frac{C_{2,k}}{\bar{\alpha}} h^{k+1} + \delta_k \|R_{h,k}\|_\infty.$$

Hence,

$$(62) \quad \|R_{h,k}\|_\infty \leq \frac{1}{1 - \delta_k} \left( C_m \tilde{C}_1 + \frac{1}{4} \frac{C_{2,k}}{\bar{\alpha}} \right) h^{k+1}.$$

Setting

$$C_k := \frac{1}{1 - \delta_k} \left( C_m \tilde{C}_1 + \frac{1}{4} \frac{C_{2,k}}{\bar{\alpha}} \right)$$

proves the theorem.

*Remarks.* (a) For  $k = 4$ , the theorem states

$$(63) \quad U_{h,4} = \bar{\Delta}_h(u + h^2 e^{(1)} + h^4 e^{(2)}) + R_{h,4},$$

where  $\|R_{h,4}\|_\infty = O(h^5)$ , as indicated in Section 1.

As the order of convergence of vector functions in the discrete  $L_2$ -norms of Pereyra et al. [7] is at least of the same order of magnitude as in the discrete maximum norm, Theorem 1 is valid with respect to the discrete  $L_2$ -norms, too.

(b) If the order of interpolation  $k$  depends on  $x \in \Gamma_h^x$ , one can only guarantee an expansion of  $U_{h,k}$  which corresponds to the smallest value of  $k$  used.

**5. Numerical Results.** There are two major classes of methods which use the asymptotic expansion of Section 4 to improve the approximations obtained: Richardson extrapolation and deferred correction methods.

In our case, the use of Richardson extrapolation is obviously prohibitive for several reasons (cf. [7]).

The deferred correction methods use, roughly speaking, approximate solutions  $E^{(1)}$  and  $E^{(2)}$  of (14) and (15) respectively instead of  $e^{(1)}$  and  $e^{(2)}$  in the expansion (2) of  $U$ . As the functions  $u$  and  $e^{(1)}$  used in the right-hand sides of (14) and (15) are unknown during the calculation,  $D_1$  and  $D_2$  are approximated by finite difference expressions  $\tilde{D}_1(L, U)$  and  $\tilde{D}_2(L, U_1, E^{(1)})$ , respectively, where  $U_1$  is the approximate solution after the first correction step given by  $U_1 := U - h^2 E^{(1)}$ .

The method just described is not the best possible. To improve the accuracy one has to recall the fact that  $E^{(1)}$  and  $E^{(2)}$  also have expansions of the form (2). Elementary considerations, which have been carried out in [5], show that this fact can be taken into account by using  $\bar{D}_2 := \tilde{D}_2(L, U_1, E^{(1)}) - 2\tilde{D}_1(L, E^{(1)})$  instead of  $\tilde{D}_2(L, U_1, E^{(1)})$  in the approximation of (15). Numerical experiments have been carried out on a Telefunken TR 440 computer at the University of Tübingen. As a model problem we used the Poisson equation  $-\Delta u = f$  in  $\Omega$ ,  $u|_\Gamma = g$  on  $\Gamma$ , where  $\Omega := \{(x_1, x_2) \in \mathbf{R}^2 \mid x_1^2 + x_2^2 < 0.999\}$ . For  $h = 0.1$ , we have 305 meshpoints, at least nine of them on each mesh line meeting  $\Omega$ . In  $\tilde{D}_1$  we alternatively used standard centered five- and seven-point formulas as given by Collatz [4] for the approximation of the fourth-order partial derivatives. As  $b_1 = b_2 = 0$ , third-order derivatives do not occur. In  $\tilde{D}_2$  we used centered seven- and centered five-point formulas (cf. [4]) for the approximation of the sixth- and fourth-order partial derivatives, respectively. Again, third- and fifth-order derivatives do not occur.

Near the boundary of  $\Omega$ , points used in  $\tilde{D}_1$  and  $\tilde{D}_2$  may lie outside  $\Omega$ . For these points provisional values were calculated by one-dimensional extrapolation of fixed degree  $\tilde{k}$  along the corresponding mesh line.

The systems of linear equations were solved by the SOR-method (cf. [8]), as this method can be used for the difference approximation of the general problem (3), (4), too. The overrelaxation parameter  $\omega$  has been determined experimentally by estimating the rate of convergence for different values of  $\omega$ . If one is only concerned with the Poisson equation, direct methods, as used by Pereyra et al. [7], seem to be advantageous.

After testing our program for problems with polynomial solutions of low order, for which the approximate solutions should be very accurate, we ran the program for problems which have the following solutions.

(i)  $u(x_1, x_2) = \sin(x_1 + x_2)$ . This is an example of a very smooth solution. It was chosen to test the general behavior of the algorithms. The accuracy obtained in the three consecutive steps was approximately  $10^{-5}$ ,  $10^{-7}$ ,  $10^{-8}$ .

(ii)  $u(x_1, x_2) = x_1^{10} + x_2^{10}$ . This is an example of a solution, which has a rather steep gradient just outside  $\Omega$ , which may cause increased errors in the polynomial extrapolation. The accuracy obtained was of the order of  $10^{-2}$  to  $10^{-3}$  and hence far below that reached for (i). The correction steps brought about only a very moderate increase in accuracy

(iii)  $u(x_1, x_2) = (r^2 - x_1^2 - x_2^2)^{5/2}$ ,  $r = 0.999$ . This problem was chosen since it does not allow for a Taylor expansion of sufficiently high order, as it is needed in the proof of Lemma 2. (This case is discussed theoretically in [5]. The result is, that there still exists an expansion of the form (2) with the error term  $O(h^{k+1})$  replaced by  $O(h^j)$ , where  $j$  depends on the smoothness of the exact solution  $u$  on  $\Omega$  (and not  $\bar{\Omega}$  as in the proof of Lemma 2).)

The accuracy was of the order of  $10^{-3}$  to  $10^{-4}$  where the first correction step brought about the major improvement, which can be expected if  $j \leq 4$ .

For each problem, the program was run several times varying  $k$ ,  $\tilde{k}$ , and the kind of approximation used in  $\tilde{D}_1$ . Generally, the improvement in the first and second correction step depended heavily on the accuracy of the approximations used in  $\tilde{D}_1$  and  $\tilde{D}_2$ , especially on the choice of  $\tilde{k}$ . When these approximations were poor, the last correction step sometimes even spoiled the solution obtained in the previous steps.

This observation is supported by test runs for (i) and (ii) which used the exact right-hand sides of (14) and (15) and which showed much better results than those using  $\tilde{D}_1$  and  $\tilde{D}_2$ . This indicates that the algorithms may be considerably improved by the use of better approximations in  $\tilde{D}_1$  and  $\tilde{D}_2$ .

Finally, comparison of the runs for  $k = 4$  and  $k = 6$  showed better results for  $k = 6$  than for  $k = 4$ . This supports the conjecture of Pereyra et al. [7] that, for  $k = 6$ , there may be an expansion of the form  $U = \bar{\Delta}_h(u + h^2e^{(1)} + h^4e^{(2)} + h^6e^{(3)}) + O(h^j)$ ,  $j > 6$ .

1. J. H. BRAMBLE & B. E. HUBBARD, "A theorem on error estimation for finite difference analogues of the Dirichlet problem for elliptic equations," *Contributions to Differential Equations*, v. 2, 1963, pp. 319–340.
2. J. H. BRAMBLE & B. E. HUBBARD, "Approximations of derivatives by finite difference methods in elliptic boundary value problems," *Contributions to Differential Equations*, v. 3, 1964, pp. 399–410.
3. P. G. CIARLET, "Discrete maximum principle for finite-difference operators," *Aequationes Math.*, v. 4, 1970, pp. 338–352.
4. L. COLLATZ, *The Numerical Treatment of Differential Equations*, 3rd ed., Die Grundlehren der math. Wissenschaften, Bd. 60, Springer-Verlag, Berlin-Göttingen-Heidelberg, 1960.
5. H. MUNZ, *Differenzenverfahren für elliptische Randwertaufgaben mit verbesserter Randinterpolation*, Diplomarbeit, Universität Tübingen, 1978. (Unpublished.)
6. V. PEREYRA, "Iterated deferred corrections for nonlinear operator equations," *Numer. Math.*, v. 10, 1967, pp. 316–323.
7. V. PEREYRA, W. PROSKUROWSKI & O. WIDLUND, "High order fast Laplace solvers for the Dirichlet problem on general regions," *Math. Comp.*, v. 31, 1977, pp. 1–16.
8. R. S. VARGA, *Matrix Iterative Analysis*, Prentice-Hall, Englewood Cliffs, N. J., 1962.
9. D. M. YOUNG, *Iterative Solution of Large Linear Systems*, Academic Press, New York and London, 1971.