# Suboptimal Enclosures for the Interval Buneman Algorithm for Arbitrary Block Dimension

## Hartmut Schwandt

The interval arithmetic Buneman algorithm is a "fast solver" for a class of block tridiagonal systems with interval coefficients. In [11], we have introduced a modification for arbitrary block dimension. In the present paper, we discuss enclosure properties depending on the block dimension of the coefficient matrix in selected cases for which optimal enclosures cannot be expected.

# Субоптимальные включения в интервальном алгоритме Бунемана при произвольной размерности блоков

## Х. Швандт

Алгоритм Бунемана, использующий интервальную арифметику, представляет собой быстрый решатель для определенного класса блочных трехдиагональных систем с интервальными коэффициентами. В [11] нами предложена модификация этого алгоритма, допускающая произвольную размерность блока. В настоящей работе рассматриваются включающие свойства алгоритма в зависимости от блочной размерности матрицы коэффициентов для некоторых случаев, в которых невозможно получение оптимальных включающих множеств.

# 1   Introduction

The interval arithmetic analogue of one of the "classical" fast Poisson solves, the Buneman algorithm [4, 5], is dedicated to the treatment of systems $(\mathbf{M}, \mathbf{b})$ of linear form with a block tridiagonal matrix $\mathbf{M} = (-\mathrm{S}, \mathbf{A}, -\mathrm{T})$, where S, T are real $p \times p$ matrices, $p \in \mathbb{N}$, where $\mathbf{A}$ is a $p \times p$ matrix with real compact intervals as coefficients and where $\mathbf{b}$ is a vector with $N$ interval components, $N = pq$, $q \in \mathbb{N}$. The specific form of $\mathbf{M}$ results from necessary commutativity conditions, which are satisfied in the main applications of *IBU* as a "linear" solver in Newton-like interval methods for nonlinear systems of equations resulting from difference methods for almost linear partial BVP [8, 9]. In contrast to the solution of linear point (i.e. noninterval) systems, an interval method like the interval Buneman algorithm *IBU* yields an interval vector including the set of solutions $SOL(\mathbf{M}, \mathbf{b}) := \{x \in \mathbb{R}^N | \mathrm{M}x = \mathrm{b}, \mathrm{M} \in \mathbf{M}, \mathrm{b} \in \mathbf{b}\} \subseteq \mathbf{x} := IBU(\mathbf{M}, \mathbf{b})$. As, in general, $SOL(\mathbf{M}, \mathbf{b})$ is not an interval vector, we can at most expect $\mathbf{x}$ to be an optimal enclosure, i.e. the tightest interval vector including this set [1, 3]. In a previous paper [11], we have introduced a modification *IBUD* for arbitrary block dimension. The original Buneman algorithm is defined for a block dimension of the form $q = 2^{r+1} - 1$ [4, 5]. A more flexible choice of $q$ is, however, often desirable if *IBU(D)* is applied as part of other algorithms. As an example we mention a domain decomposition method [12] where *IBUD* is needed for the treatment of systems on rectangular subdomains whose number and size just depend on the restriction on $q$. Suboptimal enclosures which are discussed in the present paper are admissible if only a reduced precision of the results is required or if their difference to an optimal enclosure can be neglected relative to the machine precision. Newton-like methods in which *IBU(D)* is used as a "linear" solver can be typical applications for the latter case. In analogy to the corresponding noninterval algorithm which is applicable to point systems $\mathrm{M}x = \mathrm{b}$, $\mathrm{M} = (-\mathrm{S}, \mathrm{A}, -\mathrm{T})$ of arbitrary block dimension and arbitrary block size [13], the interval method *IBUD* can be carried out for interval systems for any value of $q$, but optimal enclosures can be guaranteed under suitable conditions only for particular values for $q$ while the size of the blocks can be arbitrary. In the present paper, we modify *IBUD* for an efficient treatment of a nonoptimal block dimension and we develop an estimate for the additional width by which it can be decided whether acceptable enclosures can be obtained. The theoretical results are illustrated by numerical examples.

# 2   Notation and basic results

For the introduction of interval arithmetic and interval methods, we refer to [1], for example. We denote real numbers by $a, \ldots, z$ or by Greek letters, real point vectors and matrices by $\mathrm{a}, \ldots, \mathrm{z}$ and $\mathrm{A}, \ldots, \mathrm{Z}$, resp., intervals by $A, \ldots, Z$, real interval vectors and matrices by $\mathbf{a}, \ldots, \mathbf{z}$ and $\mathbf{A}, \ldots, \mathbf{Z}$, resp. We consider real, compact intervals: $A = [\mathrm{i}(A), \mathrm{s}(A)]$, $\mathrm{i}(A) := \min\{a | a \in A\}$, $\mathrm{s}(A) := \max\{a | a \in A\}$. For interval vectors and matrices we use the notation $\mathbf{a} = (A_i)_{i=1}^{N} = [\mathrm{i}(\mathbf{a}), \mathrm{s}(\mathbf{a})] = \big([\mathrm{i}(A_i), \mathrm{s}(A_i)]\big)_{i=1}^{N}$ and $\mathbf{A} = (A_{i,j})_{i,j=1}^{N} = [\mathrm{i}(\mathbf{A}), \mathrm{s}(\mathbf{A})] = \big([\mathrm{i}(A_{i,j}), \mathrm{s}(A_{i,j})]\big)_{i,j=1}^{N}$, where the bounds are defined componentwise. We use the componentwise ordering: $\mathrm{A} \leq \mathrm{B} \Leftrightarrow \forall i, j \in \{1, \ldots, N\} : a_{i,j} \leq b_{i,j}$. $\mathrm{I}(\mathbb{R})$, $\mathrm{V}_N\big(\mathrm{I}(\mathbb{R})\big)$, $\mathrm{M}_{NN}(\mathbb{R})$, $\mathrm{M}_{NN}\big(\mathrm{I}(\mathbb{R})\big)$ denote the sets of real compact intervals, $N$-dimensional interval vectors, $N \times N$ point and interval matrices, resp. The absolute value, midpoint, width of an interval are defined by $|A| = \max\{|\mathrm{i}(A)|, |\mathrm{s}(A)|\}$, $\mathrm{m}(A) = \big(\mathrm{i}(A) + \mathrm{s}(A)\big)/2$, $\mathrm{d}(A) = \big(\mathrm{s}(A) - \mathrm{i}(A)\big)/2$.

We further define $[x] := \min\{i \in \mathbb{N} | i \leq x\}$ for arbitrary $x \in \mathbb{R}^+$.

A real point matrix $\mathrm{A}$ is called an M matrix if $a_{i,j} \leq 0$ for $i \neq j$, if $\mathrm{A}^{-1}$ exists and if $\mathrm{A}^{-1} \geq \mathrm{O}$ (see [14], e.g.). An interval matrix $\mathbf{A}$ is called an interval M matrix if all $\mathrm{A} \in \mathbf{A}$ are M matrices. By $\rho(\mathrm{A})$ and $\mathrm{spec}(\mathrm{A})$, we denote the spectral radius and the spectrum of a quadratic real matrix $\mathrm{A}$.

**Lemma 2.1.** *If* $\mathrm{A}$, $\mathrm{B} \geq \mathrm{O}$, *if all eigenvectors of* $\mathrm{A}$ *are eigenvectors of* $\mathrm{B}$ *and if* $\mathrm{A}$ *has an eigenvector* $\mathrm{c} > 0$, *then* $\|\mathrm{A}\|_{\mathrm{c,c}} \leq \rho(\mathrm{A})$, $\|\mathrm{B}\|_{\mathrm{c,c}} \leq \rho(\mathrm{B})$ *in the matrix norm corresponding to the monotone vector norm* $\|\mathrm{x}\|_{\mathrm{c}} := \max\limits_{1 \leq i \leq N} \left\{ \dfrac{|x_i|}{c_i} \right\}$.

*Proof.* The relation $\mathrm{Ac} = \lambda\mathrm{c}$ with $\lambda \in \mathrm{spec}(\mathrm{A})$ is equivalent to $\lambda = \dfrac{1}{c_i} \sum\limits_{j=1}^{N} a_{i,j} c_j > 0 \, \forall \, 1 \leq i \leq N$. For arbitrary $\mathrm{x} \in \mathbb{R}^N$ this implies

$$\|Ax\|_c = \max_{1\le i\le N}\left\{\frac{1}{c_i}\left|\sum_{j=1}^{N}a_{i,j}x_j\right|\right\} \le \max_{1\le i\le N}\left\{\frac{1}{c_i}\sum_{j=1}^{N}a_{i,j}c_j\frac{|x_j|}{c_j}\right\}$$

$$\le \lambda\|x\|_c \le \rho(A)\|x\|_c.$$

c is also an eigenvector of B, therefore there exists $\mu \in \mathrm{spec}(B)$ such that $Bc = \mu c$, hence $\mu > 0$ and $\|Bx\|_c \le \mu\|x\|_c \le \rho(B)\|x\|_c$ as above. $\qquad\square$

**Corollary 2.2.** *Under the conditions of Lemma 2.1, the following assertions hold:*
*(a) $A \ge O$ irreducible $\Rightarrow \|A\|_{c,c} = \rho(A)$;*
*(b) $Bc = \rho(B)c \Rightarrow \|B\|_{c,c} = \rho(B)$.*

*Proof.* If $A \ge O$ is irreducible, the theorem of Perron-Frobenius implies the existence of a unique eigenvector $c > 0$ to the eigenvalue $\rho(A)$ of A ([14], Th. 2.1). The second assertion is obvious from the proof of Lemma 2.1. $\quad\square$

For elementary rules of interval analysis we refer to [1], for example. For convenience we mention some of these rules which are particularly significant for the present paper.

$$A, B, C, D \in I(\mathbb{R}), \quad A \subseteq B, C \subseteq D, \quad \bullet \in \{+, -, *, /\} \tag{1}$$
$$\Rightarrow A \bullet B \subseteq C \bullet D;$$

$$X = -X \Rightarrow X = \frac{1}{2}d(X)[-1, 1]; \tag{2}$$
$$\mathbf{x} = -\mathbf{x} \Rightarrow A\mathbf{x} = -A\mathbf{x} = |A|\mathbf{x};$$

$$A, B \ge O \Rightarrow \forall \mathbf{x} \in V_N\big(I(\mathbb{R})\big): \quad A(B\mathbf{x}) = (AB)\mathbf{x}; \tag{3}$$

$$\mathrm{i}(\mathbf{A}), \mathrm{i}(\mathbf{B}) \ge O \Rightarrow \forall \mathbf{x} \in V_N\big(I(\mathbb{R})\big): \quad (\mathbf{A} + \mathbf{B})\mathbf{x} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{x}; \tag{4}$$

$$\mathrm{i}(A) \ge 0 \Rightarrow \mathrm{d}(A) \le \mathrm{s}(A); \tag{5}$$

$$A \in M_{NN}(\mathbb{R}) \Rightarrow \forall \mathbf{x} \in V_N\big(I(\mathbb{R})\big): \quad \mathrm{d}(A\mathbf{x}) = |A|\mathrm{d}(\mathbf{x}); \tag{6}$$

$$0 \in A \Rightarrow \forall B \in I(\mathbb{R}): \quad 0 \in AB; \tag{7}$$

$$\mathrm{i}(A), \mathrm{i}(B) \ge 0 \Rightarrow \mathrm{i}(AB) = \mathrm{i}(A)\mathrm{i}(B) \ge 0; \tag{8}$$
$$\mathrm{s}(A) \le 0 \le \mathrm{i}(B) \Rightarrow \mathrm{s}(AB) = \mathrm{s}(A)\mathrm{i}(B) \le 0.$$

The following elementary properties are relevant for the use of any interval method *LES* for the treatment of systems of linear form $(\mathbf{A}, \mathbf{y})$.

$$\forall \mathbf{A} \in \mathrm{M}_{NN}\big(\mathrm{I}(\mathbb{R})\big),\ \mathbf{y}, \mathbf{z} \in \mathrm{V}_N\big(\mathrm{I}(\mathbb{R})\big) :$$
$$LES(\mathbf{A}, \mathbf{y} + \mathbf{z}) \subseteq LES(\mathbf{A}, \mathbf{y}) + LES(\mathbf{A}, \mathbf{z}); \tag{9}$$

$$\forall \mathrm{A} \in \mathrm{M}_{NN}(\mathbb{R}),\ \mathbf{y}, \mathbf{z} \in \mathrm{V}_N\big(\mathrm{I}(\mathbb{R})\big) :$$
$$LES(\mathrm{A}, \mathbf{y} + \mathbf{z}) = LES(\mathrm{A}, \mathbf{y}) + LES(\mathrm{A}, \mathbf{z}); \tag{10}$$

$$\forall \mathbf{A} \in \mathrm{M}_{NN}\big(\mathrm{I}(\mathbb{R})\big),\ \mathbf{y} \in \mathrm{V}_N\big(\mathrm{I}(\mathbb{R})\big), \alpha \in \mathbb{R} :$$
$$LES(\mathbf{A}, \alpha \mathbf{y}) = \alpha LES(\mathbf{A}, \mathbf{y}); \tag{11}$$

$$\forall \mathbf{A}, \mathbf{B} \in \mathrm{M}_{NN}\big(\mathrm{I}(\mathbb{R})\big),\ \mathbf{y}, \mathbf{z} \in \mathrm{V}_N\big(\mathrm{I}(\mathbb{R})\big) :$$
$$\mathbf{A} \subseteq \mathbf{B},\ \mathbf{y} \subseteq \mathbf{z} \Rightarrow LES(\mathbf{A}, \mathbf{y}) \subseteq LES(\mathbf{B}, \mathbf{z}); \tag{12}$$

$$\mathbf{A}_n \to \mathbf{A},\ \mathbf{y}_n \to \mathbf{y}\ (n \to \infty) \Rightarrow$$
$$LES(\mathbf{A}_n, \mathbf{y}_n) \to LES(\mathbf{A}, \mathbf{y})\ (n \to \infty). \tag{13}$$

The tightest interval vector which still includes the set of solutions $SOL(\mathbf{A}, \mathbf{y})$ is called an *optimal enclosure* of $SOL(\mathbf{A}, \mathbf{y})$. The following cases are typical for optimal enclosures:

If $\mathbf{A} \in \mathrm{M}_{NN}\big(\mathrm{I}(\mathbb{R})\big)$ is an interval M matrix, then

a) $\mathbf{A} \equiv \mathrm{A} \in \mathrm{M}_{NN}(\mathbb{R}) \Rightarrow \forall \mathbf{y} \in \mathrm{V}_N\big(\mathrm{I}(\mathbb{R})\big) :$
$$LES(\mathrm{A}, \mathbf{y}) = \big[\mathrm{A}^{-1}\mathrm{i}(\mathbf{y}), \mathrm{A}^{-1}\mathrm{s}(\mathbf{y})\big] = \mathrm{A}^{-1}\mathbf{y}$$

b) $0 \in \mathbf{y} \quad \Rightarrow\ LES(\mathbf{A}, \mathbf{y}) = \big[(\mathrm{i}(\mathbf{A}))^{-1}\mathrm{i}(\mathbf{y}), (\mathrm{i}(\mathbf{A}))^{-1}\mathrm{s}(\mathbf{y})\big];$
$\mathrm{o} \in LES(\mathbf{A}, \mathbf{y})$

c) $0 \le \mathrm{i}(\mathbf{y}) \quad \Rightarrow\ LES(\mathbf{A}, \mathbf{y}) = \big[(\mathrm{s}(\mathbf{A}))^{-1}\mathrm{i}(\mathbf{y}), (\mathrm{i}(\mathbf{A}))^{-1}\mathrm{s}(\mathbf{y})\big];$
$\mathrm{i}\big(LES(\mathbf{A}, \mathbf{y})\big) \ge \mathrm{o}$

d) $\mathrm{s}(\mathbf{y}) \le 0 \quad \Rightarrow\ LES(\mathbf{A}, \mathbf{y}) = \big[(\mathrm{i}(\mathbf{A}))^{-1}\mathrm{i}(\mathbf{y}), (\mathrm{s}(\mathbf{A}))^{-1}\mathrm{s}(\mathbf{y})\big];$
$\mathrm{s}\big(LES(\mathbf{A}, \mathbf{y})\big) \le \mathrm{o}.$

$$\tag{2.14}$$

# 3  The algorithm

The subsequent discussion is based on an algorithm for *IBUD* which has been derived in [11]. For convenience, we repeat its formulation and we briefly recall the main properties. We assume:

> *A system of linear form* $(\mathbf{M}, \mathbf{b})$, *where* $\mathbf{b} \in V_N\big(\mathrm{I}(\mathbb{R})\big)$, $\mathbf{M} = (-\mathrm{S}, \mathbf{A}, -\mathrm{T}) \in \mathrm{M}_{NN}\big(\mathrm{I}(\mathbb{R})\big)$ *block tridiagonal with* $q$ *block rows,* $\mathrm{S}, \mathrm{T} \in \mathrm{M}_{pp}(\mathbb{R})$, $\mathbf{A} \in \mathrm{M}_{pp}\big(\mathrm{I}(\mathbb{R})\big)$, $N = pq,\ p, q \in \mathbb{N}$. $\tag{3.1}$

The system under consideration is one with a block tridiagonal interval matrix with, except for the first and last, identical block rows. In [11], we followed principally the ideas of the Buneman algorithm. Starting from a point system $\mathrm{Mx} = \mathrm{y}$ with $\mathrm{M} \in \mathbf{M}$, $\mathrm{b} \in \mathbf{b}$, block cyclic reduction first yields reduced systems $\mathrm{M}^{(r)}\mathrm{x}^r = \mathrm{b}^r$, $\mathrm{M}^{(r)} \in \mathrm{M}_{pj_r, pj_r}(\mathbb{R})$, $\mathrm{x}^r, \mathrm{b}^r \in \mathrm{V}_{pj_r}(\mathbb{R})$, $0 \le r \le r_q$, $r_q = [\log_2(q)] + 1$, where

$$
\mathrm{M}^{(r)} = \begin{pmatrix}
\mathrm{A}^{(r)} & -\mathrm{T}^{(r)} & & & & \\
-\mathrm{S}^{(r)} & \mathrm{A}^{(r)} & -\mathrm{T}^{(r)} & & & \\
& \ddots & \ddots & \ddots & & \\
& & -\mathrm{S}^{(r)} & \mathrm{A}^{(r)} & -\mathrm{T}^{(r)} & \\
& & & -\mathrm{S}^{(r)} & \mathrm{B}^{(r)}\big(\mathrm{C}^{(r)}\big)^{-1}
\end{pmatrix},
$$

$$
(2)
$$

$$
\mathrm{x}^r = \begin{pmatrix} \mathrm{x}_{2^r} \\ \mathrm{x}_{2^{r+1}} \\ \vdots \\ \mathrm{x}_{j_r - 2^r} \\ \mathrm{x}_{j_r} \end{pmatrix}, \qquad
\mathrm{b}^r = \begin{pmatrix} \mathrm{b}^r_{2^r} \\ \mathrm{b}^r_{2^{r+1}} \\ \vdots \\ \mathrm{b}^r_{j_r - 2^r} \\ \mathrm{b}^r_{j_r} \end{pmatrix}.
$$

In the sequel, superscripts on real numbers, intervals, real and interval vectors and matrices indicate reduction or solution steps of (block) cyclic reduction algorithms. Parentheses are set to avoid any confusion with powers, except for vectors: $\lambda_i^{(r)}$, $\mathrm{y}_j^r$, $X^{(r)}$, $\mathbf{A}^{(r)}$, e.g.

An arbitrary block dimension $q$ can be handled due to the introduction of the matrices $\mathrm{B}^{(r)}$, $\mathrm{C}^{(r)}$ in the last block equation [13]. The instability of this block cyclic reduction procedure is avoided in the Buneman algorithm and its variants by replacing the $\mathrm{b}_j^r$ by auxiliary vectors and by a factorization of the matrices $\mathrm{A}^{(r)}$, $\mathrm{B}^{(r)}$, $\mathrm{C}^{(r)}$. For $0 \le r \le r_q$, we obtain

$$
\begin{aligned}
\mathrm{A}^{(r)} &= \prod_{i=1}^{2^r} \big(\mathrm{A} - \alpha_i^{(r)}\mathrm{U}\big), & \alpha_i^{(r)} &= 2\cos\left(\tfrac{2i-1}{2^{r+1}}\pi\right), \\
\mathrm{B}^{(r)} &= \prod_{i=1}^{k_r} \big(\mathrm{A} - \lambda_i^{(r)}\mathrm{U}\big), & \lambda_i^{(r)} &= 2\cos\left(\tfrac{i}{k_r+1}\pi\right), & (3.3) \\
\mathrm{C}^{(r)} &= \prod_{i=1}^{l_r} \big(\mathrm{A} - \mu_i^{(r)}\mathrm{U}\big), & \mu_i^{(r)} &= 2\cos\left(\tfrac{i}{l_r+1}\pi\right)
\end{aligned}
$$

where $k_0 = 1$, $l_0 = 0$, and

$$k_{r+1} = \begin{cases} k_r + 2^r \\ k_r + 2^{r+1}. \end{cases} \qquad l_{r+1} = \begin{cases} l_r & \text{if } q_r \text{ even,} \\ k_r & \text{if } q_r \text{ odd.} \end{cases} \qquad (3.4)$$

Starting from $M \in \mathbf{M}$, $b \in \mathbf{b}$, the interval arithmetic evaluation [1] of the resulting real point system, mainly by the repeated application of (2.1) and (2.12), leads to the following algorithm (3.5) from [11] to which we refer in the sequel. In (3.5) $L\tilde{E}S$ denotes a partial algorithm for the treatment of partial systems with the matrices $\mathbf{A}^{(r)}$. $L\tilde{E}S$ will be defined in detail after (3.5).

**Modified *IBU* for arbitrary block dimension *(IBUD)*** (3.5)

$q_0 = j_0 = q$, $r_q = [\log_2(q)] + 1$; $r_s := \min\{r | 0 \leq r \leq r_q, q_r \text{ even}\}$

$S^{(0)} = S$; $T^{(0)} = T$; $U^{(0)} = U = \sqrt{ST}$; $\mathbf{x}_0 = \mathbf{x}_{q+1} = \text{o}$;

<u>reduction phase</u>

**for** $j := 1$ **to** $q$ **do**     $\mathbf{p}_j^0 = \text{o}$; $\mathbf{q}_j^0 = \mathbf{b}_j$;

**for** $r := 0$ **to** $r_s - 1$ **do**

$\qquad S^{(r+1)} = \left(S^{(r)}\right)^2$; $T^{(r+1)} = \left(T^{(r)}\right)^2$; $U^{(r+1)} = S^{(r)}T^{(r)}$;

$\qquad q_{r+1} = q_r \text{ div } 2$; $j_{r+1} = j_r - 2^r$; $k_{r+1} = k_r + 2^{r+1}$; $l_{r+1} = k_r$;

$\qquad$ **for** $j := 2^{r+1}$ **step** $2^{r+1}$ **to** $j_{r+1}$ **do**

$$\begin{aligned} \mathbf{p}_j^{r+1} &= \mathbf{p}_j^r + L\tilde{E}S\left(\mathbf{A}^{(r)}, S^{(r)}\mathbf{p}_{j-2^r}^r + T^{(r)}\mathbf{p}_{j+2^r}^r + \mathbf{q}_j^r\right); \\ \mathbf{q}_j^{r+1} &= S^{(r)}\mathbf{q}_{j-2^r}^r + T^{(r)}\mathbf{q}_{j+2^r}^r + 2U^{(r+1)}\mathbf{p}_j^{r+1} \end{aligned} \qquad (3.6a)$$

**for** $r := r_s$ **to** $r_q - 1$ **do**

$\qquad S^{(r+1)} = \left(S^{(r)}\right)^2$; $T^{(r+1)} = \left(T^{(r)}\right)^2$; $U^{(r+1)} = S^{(r)}T^{(r)}$; $q_{r+1} = q_r \text{ div } 2$;

$\qquad$ **if** $q_r$ even **then**

$\qquad\qquad j_{r+1} = j_r$; $k_{r+1} = k_r + 2^r$; $l_{r+1} = l_r$;

$$\begin{aligned} \mathbf{p}_{j_{r+1}}^{r+1} &= \mathbf{p}_{j_{r+1}}^r + L\tilde{E}S\left(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\left(S^{(r)}\mathbf{p}_{j_{r+1}-2^r}^r + \mathbf{q}_{j_{r+1}}^r\right)\right); \\ \mathbf{q}_{j_{r+1}}^{r+1} &= S^{(r)}\mathbf{q}_{j_{r+1}-2^r}^r + U^{(r+1)}\mathbf{p}_{j_{r+1}}^{r+1} \end{aligned} \qquad (3.6b)$$

$\qquad$ **else**

$\qquad\qquad j_{r+1} = j_r - 2^r$; $k_{r+1} = k_r + 2^{r+1}$; $l_{r+1} = k_r$;

$$\mathbf{p}_{j_{r+1}}^{r+1} = \mathbf{p}_{j_{r+1}}^r + L\tilde{E}S\left(\mathbf{A}^{(r)}, S^{(r)}\mathbf{p}_{j_{r+1}-2^r}^r + T^{(r)}\mathbf{p}_{j_{r+1}+2^r}^r + \mathbf{q}_{j_{r+1}}^r\right);$$
$$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad (3.6c)$$
$$\mathbf{q}_{j_{r+1}}^{r+1} = S^{(r)}\mathbf{q}_{j_{r+1}-2^r}^r + U^{(r+1)}\mathbf{p}_{j_{r+1}}^{r+1}$$
$$\quad\quad\quad + L\tilde{E}S\left(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\mathbf{A}^{(r)}\left(T^{(r)}\mathbf{q}_{j_{r+1}+2^r}^r + U^{(r+1)}\mathbf{p}_{j_{r+1}}^{r+1}\right)\right);$$

**for** $j := 2^{r+1}$ **step** $2^{r+1}$ **to** $j_{r+1} - 2^{r+1}$ **do**

$$\mathbf{p}_j^{r+1} = \mathbf{p}_j^r + L\tilde{E}S\left(\mathbf{A}^{(r)}, S^{(r)}\mathbf{p}_{j-2^r}^r + T^{(r)}\mathbf{p}_{j+2^r}^r + \mathbf{q}_j^r\right);$$
$$\mathbf{q}_j^{r+1} = S^{(r)}\mathbf{q}_{j-2^r}^r + T^{(r)}\mathbf{q}_{j+2^r}^r + 2U^{(r+1)}\mathbf{p}_j^{r+1}$$

solution phase

**if** $r_s = r_q$ **then** $\mathbf{x}_{2^{r_q}} = \mathbf{p}_{2^{r_q}}^{r_q} + L\tilde{E}S\left(\mathbf{A}^{(r_q)}, \mathbf{q}_{2^{r_q}}^{r_q}\right)$

$\quad$ **else** $\mathbf{x}_{2^{r_q}} = \mathbf{p}_{2^{r_q}}^{r_q} + L\tilde{E}S\left(\mathbf{B}^{(r_q)}, \mathbf{C}^{(r_q)}\mathbf{q}_{2^{r_q}}^{r_q}\right);$

**for** $r := r_q - 1$ **step** $-1$ **to** $\min\{r_s, r_q - 1\} + 1$ **do**

$\quad$ **for** $j := 2^r$ **step** $2^{r+1}$ **to** $j_r - 2^r$ **do**

$$\quad\quad \mathbf{x}_j = \mathbf{p}_j^r + L\tilde{E}S\left(\mathbf{A}^{(r)}, S^{(r)}\mathbf{x}_{j-2^r} + T^{(r)}\mathbf{x}_{j+2^r} + \mathbf{q}_j^r\right);$$

$\quad$ **if** $q_r$ odd **then**

$$\quad\quad \mathbf{x}_{j_r} = \mathbf{p}_{j_r}^r + L\tilde{E}S\left(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\left(S^{(r)}\mathbf{x}_{j_r-2^r} + \mathbf{q}_{j_r}^r\right)\right);$$

**for** $r := \min\{r_s, r_q - 1\}$ **step** $-1$ **to** $0$ **do**

$\quad$ **for** $j := 2^r$ **step** $2^{r+1}$ **to** $j_r$ **do**

$$\quad\quad \mathbf{x}_j = \mathbf{p}_j^r + L\tilde{E}S\left(\mathbf{A}^{(r)}, S^{(r)}\mathbf{x}_{j-2^r} + T^{(r)}\mathbf{x}_{j+2^r} + \mathbf{q}_j^r\right).$$

For the solution of the partial systems, we define the abstract method $L\tilde{E}S$ as follows. Instead of $\mathbf{A}^{(r)}$, $\mathbf{B}^{(r)}$, $\mathbf{C}^{(r)}$, we use interval extensions [1] $\mathbf{A} - \alpha\mathbf{U}$ of the matrices $A - \alpha U$. *LES* denotes any interval method which is applicable to coefficient matrices $A - \alpha U$. Typical examples are the interval Gauss algorithm *IGA* [1], interval arithmetic cyclic reduction *ICR* for tridiagonal systems [10], the interval Cholesky algorithm [2] or *IBU(D)* itself in discretizations of three-dimensional problems.

$$\mathbf{z} = L\tilde{E}S\left(\mathbf{A}^{(r)}, \mathbf{y}\right): \quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad (3.7)$$

$\quad$ $\mathbf{z}_0 = \mathbf{y}$

$\quad$ **for** $i := 1$ **to** $2^r$ **do** $\mathbf{z}_i = LES\left(\mathbf{A} - \alpha_i^{(r)}\mathbf{U}, \mathbf{z}_{i-1}\right)$

$\quad$ $\mathbf{z} := \mathbf{z}_{2^r}$

$$\mathbf{z} = L\tilde{E}S\left(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\mathbf{y}\right) : \tag{3.8}$$

$\mathbf{z}_0 = \mathbf{y}$

**for** $i := 1$ **to** $l_r$ **do**

$\qquad \mathbf{z}_i = \mathbf{z}_{i-1} + \left(\lambda_i^{(r)} - \mu_i^{(r)}\right)\mathrm{U}LES\!\left(\mathbf{A} - \lambda_i^{(r)}\mathrm{U}, \mathbf{z}_{i-1}\right)$

**for** $i := l_r + 1$ **to** $k_r$ **do** $\mathbf{z}_i = LES\!\left(\mathbf{A} - \lambda_i^{(r)}\mathrm{U}, \mathbf{z}_{i-1}\right)$

$\mathbf{z} := \mathbf{z}_{k_r}.$

As already mentioned in [11], nonoptimal enclosures are caused by step (3.6c). Therefore, we have determined in that paper all values of the block dimension $q$ for which step (3.6c) is never used. In the present context, our goal is to consider nonoptimal values of $q$. We use the following algorithm for (3.6c) which will be justified in the next section:

$$\mathbf{z} = L\tilde{E}S\left(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\mathbf{A}^{(r)}\mathbf{y}\right) :$$

$\mathbf{z}_0 = \mathbf{y}$

**for** $i := 1$ **to** $k_r - 1$ **do**

$\qquad \mathbf{z}_i = \mathbf{z}_{i-1} + \left(\lambda_i^{(r)} - \gamma_{j(i)}^{(r)}\right)\mathrm{U}LES\!\left(\mathbf{A} - \lambda_i^{(r)}\mathrm{U}, \mathbf{z}_{i-1}\right) \tag{3.9}$

$\qquad \mathbf{z}_{k_r} = \mathbf{z}_{k_r-1} + \left(\lambda_{k_r}^{(r)} - \gamma_{j(k_r)}^{(r)}\right)\mathrm{U}LES\!\left(\mathbf{A} - \lambda_{k_r}^{(r)}\mathrm{U}, \mathbf{z}_{k_r-1}\right)$

$\mathbf{z} := \mathbf{z}_{k_r}$

where

$$\left(\gamma_i^{(r)}\right)_{i=1}^{k_r} := \mathrm{merge}\left\{\left(\mu_i^{(r)}\right)_{i=1}^{l_r}, \left(\alpha_i^{(r)}\right)_{i=1}^{2^r}\big|\gamma_{i+1}^{(r)} \leq \gamma_i^{(r)}\right\}, \gamma_1^{(r)} = \alpha_1^{(r)} \tag{3.10}$$

i.e. the $\gamma_i^{(r)}$ are ordered by size. In the sequel, we use the ordering

$$j(i) = i + 1 \quad (1 \leq i \leq k_r - 1); \quad j(k_r) = 1 \tag{3.11}$$

which will be discussed in the next section under the aspect of controllable enclosures. Principally, any ordering is admissible. Note that we have shown in [11] for (3.8) that

$$\forall r \in \{1, \ldots, r_q\} \quad \forall i \in \{1, \ldots, l_r\} : \quad \lambda_i^{(r)} \geq \mu_i^{(r)}. \tag{3.12}$$

We summarize the main properties of *IBUD* from [11] in the following theorem. $\alpha$ is defined as one of the roots $\alpha_i^{(r)}$, $\lambda_i^{(r)}$, $\mu_i^{(r)}$ of (3.3).

**Theorem 3.1.** *Assume (3.1). Under the conditions*

> (a) $\forall A \in \mathbf{A}$: $AS = SA$, $AT = TA$, and $ST = TS$
>
> (b) $U = \sqrt{ST} \geq O$, *i.e.* $U^2 = ST$, *and* $U^{-1}$ *exists*
>
> (c) *LES is an interval method for systems with the coefficient structure of* $(\mathbf{A} - \alpha U, \mathbf{y})$
>
> (d) *LES is applicable to interval M matrices and satisfies (2.9)–(2.13)*
>
> (e) *LES yields optimal enclosures for* $SOL(\mathbf{A} - \alpha U, \mathbf{y})$ *according to (2.14, a–d)*
>
> (f) $\mathbf{M}$ *and* $\mathbf{A} + \alpha_1^{r_q} U$, $\mathbf{A} - \alpha_1^{r_q} U$ *are interval M matrices,*

*the following assertions hold:*

> 1) *IBUD can be applied to* $(\mathbf{M}, \mathbf{b})$ *and satisfies (2.9)–(2.13)*
>
> 2) $SOL(\mathbf{M}, \mathbf{b}) \subseteq LES(\mathbf{M}, \mathbf{b})$
>
> 3) $\forall r \in \{0, \ldots, r_s - 1\} : A^{(r)} = B^{(r)} \left(C^{(r)}\right)^{-1}$; *IBU and IBUD coincide in these steps.*

*If, in addition,*

> (g) $q = 2^n(2^m + 1) - 1$, $n, m \in \mathbb{N}$,

*then*

> 4) *IBUD yields optimal enclosures of (3.1) according to (2.14, a–d).*

The assertions 1 and 2 also hold for interval H matrices [11]. Note that for $\mathbf{M} = (-S, \mathbf{A}, -T)$ condition (f) implies $S, T \geq O$, $U \geq O$. (b) is satisfied, for example, for $S = T$, T invertible, or $S = \alpha I$, $T = \beta I$, $\alpha, \beta > 0$ and (a) for $\mathbf{A} = A + D I$. Typical problems result from standard discretizations of elliptic BVP with difference methods. As simple characteristical examples (see [9]), we mention discretizations of Dirichlet problems $Lu = f$ on a rectangle $\Omega$, $u = g$ on $\partial\Omega$, with central difference quotients on a grid with constant mesh size:

> (i) the nine point formula for $Lu = u_{xx} + u_{yy}$ leading to $A = (-4, 20, -4)$, $S = T = (1, 4, 1)$,

(ii) a five point formula for $Lu = \big(a(x)u_x\big)_x + bu_{yy}$, $a, b > 0$, leading to
   $A = (-\gamma_{i-1}, 2\beta + \gamma_{i-1} + \gamma_i, -\gamma_i)$, $S = T = \beta I$,

(iii) a five point formula for $Lu = a(x)u_{xx} + bu_{yy} + c(x)u_x + du_y + e(x)u$,
   $a, b > 0$, $e \geq 0$, under appropriate conditions on the mesh size leading
   to $A = (-\gamma_i, \varepsilon_i + \alpha + \beta, -\delta_i)$, $S = \alpha I$, $T = \beta I$.

Interval matrices **M**, i.e. **A** $=$ A $+ D$I result, for example, from the application of $IBU(D)$ in Newton-like methods where $D$ is due to the interval arithmetic treatment of the nonlinearity of problems like $f(x) = Mx + \varphi(x) = o$ with a diagonal function $\varphi$, where $\varphi' \geq 0$ on $\Omega$ (see [9]).

# 4   Suboptimal enclosures

Theorem 3.1 shows that $IBUD$ can be carried out for arbitrary $q$ under appropriate conditions while optimal enclosures can be expected only for particular $q$. The nonoptimality is reflected by an increased interval width, due to a change of sign in one or more of the terms $\lambda_i^{(r)} - \gamma_{j(i)}^{(r)}$ in (3.6c). Based on the ordering (3.11), we prove a qualitative result for the dependence of the additional width for nonoptimal $q$ on the diagonal dominance of **A** and we develop an estimate of this width for a restricted set of values for $q$ which are candidates for a moderate increase of the width. In a noninterval system $Mx = y$, the problem posed by (3.6c) is reduced to the solution of systems of the form

$$\prod_{i=1}^{k_r} \big(A - \lambda_i^{(r)}U\big)x = \prod_{i=1}^{l_r} \big(A - \mu_i^{(r)}U\big) \prod_{i=1}^{2^r} \big(A - \alpha_i^{(r)}U\big)y. \qquad (4.1)$$

From the definition of the $k_r$ and $l_r$, it follows that

$$(a)\, k_r = l_r + 2^r \quad (b)\, 0 \leq l_r \leq 2^r \quad (c)\, 2^r \leq k_r < 2^{r+1}. \qquad (4.2)$$

By merging the $\alpha_i^{(r)}$ and $\mu_i^{(r)}$ appropriately $\big($compare (3.10)$\big)$, (4.1) can be simplified to

$$\prod_{i=1}^{k_r} \big(A - \lambda_i^{(r)}U\big)x = \prod_{i=1}^{k_r} \big(A - \gamma_i^{(r)}U\big)y. \qquad (4.3)$$

In the interval case, we have to deal with matrices $\mathbf{A} - \alpha\mathrm{U}$. After some manipulations in (4.1) [13] and a subsequent interval extension of (4.3), we get (3.9). This formula ensures that multiplications by $\mathbf{A}$ on the right-hand sides of interval systems are avoided. In a first step, we justify the ordering (3.11). In the sequel, we assume $r$, $0 \leq r \leq r_q$, to be arbitrary, but fixed.

**Lemma 4.1.** *In the case of nontrivial use of (3.6c), i.e.* $q \neq 2^n(2^m+1)-1$, $n, m \in \mathbb{N}$, *the following relation holds:*

$$\forall r \in \{1, \ldots, r_q\} \, \forall i \in \{1, \ldots, k_r\} : \quad \lambda_i^{(r)} < \alpha_i^{(r)}.$$

*Proof.* According to (4.2) we note $k_r + 1 \leq 2^{r+1}$ where equality only holds if $q_s$ is odd for all $s \leq r$, i.e. $k_r = 2^{r+1} - 1 = \sum_{i=0}^{r} 2^i$. Therefore, $k_r + 1 \leq 2^{r+1}$ which is equivalent to $\frac{i}{k_r+1} > \frac{1}{2^{r+1}}$, i.e. the above assertion holds. $\qquad\square$

**Corollary 4.2.** *There does* **not** *exist an ordering such that*

$$\forall i \in \{1, \ldots, k_r\} :$$
$$\exists j \in \{1, \ldots, 2^r\} : \lambda_i^{(r)} \geq \alpha_j^{(r)} \quad \text{or} \quad \exists k \in \{1, \ldots, l_r\} : \lambda_i^{(r)} \geq \mu_k^{(r)}.$$

Lemma 4.1 and Corollary 4.2 show that optimal enclosures cannot be expected by *IBUD* if (3.6c) and (3.9), resp., are involved. In the sequel, we use the abbreviations

$$l^{(i)} := \frac{i}{l_r + 1}, \quad a^{(i)} := \frac{2i - 1}{2^{r+1}}, \quad k^{(i)} := \frac{i}{k_r + 1}.$$

**Lemma 4.3.** *Each of the intervals* $\left[-2, \lambda_{k_r}^{(r)}\right]$, $\left[\lambda_1^{(r)}, 2\right]$, $\left[\lambda_i^{(r)}, \lambda_{i+1}^{(r)}\right]$, $1 \leq i \leq k_r - 1$, *contains at most one* $\alpha_j^{(r)}$ *and one* $\mu_k^{(r)}$.

*Proof.* We first note $\frac{1}{k_r+1} \leq \frac{1}{2^r} \leq \frac{1}{l_r+1}$ which implies for all $i \in \{1, \ldots, k_r-1\}$, $m \in \{1, \ldots, 2^r - 1\}$, $j \in \{1, \ldots, l_r - 1\}$ that $k^{(i+1)} - k^{(i)} \leq a^{(m+1)} - a^{(m)} \leq l^{(j+1)} - l^{(j)}$. On the other hand, (4.2) yields $\frac{1}{2^{r+1}} \leq \frac{1}{k_r+1} \leq \frac{1}{l_r+1}$, hence $\alpha_1^{(r)} \geq \lambda_1^{(r)} \geq \mu_1^{(r)}$. The corresponding relation for the other boundary interval follows from the symmetry of the cosine. $\qquad\square$

**Lemma 4.4.**

$$x, y, z \in \mathbb{R}^+, \ x = y + z \Rightarrow [x] \geq [y] + [z].$$

*Proof.* Let $f(s)$ be the fractional part of an arbitrary real number, i.e. $s = f(s) + [s]$. The assertion follows from

$$
\begin{array}{lll}
[x] = [y] + [z], & f(x) = f(y) + f(z) & \text{if } f(y) + f(z) \leq 1 \\
[x] = [y] + [z] + 1, & f(x) = f(y) + f(z) - 1 & \text{if } f(y) + f(z) > 1.
\end{array}
$$

$\square$

**Lemma 4.5.** *Define* $\text{card}(k^{(i)})$ *to be the number of* $k^{(i)}$ *(analogously* $l^{(i)}$, $a^{(i)}$*) in a given interval* $[0, t]$, $0 \leq t \leq 1$. *Then*

$$-1 \leq \text{card}(k^{(i)}) - \big(\text{card}(l^{(i)}) + \text{card}(a^{(i)})\big).$$

*Proof.* According to their definition, the distance $k^{(i+1)} - k^{(i)}$ between two consecutive numbers $k^{(i)}$, $k^{(i+1)}$ (similarly $a^{(i+1)} - a^{(i)}$, $l^{(i+1)} - l^{(i)}$) is always the same. This implies $\text{card}(k^{(i)}) = [tk_r]$ and $\text{card}(l^{(i)}) = [tl_r]$. In the case of the $a^{(i)}$, we have to take account of $a^{(1)} - 0 = 1/2^{r+1}$ instead of $1/2^r$ and we get $[t2^r] \leq \text{card}(a^{(i)}) \leq [t2^r] + 1$. The assertion then follows from Lemma 4.4. $\square$

**Corollary 4.6.** *Assume an arbitrary subinterval* $I \subseteq [-2, 2]$. *The number of* $\lambda_i^{(r)}$ *included in* $I$ *differs by at most 1 from the number of* $\alpha_j^{(r)}$, $\mu_k^{(r)}$ *included in* $I$.

We summarize the result of the preceding lemmata in

**Theorem 4.7.** *In the case of nontrivial use of (3.6c), i.e.* $q \neq 2^n(2^m+1)-1$, $n, m \in \mathbb{N}$, *we get*

$$\lambda_{k_r}^{(r)} < \gamma_1^{(r)} \quad \text{and} \quad \lambda_i^{(r)} \geq \gamma_{i+1}^{(r)}, \quad 1 \leq i \leq k_r - 1$$

*for the ordering (3.11) where* $\big(\mu_i^{(r)}\big)_{i=1}^{l_r}$, $\big(\alpha_i^{(r)}\big)_{i=1}^{2^r}$ *are merged in* $\big(\gamma_i^{(r)}\big)_{i=1}^{k_r}$ *according to (3.10).*

*Proof.* According to Lemma 4.1, we have $\lambda_{k_r}^{(r)} \leq \lambda_1^{(r)} < \gamma_1^{(r)} = \alpha_1^{(r)}$ and $\alpha_2^{(r)}, \mu_2^{(r)} \leq \lambda_1^{(r)}$. Therefore, the interval $\left[ -2, \lambda_1^{(r)} \right]$ contains $k_r - 1$ $\gamma_i^{(r)}$, in particular $\gamma_2^{(r)}$, hence $\lambda_1^{(r)} \geq \gamma_2^{(r)}$. Assume now, that $\lambda_j^{(r)} \geq \gamma_{j+1}^{(r)}$ for $1 \leq j \leq i$. According to Lemma 4.5, the interval $\left( \lambda_{i+1}^{(r)}, 2 \right]$ contains at most $i + 1$ $\gamma_j^{(r)}$, i.e. $\gamma_j^{(r)} \in \left[ -2, \lambda_{i+1}^{(r)} \right]$ for $i + 2 \leq j \leq k_r$ or $\lambda_{i+1}^{(r)} \geq \gamma_{i+2}^{(r)}$. Finally, there remain $\lambda_{k_r}^{(r)}$ and $\gamma_1^{(r)}$ with $\lambda_{k_r}^{(r)} < \gamma_1^{(r)}$.  □

Theorem 4.7 yields a justification for the ordering (3.11) of the $\lambda_i^{(r)}$ and $\gamma_j^{(r)}$ which is only applied in (3.9). This will enable us to develop easily computable estimates for the width in (3.9) (see Lemma 4.9) which will be integrated into an estimate for the overall width of *IBUD* for nonoptimal $q$. In the noninterval case the terms $\left| \lambda_i^{(r)} - \gamma_{j(i)}^{(r)} \right|$ are minimized in order to minimize the roundoff error [13]. In the interval case, an optimal enclosure requires $\lambda_i^{(r)} \geq \gamma_{j(i)}^{(r)}$ which cannot hold for all $i$. (3.11) can at least ensure that only one change of the sign in the sequence of the $\lambda_i^{(r)} - \gamma_{j(i)}^{(r)}$, i.e. only one inequality $\lambda_i^{(r)} < \gamma_{j(i)}^{(r)}$ occurs. We do not claim that this strategy is the best one as principally any ordering is admissible. As an example we mention any minimization of the terms $\left| \lambda_i^{(r)} - \gamma_{j(i)}^{(r)} \right|$. But the proof of Lemma 4.9 will illustrate that the complexity of estimates for the width for orderings of this type seems not to be reasonable.

After having defined the final form of (3.9) by Theorem 4.7, we are able to estimate the additional width of the enclosure vector $\mathbf{x} = LES(\mathbf{M}, \mathbf{b})$ in the case of nonoptimal enclosures. In order to achieve this task, we determine an enclosure

$$\mathbf{x} \subseteq \tilde{\mathbf{x}} + \hat{\mathbf{x}} \tag{4.4}$$

where $\tilde{\mathbf{x}} \equiv \mathbf{x}_{opt}$ denotes the optimal enclosure according to Theorem 3.1 and where $\hat{\mathbf{x}}$ denotes an enclosure of the additional width occurring for nonoptimal $q$. Then we derive an estimate for $\mathrm{d}(\hat{\mathbf{x}})$ depending on $\mathrm{d}(\tilde{\mathbf{x}})$ in a suitable norm.

The larger enclosures are caused by the application of (3.9) in (3.6c). Therefore, we first derive an estimate of the form (4.4) for (3.9). In the sequel, we denote all subvectors leading to $\tilde{\mathbf{x}}$ accordingly by "˜" and the corresponding enclosures of the additional width by "ˆ". We note some

conditions which are needed in the subsequent Lemmata and Theorems:

> *(a)–(f) of Theorem 3.1*
> *(g)*   $\forall A \in \mathbf{A} : \text{spec}(A) = (\lambda_i)_{i=1}^p, \text{spec}(U) = (\omega_i)_{i=1}^p \subset \mathbb{R}$
>      *and* $\forall i \in \{1, \ldots, p\} : \lambda_i \geq 2\omega_i > 0$         (4.5)
> *(h)*   $\forall A \in \mathbf{A} :$   *all eigenvectors of A are eigenvectors of* S, T, *and* U
> *(i)*   $\forall A \in \mathbf{A} :$ A *is irreducible or* A *has an eigenvector* c $> 0$.

The conditions (g) and (h) enable us to treat nonsymmetric matrices like in example (iii) after Theorem 3.1. Obviously, (g) and (h) are satisfied for symmetric positive definite and commuting matrices A, S, T where $S = T$.

Before we start the discussion of the behaviour of *IBUD* with respect to the occurrence of (3.9), we collect in the following lemma some auxiliary results which are needed in later estimates.

**Lemma 4.8.** *Assume (3.1)–(3.11), (4.5), $A \in \mathbf{A}$. Define $\theta \equiv \theta(\lambda, \omega) := \text{arcosh}\big(\lambda/(2\omega)\big)$ for $\lambda \geq 2\omega$,*

$$\text{sh}(m, n, \theta) := \frac{\sinh\{m\theta\}}{\sinh\{n\theta\}} \text{ for } \theta > 0, \ \text{sh}(m, n, 0) := 1 \text{ and } \lambda_{min} := 1/\rho(A^{-1}),$$
$$\omega_{max} := \rho(U), \ s := \rho(S), \ t := \rho(T), \theta_i := \theta(\lambda_i, \omega_i), \ \theta_{min} := \theta(\lambda_{min}, \omega_{max}).$$

*Then*

(a) $(A - \alpha U)^{-1} \geq O$ *for all roots $\alpha$ in (3.3)*

$$S^{(r)}, T^{(r)}, U^{(r)}, \big(A^{(r)}\big)^{-1}, \big(B^{(r)}\big)^{-1}C^{(r)}, \big(A - \gamma_1^{(r)}U\big)^{-1}U\big(B^{(r)}\big)^{-1}C^{(r)}A^{(r)} \geq O$$

(b) *All eigenvectors of A are eigenvectors of* $A^{(r)}, B^{(r)}, C^{(r)}, S^{(r)}, T^{(r)}, U^{(r)}$

(c) $A^{(r)}, B^{(r)}, C^{(r)}, S^{(r)}, T^{(r)}, U^{(r)}$ *commute*

(d) *For $0 \leq m < n : \text{sh}(m, n, \theta) \downarrow 0 \ (\theta \uparrow \infty)$ and $\text{sh}(m, n, \theta) \leq 1$*

(e) $\big|\lambda_{k_r}^{(r)} - \gamma_1^{(r)}\big| \leq 4$

(f) $a_\gamma^{(r)} := \rho\Big(\big(A - \gamma_1^{(r)}U\big)^{-1}U\Big) = \dfrac{1}{(\lambda_{min}/\omega_{max}) - \gamma_1^{(r)}}$

(g) $b_S^{(r)} := \rho\left(\left(B^{(r)}\right)^{-1}C^{(r)}S^{(r)}\right)$

$\qquad = \left(\dfrac{s}{\omega_{max}}\right)^{2^r} \text{sh}\{l_r + 1, k_r + 1, \theta_{min}\} \leq \left(\dfrac{s}{\omega_{max}}\right)^{2^r}$

$\quad b_T^{(r)} := \rho\left(\left(B^{(r)}\right)^{-1}C^{(r)}T^{(r)}\right)$

$\qquad = \left(\dfrac{t}{\omega_{max}}\right)^{2^r} \text{sh}\{l_r + 1, k_r + 1, \theta_{min}\} \leq \left(\dfrac{t}{\omega_{max}}\right)^{2^r}$

$\quad b_{TU}^{(r)} := \rho\left(\left(B^{(r)}\right)^{-1}C^{(r)}T^{(r-1)}U^{(r-1)}\right)$

$\qquad = \left(\dfrac{t}{\omega_{max}}\right)^{2^{r-1}} \text{sh}\{l_r + 1, k_r + 1, \theta_{min}\} \leq \left(\dfrac{t}{\omega_{max}}\right)^{2^{r-1}}$

$\quad b_U^{(r)} := \rho\left(\left(B^{(r)}\right)^{-1}C^{(r)}U^{(r)}\right) = \text{sh}\{l_r + 1, k_r + 1, \theta_{min}\} \leq 1$

(h) $a_S^{(r)} := \rho\left(\left(A^{(r)}\right)^{-1}S^{(r)}\right) = \left(\dfrac{s}{\omega_{max}}\right)^{2^r} \dfrac{1}{2\cosh\{2^r\theta_{min}\}} \leq \dfrac{1}{2}\left(\dfrac{s}{\omega_{max}}\right)^{2^r}$

$\quad a_T^{(r)} := \rho\left(\left(A^{(r)}\right)^{-1}T^{(r)}\right) = \left(\dfrac{t}{\omega_{max}}\right)^{2^r} \dfrac{1}{2\cosh\{2^r\theta_{min}\}} \leq \dfrac{1}{2}\left(\dfrac{t}{\omega_{max}}\right)^{2^r}$

$\quad a_U^{(r)} := \rho\left(\left(A^{(r)}\right)^{-1}U^{(r)}\right) = \dfrac{1}{2\cosh\{2^r\theta_{min}\}} \leq \dfrac{1}{2}$

$\quad x^{(r)} := \rho\left(X^{(r)}\right) = \dfrac{1}{\prod\limits_{i=0}^{r}(2\omega^{2^r}\cosh\{2^r\theta_{min}\})} \leq \dfrac{1}{2^{r+1}}$

(i) $b_{\gamma A}^{(r)} := \rho\left(\left(A - \gamma_1^{(r)}U\right)^{-1}U\left(B^{(r)}\right)^{-1}C^{(r)}A^{(r)}\right)$

$\qquad = b_A^{(r)} \dfrac{\omega_{max}}{\lambda_{min} - \gamma_1^{(r)}\omega_{max}} \leq \dfrac{1}{\lambda_{min}/\omega_{max} - \gamma_1^{(r)}}$

$\qquad$ with $b_A^{(r)} := 1 - \text{sh}\{2^r - (l^r + 1), 2^r + (l^r + 1), \theta_{min}\}$

(j) $w^{(r)} := \sum\limits_{k=1}^{2^r}\{s^{2^r-1+k}t^{2^r-k} + s^{2^r-k}t^{2^r-1+k}\} \leq 2^{r+1}\max\{s,t\}^{2^{r+1}-1}$

(k) $a_\gamma^{(r)}, a_S^{(r)}, a_T^{(r)}, a_U^{(r)}, b_{\gamma A}^{(r)}, b_S^{(r)}, b_T^{(r)}, b_U^{(r)}, b_{TU}^{(r)}, x^{(r)} \downarrow 0\ (\lambda \uparrow \infty).$

*Proof.* We first note that each $A \in \mathbf{A}$ is an M matrix as $\mathbf{M}$ is an interval M matrix. Therefore, $A^{-1} \geq O$ and with (4.5, g) $\lambda_{min} = \min\{\lambda | \lambda \in \mathrm{spec}(A)\} = 1/\rho(A^{-1})$. In (a) $S^{(r)}, T^{(r)}, U^{(r)} \geq O$ follows from the definition of the corresponding matrices. (4.5, f) and the definition of the roots in (3.3) imply $A - \alpha_1^{(r)}U \leq A - \alpha U \leq A + \alpha_1^{(r)}U$. Then $A - \alpha U$ is also an M matrix and $O \leq (A + \alpha_1^{(r)}U)^{-1} \leq (A - \alpha U)^{-1} \leq (A - \alpha_1^{(r)}U)^{-1}$. This also implies $(A^{(r)})^{-1} \geq O$. We further get

$$
(A - \gamma_1^{(r)}U)^{-1}U(B^{(r)})^{-1}C^{(r)}A^{(r)} = \prod_{i=1}^{k_r-1}(A - \gamma_{i+1}^{(r)}U)^{-1}\prod_{i=1}^{k_r}(A - \lambda_i^{(r)}U)^{-1}
$$
$$
= \prod_{i=1}^{k_r-1}\{I + (\lambda_i^{(r)} - \gamma_{i+1}^{(r)})U(A - \lambda_i^{(r)}U)^{-1}\}(A - \lambda_{k_r}^{(r)}U)^{-1} \geq O
$$

$$(4.6)$$

and, with (3.12),

$$
(B^{(r)})^{-1}C^{(r)} = \prod_{i=1}^{l_r}(A - \mu_i^{(r)}U)\prod_{i=1}^{k_r}(A - \lambda_i^{(r)}U)^{-1}
$$
$$
= \prod_{i=1}^{l_r}\{I + (\lambda_i^{(r)} - \mu_i^{(r)})U(A - \lambda_i^{(r)}U)^{-1}\}\prod_{i=l_r+1}^{k_r}(A - \lambda_i^{(r)}U)^{-1} \geq O.
$$

$$(4.7)$$

Parts (b), (c) follow trivially from the definition of the above matrices. For the proof of (d), we recall that with the substitution $\theta := \mathrm{arcosh}(x)$, $x := \lambda/(2\omega)$, $x \geq 1$, we can rewrite sh with Chebyshev polynomials of the second kind

$$
\sinh(m\theta) = \frac{\sinh\{m\theta\}}{\sinh\{\theta\}} = \prod_{i=1}^{m-1}(2x - \lambda_i^{(m)}) = U_{m-1}(x), \quad \lambda_i^{(m)} = 2\cos\frac{i\pi}{m}
$$

as

$$
\mathrm{sh}(m, n, \theta) = \frac{\sinh\{m\theta\}}{\sinh\{n\theta\}} = \frac{\displaystyle\prod_{i=1}^{m-1}(2x - \lambda_i^{(m)})}{\displaystyle\prod_{i=1}^{n-1}(2x - \lambda_i^{(n)})} = \frac{U_{m-1}(x)}{U_{n-1}(x)}.
$$

Note that $\lambda_i^{(m)} < 2$ and $\lambda \geq 2\omega$ imply $2x > \lambda_i^{(m)}$. Differentiation with respect to $x$ yields

$$\frac{d}{dx}\left(\frac{\mathrm{U}_{m-1}(x)}{\mathrm{U}_{n-1}(x)}\right) = 2\frac{\prod\limits_{i=1}^{m-1}\left(2x - \lambda_i^{(m)}\right)}{\prod\limits_{i=1}^{n-1}\left(2x - \lambda_i^{(n)}\right)}\left(\sum_{i=1}^{m-1}\frac{1}{\left(2x - \lambda_i^{(m)}\right)} - \sum_{i=1}^{n-1}\frac{1}{\left(2x - \lambda_i^{(n)}\right)}\right) < 0.$$

Because of $\theta \uparrow \infty \Leftrightarrow x \uparrow \infty$, the monotone convergence to 0 follows from $m < n$, in particular from $\lambda_i^{(m)} < \lambda_i^{(n)}$. (e) follows from the definition (3.3) of $\gamma_1^{(r)}$, $\lambda_{k_r}^{(r)}$. As A and U have the same eigenvectors, $\left(\mathrm{A} - \gamma_1^{(r)}\mathrm{U}\right)^{-1}\mathrm{U}$ has the eigenvalues $\omega_i/\left(\lambda_i - \gamma_1^{(r)}\omega_i\right)$, where $\lambda_i - \gamma_1^{(r)}\omega_i > \lambda_i - 2\omega_i \geq 0$. Because of $\left(\mathrm{A} - \gamma_1^{(r)}\mathrm{U}\right)^{-1}\mathrm{U} \geq \mathrm{O}$, the theorem of Perron-Frobenius ([14], Th. 2.1) shows that $a_\gamma^{(r)}$ is ab eigenvalue, where $\lambda_i = \lambda_{min} = 1/\rho(\mathrm{A}^{-1})$, $\omega_i = \omega_{max} = \rho(\mathrm{U})$, hence (f) holds.

According to the factorization (3.3), $\mathrm{A}^{(r)}$, $\mathrm{B}^{(r)}$, $\mathrm{C}^{(r)}$ have the nonnegative eigenvalues $\omega_i^{2^r}\cosh\{2^r\theta_i\}$, $\omega_i^{k_r}\sinh\{(k_r+1)\theta_i\}/\sinh\{\theta_i\}$, $\omega_i^{l_r}\sinh\{(l_r+1)\theta_i\}/\sinh\{\theta_i\}$ $(1 \leq i \leq p)$, resp., ([5, 13]), whose maximum is reached with (d) for $\theta_{min} = \theta(\lambda_{min}, \omega_{max})$. With part (b), we get

$$\rho\left(\left(\mathrm{B}^{(r)}\right)^{-1}\mathrm{C}^{(r)}\mathrm{S}^{(r)}\right) = \rho\left(\left(\mathrm{B}^{(r)}\right)^{-1}\mathrm{C}^{(r)}\right)\rho\left(\mathrm{S}^{(r)}\right) = s^{2^r}\omega_{max}^{l_r-k_r}\mathrm{sh}\{l_r+1, k_r+1, \theta_{min}\}.$$

The other relations in (g) can be shown similarly. (h) follows from $\cosh\{2^r\theta_i\} \geq 1$. According to (a) and to the ([14], 2.1), $\rho\left(\left(\mathrm{A} - \right.\right.$

$\gamma_1^{(r)} U)^{-1} U (B^{(r)})^{-1} C^{(r)} A^{(r)} \Big)$ is an eigenvalue, hence

$$
\rho\Big( \big(A - \gamma_1^{(r)} U\big)^{-1} U \big(B^{(r)}\big)^{-1} C^{(r)} A^{(r)} \Big)
$$

$$
= \prod_{i=1}^{k_r - 1} \left\{ 1 + \big(\lambda_i^{(r)} - \gamma_{i+1}^{(r)}\big) \frac{\omega_{max}}{\lambda_{min} - \lambda_i^{(r)} \omega_{max}} \right\} \frac{\omega_{max}}{\lambda_{min} - \lambda_{k_r}^{(r)} \omega_{max}}
$$

$$
= \frac{\omega_{max}}{\lambda_{min} - \gamma_1^{(r)} \omega_{max}} \prod_{i=1}^{k_r} \frac{\lambda_{min} - \gamma_i^{(r)} \omega_{max}}{\lambda_{min} - \lambda_i^{(r)} \omega_{max}} \tag{4.8}
$$

$$
= \big(1 - \mathrm{sh}\{2^r - (l^r + 1), 2^r + (l^r + 1), \theta_{min}\}\big) \frac{\omega_{max}}{\lambda_{min} - \gamma_1^{(r)} \omega_{max}}
$$

$$
\leq \frac{1}{(\lambda_{min}/\omega_{max}) - \gamma_1^{(r)}}
$$

$\big($compare (4.6)$\big)$. Assertion (j) is trivial. The monotonicity in $\lambda \equiv \lambda_{min}$ follows from (d). For the matrix in (4.8), it can be directly deduced from the expression after the first equality. $\qquad\square$

After these preparations, we start with the estimation of the additional width caused by one occurrence of (3.9).

**Lemma 4.9.** *Assume (3.1)–(3.11) and (4.5, a–f). Let $\tilde{\mathbf{z}}$ be the vector in (3.9) corresponding to an optimal enclosure for the global system $(\mathbf{M}, \mathbf{b})$ in one of the cases (a) $\mathbf{M} \equiv M$ and arbitrary $\mathbf{b}$, (b) $0 \in \mathbf{b}$, (c) $\mathrm{i}(\mathbf{b}) \geq o$, or (d) $\mathrm{s}(\mathbf{b}) \leq o$. Then the vector $\mathbf{z}$ computed by (3.9) in IBUD satisfies*

(a) $\quad \mathbf{z} = \tilde{\mathbf{z}} + \big|\lambda_{k_r}^{(r)} - \gamma_1^{(r)}\big| \big(A - \gamma_i^{(r)} U\big)^{-1} U \big(B_d^{(r)}\big)^{-1} C_d^{(r)} A_d^{(r)} \mathrm{d}(\mathbf{y})[-1, 1]$

(b) $\quad \mathbf{z} = \tilde{\mathbf{z}} + \big|\lambda_{k_r}^{(r)} - \gamma_1^{(r)}\big| \big(\mathrm{i}(\mathbf{A}) - \gamma_i^{(r)} U\big)^{-1} U \big(B_{inf}^{(r)}\big)^{-1} C_{inf}^{(r)} A_{inf}^{(r)} \mathrm{d}(\mathbf{y})[-1, 1]$

(c) $\quad \begin{aligned} \mathbf{z} = \tilde{\mathbf{z}} + \big|\lambda_{k_r}^{(r)} - \gamma_1^{(r)}\big| \Big( & \big(\mathrm{i}(\mathbf{A}) - \gamma_i^{(r)} U\big)^{-1} U \big(B_{inf}^{(r)}\big)^{-1} C_{inf}^{(r)} A_{inf}^{(r)} \mathrm{s}(\mathbf{y}) \\ & - \big(\mathrm{s}(\mathbf{A}) - \gamma_i^{(r)} U\big)^{-1} U \big(B_{sup}^{(r)}\big)^{-1} C_{sup}^{(r)} A_{sup}^{(r)} \mathrm{i}(\mathbf{y}) \Big)[-1, 1] \end{aligned}$

(d) $\quad$ *like (c) with $\mathrm{i}(\mathbf{A})$ and $\mathrm{s}(\mathbf{A})$ exchanged*

*where $A_{inf}^{(r)} := \prod_{i=1}^{2^r} \big(\mathrm{i}(\mathbf{A}) - \alpha_i^{(r)} U\big)$, $A_{sup}^{(r)} := \prod_{i=1}^{2^r} \big(\mathrm{s}(\mathbf{A}) - \alpha_i^{(r)} U\big)$, $A_d^{(r)} :=$*

*$\prod_{i=1}^{2^r} \big(A - \alpha_i^{(r)} U\big)$, and similarly defined $B_d^{(r)}$, $C_d^{(r)}$, $B_{inf}^{(r)}$, $C_{inf}^{(r)}$, $B_{sup}^{(r)}$, $C_{sup}^{(r)}$.*

*Proof.* (3.9) is characterized by the treatment of a series of systems of the form

$$\mathbf{z} = \mathbf{y} + \alpha U \, LES(\mathbf{B}, \mathbf{y}) \tag{4.9}$$

with, according to (3.10) and (4.5, f), interval M matrices $\mathbf{B}$ and $U \geq O$. The applicability of *LES* follows from (4.5, c, d). *LES* yields optimal enclosures according to (4.5, e), i.e. (2.14).

Assume one of cases (a)–(d) for the main system $(\mathbf{M}, \mathbf{b})$. With $S^{(r)}, T^{(r)}, U^{(r)} \geq O$ and (4.5, e), it can be shown by induction that in any partial system $LES(\mathbf{A} - \alpha U, \mathbf{w})$ in *IBUD* before the first occurrence of (3.6c)/(3.9) (including the latter) the same case applies. For (a) this is obvious. For (b) this follows with (2.7) from the relations $o \in \mathbf{b} \Rightarrow o \in \mathbf{w} \Rightarrow o \in LES(\mathbf{A} - \alpha U, \mathbf{w}) \Rightarrow o \in \mathbf{y}$ in (3.9). (c), (d) are treated similarly with (2.8).

If (a) holds, we obtain for $\alpha \geq 0$ an optimal enclosure $\mathbf{z} = \tilde{\mathbf{z}}$ and with the abbreviation $\mathbf{B} = \mathbf{A} - \alpha U$ and $B = A - \alpha U$, resp.,

$$
\begin{aligned}
i(\tilde{\mathbf{z}}) &= i(\mathbf{y}) + \alpha U B^{-1} i(\mathbf{y}) = (I + \alpha U B^{-1}) i(\mathbf{y}) \\
s(\tilde{\mathbf{z}}) &= s(\mathbf{y}) + \alpha U B^{-1} s(\mathbf{y}) = (I + \alpha U B^{-1}) s(\mathbf{y}) \\
\tilde{\mathbf{z}} &= (I + \alpha U B^{-1}) \mathbf{y} \\
d(\tilde{\mathbf{z}}) &= (I + \alpha U B^{-1}) d(\mathbf{y}).
\end{aligned} \tag{4.10}
$$

For $\alpha < 0$ an optimal enclosure still has the form (4.10), but *IBUD* yields

$$
\begin{aligned}
i(\mathbf{z}) &= i(\mathbf{y}) - |\alpha| U B^{-1} s(\mathbf{y}) \\
s(\mathbf{z}) &= s(\mathbf{y}) - |\alpha| U B^{-1} i(\mathbf{y}) \\
\mathbf{z} &= \mathbf{y} - |\alpha| U B^{-1} \mathbf{y} \\
d(\mathbf{z}) &= (I + |\alpha| U B^{-1}) d(\mathbf{y}).
\end{aligned} \tag{4.11}
$$

The comparison of (4.10) and (4.11) yields

$$
\begin{aligned}
i(\mathbf{z}) &= i(\tilde{\mathbf{z}}) - |\alpha| U B^{-1} d(\mathbf{y}) \\
s(\mathbf{z}) &= s(\tilde{\mathbf{z}}) + |\alpha| U B^{-1} d(\mathbf{y}) \\
d(\mathbf{z}) &= d(\tilde{\mathbf{z}}) + 2|\alpha| U B^{-1} d(\mathbf{y}).
\end{aligned} \tag{4.12}
$$

For (b), we replace $B$ by $i(\mathbf{B})$ in (4.10)–(4.12), for (c) $\big($similarly for (d)$\big)$, we get

$$
\begin{aligned}
i(\tilde{\mathbf{z}}) &= i(\mathbf{y}) - |\alpha| U \, i(\mathbf{B})^{-1} s(\mathbf{y}) \\
s(\tilde{\mathbf{z}}) &= s(\mathbf{y}) - |\alpha| U \, s(\mathbf{B})^{-1} i(\mathbf{y}) \\
d(\tilde{\mathbf{z}}) &= d(\mathbf{y}) + |\alpha| U \, d\big(LES(\mathbf{B}, \mathbf{y})\big),
\end{aligned} \tag{4.13}
$$

$$\begin{aligned}
\mathrm{i}(\mathbf{z}) &= \mathrm{i}(\tilde{\mathbf{z}}) - |\alpha| \mathrm{U}\, \mathrm{d}\big(\mathit{LES}(\mathbf{B},\mathbf{y})\big) \\
\mathrm{s}(\mathbf{z}) &= \mathrm{s}(\tilde{\mathbf{z}}) + |\alpha| \mathrm{U}\, \mathrm{d}\big(\mathit{LES}(\mathbf{B},\mathbf{y})\big) \\
\mathrm{d}(\mathbf{z}) &= \mathrm{d}(\tilde{\mathbf{z}}) + |\alpha| \mathrm{U}\, \mathrm{d}\big(\mathit{LES}(\mathbf{B},\mathbf{y})\big).
\end{aligned} \tag{4.14}$$

The above relations can be summarized in

$$\begin{aligned}
\mathbf{z} &= \tilde{\mathbf{z}} + |\alpha| \mathrm{U}\, \mathrm{d}\big(\mathit{LES}(\mathbf{B},\mathbf{y})\big)[-1,1] && \text{for (c), (d)} \\
\mathbf{z} &= \tilde{\mathbf{z}} + |\alpha| \mathrm{U}\, \mathrm{i}(\mathbf{B})^{-1}\mathrm{d}(\mathbf{y})[-1,1] && \text{for (b)} \\
\mathbf{z} &= \tilde{\mathbf{z}} + |\alpha| \mathrm{U} B^{-1}\mathrm{d}(\mathbf{y})[-1,1] && \text{for (a).}
\end{aligned} \tag{4.15}$$

We first consider (a). For $\alpha = \lambda_i^{(r)} - \gamma_{i+1}^{(r)} \geq 0$, $\mathbf{z} = \mathbf{z}_i$, $\mathbf{y} = \mathbf{z}_{i-1}$, and the interval M matrix $\mathbf{B} = \mathbf{A} - \lambda_i^{(r)}\mathrm{U} \equiv A - \lambda_i^{(r)}\mathrm{U} \equiv B$ (4.10) yields

$$\begin{aligned}
\tilde{\mathbf{z}}_i &= \tilde{\mathbf{z}}_{i-1} + \big(\lambda_i^{(r)} - \gamma_{i+1}^{(r)}\big)\mathrm{U}\big(A - \lambda_i^{(r)}\mathrm{U}\big)^{-1}\tilde{\mathbf{z}}_{i-1} \\
&= \Big\{I + \big(\lambda_i^{(r)} - \gamma_{i+1}^{(r)}\big)\mathrm{U}\big(A - \lambda_i^{(r)}\mathrm{U}\big)^{-1}\Big\}\tilde{\mathbf{z}}_{i-1} \\
&= \Big(\big(A - \gamma_{i+1}^{(r)}\mathrm{U}\big)\big(A - \lambda_i^{(r)}\mathrm{U}\big)^{-1}\Big)\tilde{\mathbf{z}}_{i-1} \\
&= \prod_{j=1}^{i}\Big(\big(A - \gamma_{i+1}^{(r)}\mathrm{U}\big)\big(A - \lambda_j^{(r)}\mathrm{U}\big)^{-1}\Big)\mathbf{y}
\end{aligned} \tag{4.16}$$

$$\mathrm{d}(\tilde{\mathbf{z}}_i) = \prod_{j=1}^{i}\Big(\big(A - \gamma_{j+1}^{(r)}\mathrm{U}\big)\big(A - \lambda_j^{(r)}\mathrm{U}\big)^{-1}\Big)\mathrm{d}(\mathbf{y})$$

for $1 \leq i \leq k_r$ because of (2.4), (2.6) and

$$I + \big(\lambda_i^{(r)} - \gamma_{i+1}^{(r)}\big)\mathrm{U}\big(A - \lambda_i^{(r)}\mathrm{U}\big)^{-1} = \big(A - \lambda_i^{(r)}\mathrm{U}\big)^{-1}\big(A - \gamma_{i+1}^{(r)}\mathrm{U}\big) \geq O \tag{4.17}$$

with optimal enclosures $\mathbf{z}_i = \tilde{\mathbf{z}}_i$ for $1 \leq i \leq k_r - 1$. With $\alpha = \lambda_{k_r}^{(r)} - \gamma_1^{(r)} < 0$, $\mathbf{z} = \mathbf{z}_{k_r}$, $\mathbf{y} = \mathbf{z}_{k_r - 1}$, we further get

$$\mathbf{z}_{k_r} = \tilde{\mathbf{z}}_{k_r} + \big|\lambda_{k_r}^{(r)} - \gamma_1^{(r)}\big|\mathrm{U}\big(A - \lambda_{k_r}^{(r)}\mathrm{U}\big)^{-1}\mathrm{d}(\mathbf{z}_{k_r - 1})[-1,1]. \tag{4.18}$$

With (4.16), (4.18), (2.2), (2.3), and $\mathbf{z}_{k_r - 1} = \tilde{\mathbf{z}}_{k_r - 1}$, the computed enclosure $\mathbf{z} = \mathbf{z}_{k_r}$ satisfies

$$\mathbf{z} = \tilde{\mathbf{z}} + \left(\big|\lambda_{k_r}^{(r)} - \gamma_1^{(r)}\big|\mathrm{U}\prod_{i=2}^{k_r}\big(A - \gamma_i^{(r)}\mathrm{U}\big)\prod_{i=1}^{k_r}\big(A - \lambda_i^{(r)}\mathrm{U}\big)^{-1}\right)\mathrm{d}(\mathbf{y})[-1,1] \tag{4.19}$$

with $\gamma_{k_r+1}^{(r)} := \gamma_1^{(r)}$. With (3.3), this equivalent to the assertion for case (a). The other two cases are treated similarly using (4.15). For case (b), we note with (2.7) $o \in \tilde{\mathbf{z}}_{i-1} \Rightarrow o \in \tilde{\mathbf{z}}_i$ in (4.16). For (c), the repeated application of (2.14, c) yields $o \leq i(\tilde{\mathbf{z}}_{i-1}) \Rightarrow o \leq i(\tilde{\mathbf{z}}_i)$ with (2.4), (2.8) and

$$
\begin{aligned}
\tilde{\mathbf{z}}_i &= \tilde{\mathbf{z}}_{i-1} + \left(\lambda_i^{(r)} - \gamma_{i+1}^{(r)}\right)U\big[\left(s(\mathbf{A}) - \lambda_i^{(r)}U\right)^{-1}i(\tilde{\mathbf{z}}_{i-1}), \\
&\qquad\qquad\qquad \left(i(\mathbf{A}) - \lambda_i^{(r)}U\right)^{-1}s(\tilde{\mathbf{z}}_{i-1})\big] \\
&= \Big[\left\{I + \left(\lambda_i^{(r)} - \gamma_{i+1}^{(r)}\right)U\left(s(\mathbf{A}) - \lambda_i^{(r)}U\right)^{-1}\right\}i(\tilde{\mathbf{z}}_{i-1}), \\
&\qquad \left\{I + \left(\lambda_i^{(r)} - \gamma_{i+1}^{(r)}\right)U\left(i(\mathbf{A}) - \lambda_i^{(r)}U\right)^{-1}\right\}s(\tilde{\mathbf{z}}_{i-1})\Big] \\
&= \Big[\left(s(\mathbf{A}) - \gamma_{i+1}^{(r)}U\right)\left(s(\mathbf{A}) - \lambda_i^{(r)}U\right)^{-1}i(\tilde{\mathbf{z}}_{i-1}), \\
&\qquad \left(i(\mathbf{A}) - \gamma_{i+1}^{(r)}U\right)\left(i(\mathbf{A}) - \lambda_i^{(r)}U\right)^{-1}s(\tilde{\mathbf{z}}_{i-1})\Big] \\
&= \Big[\prod_{j=1}^{i}\left(\left(s(\mathbf{A}) - \gamma_{j+1}^{(r)}U\right)\left(s(\mathbf{A}) - \lambda_j^{(r)}U\right)^{-1}\right)i(\mathbf{y}) \\
&\qquad \prod_{j=1}^{i}\left(\left(i(\mathbf{A}) - \gamma_{j+1}^{(r)}U\right)\left(i(\mathbf{A}) - \lambda_j^{(r)}U\right)^{-1}\right)s(\mathbf{y})\Big].
\end{aligned}
\tag{4.20}
$$

For (d), we apply (2.14, d) exchanging $i(\mathbf{A})$ and $s(\mathbf{A})$ and using $s(\tilde{\mathbf{z}}_{i-1}) \leq o \Rightarrow s(\tilde{\mathbf{z}}_i) \leq o$ and (2.8).                                                           □

Orderings with more than one inequality $\lambda_i^{(r)} < \gamma^{(r)j(i)}$ are possible, but they require a nested application of (4.18) also for $i \neq k_r$. In this case, the assertions of Lemma 4.9 become extremely complex. Lemma 4.9 indicates the additional width in one application of (3.9) with the ordering (3.11). The following theorem describes the propagation of this width across *IBUD*. We define

$$
r_t := \min\{r\,|\,r_s + 1 \leq r \leq r_q, r \text{ odd}\}.
\tag{4.21}
$$

**Theorem 4.10.** *Assume (4.5), (3.1)–(3.11) and (a) $\mathbf{M} \equiv M$, i.e. $\mathbf{A} \equiv A$, $\mathbf{b}$ arbitrary or (b) $o \in \mathbf{b}$ or (c) $o \leq i(\mathbf{b})$ or (d) $s(\mathbf{b}) \leq o$. Then*

*1) $\exists c(q) \geq 0$:*

$$
\max_{1\leq j\leq q}\{\|d(\hat{\mathbf{x}}_j)\|_c\} \leq c(q)\begin{cases} \max\limits_{1\leq j\leq q}\{\|d(\tilde{\mathbf{x}}_j)\|_c\} & a), b) \\[2mm] \max\limits_{1\leq j\leq q}\{\|\,|s(\tilde{\mathbf{x}}_j)|\,\|_c\} & c), d) \end{cases}
$$

$$\text{where } \|x\|_c := \max_{1 \le i \le p} \left\{ \frac{|x_i|}{c_i} \right\}, \; Ac = \rho(A)c \text{ for } (a) \text{ and } i(\mathbf{A})c = \rho\big(i(\mathbf{A})\big)c$$

*otherwise*

2) $c(q) \downarrow 0 \, (d \uparrow \infty)$ for $\mathbf{A}$ replaced by $\mathbf{A} + d\mathbf{I}$.

*If, in addition,*

(j) $q = 2^n \{2^m(2^l + 1) + 1\} - 1$, $n \ge 0$, $m, l > 0$

*then*

3) $r_s = n$, $r_t = n + m$, $r_q = n + m + l$

4) the assertions 1), 2) hold for

$$
\begin{aligned}
c(q) \; := \; & 2 \max \left\{ 1, \max \left\{ \prod_{r=i}^{r_q - 1} \big(a_S^{(r)} + a_T^{(r)}\big) \Big| 0 \le i \le r_q - 1 \right\} \right\} \\
& * \prod_{r=r_t+2}^{r_q} \big(1 + b_U^{(r)}\big) |\lambda_{k_{r_t}}^{(r_t)} - \gamma_1^{(r_t)}| \rho\Big(\big(A - \gamma_1^{(r_t)}U\big)^{-1}U\Big) \\
& * \Big\{ b_U^{(r_t+1)} b_S^{(r_t)} \big(x^{(r_t-1)} w^{(r_t-1)}\big) + b_T^{(r_t)} b_S^{(r_t-1)} \big(x^{(r_t-2)} w^{(r_t-2)}\big) \\
& + b_{TU}^{(r_t+1)} b_A^{(r_t)} \big(1 + a_U^{(r_t)}\big) \prod_{r=r_s+1}^{r_t-1} \big(1 + b_U^{(r)}\big) a_T^{(r_s)} \Big\}
\end{aligned}
$$

where all constants are defined by Lemma 4.8 ($\mathbf{A}$ replaced by $i(\mathbf{A})$ for (b), (c), (d)).

*If, in addition,*

(k) $a_S^{(r)} + a_T^{(r)} \le 1 \quad \forall i \in \{0, \ldots, r_q - 1\}$ and $\rho(S) \le \rho(T)$

*then*

5) $\max\limits_{1 \le j \le q} \{\|d(\hat{\mathbf{x}}_j)\|_c\} = \|d(\tilde{\mathbf{x}}_{2^{r_q}})\|_c$.

*Proof.* In the following proof, we use (2.9)–(2.13) without mentioning it explicitly. Assertion 3) obviously follows from the definition of $q$ in (j) as $q = 2^{n+m+l} + 2^{n+m} + 2^n - 1$. The first deviation from an optimal enclosure occurs for $r = r_t + 1$, i.e. the first step in the reduction phase after a change

from an even $r$ to an odd $r$: The width increases while treating the linear system in the computation of $\mathbf{q}_{j_{r+1}}^{r+1}$ if $q_r$ is odd. More precisely, for $r = r_t$ and $j = j_{r_t+1}$, we note with (3.6c)

$$
\begin{aligned}
\mathbf{q}_j^{r+1} &= \mathrm{U}^{(r+1)}\mathbf{p}_j^{r+1} + \mathrm{S}^{(r)}\mathbf{q}_{j-2^r}^r \\
&\quad + L\tilde{E}S\big(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\mathbf{A}^{(r)}\{\mathrm{U}^{(r+1)}\mathbf{p}_j^{r+1} + \mathrm{T}^{(r)}\mathbf{q}_{j+2^r}^r\}\big) \qquad (22) \\
&= \tilde{\mathbf{q}}_j^{r+1} + \hat{\mathbf{q}}_j^{r+1}
\end{aligned}
$$

$$
\begin{aligned}
\tilde{\mathbf{q}}_j^{r+1} &= \mathrm{U}^{(r+1)}\tilde{\mathbf{p}}_j^{r+1} + \mathrm{S}^{(r)}\tilde{\mathbf{q}}_{j-2^r}^r \\
&\quad + \big(\mathbf{B}^{(r)}\big)^{-1}\mathbf{C}^{(r)}\mathbf{A}^{(r)}\{\mathrm{U}^{(r+1)}\tilde{\mathbf{p}}_j^{r+1} + \mathrm{T}^{(r)}\tilde{\mathbf{q}}_{j+2^r}^r\}
\end{aligned} \qquad (23)
$$

$$
\begin{aligned}
\hat{\mathbf{q}}_j^{r+1} &:= \big|\lambda_{k_r}^{(r)} - \gamma_1^{(r)}\big|\mathrm{U}L\tilde{E}S\Big((\mathbf{A} - \gamma_1^{(r)}\mathrm{U})\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\mathbf{A}^{(r)}\big(\mathrm{U}^{(r+1)}\mathrm{d}(\tilde{\mathbf{p}}_j^{r+1}) \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad + \mathrm{T}^{(r)}\mathrm{d}(\tilde{\mathbf{q}}_{j+2^r}^r)\big)[-1,1]\Big) \\
&= \big|\lambda_{k_r}^{(r)} - \gamma_1^{(r)}\big|\mathrm{U}\big(\mathbf{A} - \gamma_1^{(r)}\mathrm{U}\big)^{-1}\big(\mathbf{B}^{(r)}\big)^{-1}\mathbf{C}^{(r)}\mathbf{A}^{(r)}\big(\mathrm{U}^{(r+1)}\mathrm{d}(\tilde{\mathbf{p}}_j^{r+1}) \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad + \mathrm{T}^{(r)}\mathrm{d}(\tilde{\mathbf{q}}_{j+2^r}^r)\big)[-1,1].
\end{aligned}
$$

For the case (a) which we consider first in order to facilitate the discussion. For the same reason, we further assume even $q_r$ for $r \geq r_t + 1$ which is equivalent to (j) and to the property that (3.6c) is used only once in the whole algorithm (3.5) for *IBUD*. (j) can be expressed by

$$
q_r \text{ is } \begin{cases} \text{even} & \text{for } r_s \leq r \leq r_t - 1 \text{ and } r_t + 1 \leq r \leq r_q \\ \text{odd} & \text{for } r = r_t \text{ and } 0 \leq r \leq r_s - 1. \end{cases} \qquad (4.24)
$$

According to the definition of $r_s$ and $r_t$, $\mathbf{p}_{j_r}^r$, $\mathbf{q}_{j_r}^r$ are computed with (3.6a) for $0 \leq r \leq r_s$, with (3.6b) for $r_s+1 \leq r \leq r_t$, for $r = r_t+1$ with (3.6c) and for $r_t+2 \leq r \leq r_q$ with (3.6b) or (3.6c) in the general case, under restriction (j) only with (3.6b). In the sequel, we extensively use the commutativity of the matrices $\mathrm{A}^{(r)}$, $\mathrm{B}^{(r)}$, $\mathrm{S}^{(r)}$, $\mathrm{T}^{(r)}$, $\mathrm{U}^{(r)}$, and of all matrices derived from them, which results from the commutativity of A, S, and T, hence also of U, without mentioning it explicitly. With (j), we only admit even $q_r$ for $r \geq r_t + 1$. In this case $j_r = j_{r+1}$, hence $j_r - 2^{r+1} + 2^r < j_r = j_{r+1}$, i.e. the additional width introduced in (3.6c) propagates only to $\mathbf{p}_{j_r}^r$ and $\mathbf{q}_{j_r}^r$ for $r \geq r_t + 2$. We then get

for $r = r_t + 2$, $j = j_r$                                            (25)

$$\begin{aligned}
\mathbf{p}_j^r &= \mathbf{p}_j^{r-1} + L\tilde{E}S\big(\mathbf{B}^{(r-1)}, \mathbf{C}^{(r-1)}\{\mathrm{S}^{(r-1)}\mathbf{p}_{j-2^{r-1}}^{r-1} + \mathbf{q}_j^{r-1}\}\big) \\
&= \tilde{\mathbf{p}}_j^{r-1} + L\tilde{E}S\big(\mathbf{B}^{(r-1)}, \mathbf{C}^{(r-1)}\{\mathrm{S}^{(r-1)}\tilde{\mathbf{p}}_{j-2^{r-1}}^{r-1} + (\tilde{\mathbf{q}}_j^{r-1} + \hat{\mathbf{q}}_j^{r-1})\}\big) \\
&\subseteq \tilde{\mathbf{p}}_j^r + \hat{\mathbf{p}}_j^r \\
\tilde{\mathbf{p}}_j^r &:= \tilde{\mathbf{p}}_j^{r-1} + L\tilde{E}S\big(\mathbf{B}^{(r-1)}, \mathbf{C}^{(r-1)}\{\mathrm{S}^{(r-1)}\tilde{\mathbf{p}}_{j-2^{r-1}}^{r-1} + \tilde{\mathbf{q}}_j^{r-1}\}\big) \\
\hat{\mathbf{p}}_j^r &:= L\tilde{E}S\big(\mathbf{B}^{(r-1)}, \mathbf{C}^{(r-1)}\hat{\mathbf{q}}_j^{r-1}\big) = \big(\mathbf{B}^{(r-1)}\big)^{-1}\mathbf{C}^{(r-1)}\hat{\mathbf{q}}_j^{r-1}
\end{aligned}$$

and

$$\text{for } r > r_t + 2, \; j = j_r \tag{26}$$
$$\hat{\mathbf{p}}_j^r := \hat{\mathbf{p}}_j^{r-1} + \big(\mathrm{B}^{(r-1)}\big)^{-1}\mathbf{C}^{(r-1)}\hat{\mathbf{q}}_j^{r-1}$$

as well as

$$\text{for } r \geq r_t + 2, \; j = j_r \tag{27}$$
$$\begin{aligned}
\mathbf{q}_j^r &= \mathrm{S}^{(r-1)}\mathbf{q}_{j-2^{r-1}}^{r-1} + \mathrm{U}^{(r)}\mathbf{p}_j^r \\
&\subseteq \mathrm{S}^{(r-1)}\tilde{\mathbf{q}}_{j-2^{r-1}}^{r-1} + \mathrm{U}^{(r)}(\tilde{\mathbf{p}}_j^r + \hat{\mathbf{p}}_j^r) \subseteq \tilde{\mathbf{q}}_j^r + \hat{\mathbf{q}}_j^r \\
\tilde{\mathbf{q}}_j^r &= \mathrm{S}^{(r-1)}\tilde{\mathbf{q}}_{j-2^{r-1}}^{r-1} + \mathrm{U}^{(r)}\tilde{\mathbf{p}}_j^r \\
\hat{\mathbf{q}}_j^r &:= \mathrm{U}^{(r)}\hat{\mathbf{p}}_j^r.
\end{aligned}$$

Note that $\hat{\mathbf{p}}_{j_r}^r$, $\hat{\mathbf{q}}_{j_r}^r$ are symmetric interval vectors, i.e. (2.14, b) applies for each application of $LES$ in $L\tilde{E}S$ and the resulting vectors are again symmetric.

The solution phase is started with ($r_t < r_q$, otherwise $q$, i.e. the enclosure, would be optimal)

$$r = r_q, \; j = j_r = 2^{r_q} \tag{28}$$
$$\begin{aligned}
\mathbf{x}_j &= \mathbf{p}_j^r + L\tilde{E}S\big(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\mathbf{q}_j^r\big) \\
&\subseteq (\tilde{\mathbf{p}}_j^r + \hat{\mathbf{p}}_j^r) + L\tilde{E}S\big(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}(\tilde{\mathbf{q}}_j^r + \hat{\mathbf{q}}_j^r)\big) \subseteq \tilde{\mathbf{x}}_j + \hat{\mathbf{x}}_j \\
\tilde{\mathbf{x}}_j &= \tilde{\mathbf{p}}_j^r + L\tilde{E}S\big(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\tilde{\mathbf{q}}_j^r\big) \\
\hat{\mathbf{x}}_j &= \hat{\mathbf{p}}_j^r + L\tilde{E}S\big(\mathbf{B}^{(r)}, \mathbf{C}^{(r)}\hat{\mathbf{q}}_j^r\big) = \hat{\mathbf{p}}_j^r + \big(\mathrm{B}^{(r)}\big)^{-1}\mathbf{C}^{(r)}\hat{\mathbf{q}}_j^r.
\end{aligned}$$

In the steps $r = r_q - 1$ to $r = 0$, the additional width propagates to all $\mathbf{x}_j$. To continue the enclosure of this error, we first note that, according to (4.24), $j_r = j_{r+1}$, i.e. $j_{r+1} - 2^{r+1} < j_{r+1}$ for $r \geq r_t + 1$. Therefore, $\mathbf{p}_{j_r}^r$, $\mathbf{q}_{j_r}^r$

are not optimal for $r \geq r_t + 1$ and optimal for $r \leq r_t$. In the solution phase, only for $r = r_q$ they influence the computation of $\mathbf{x}$ under the assumption (4.24), as $\mathbf{x}_{j_r} = \mathbf{x}_{j_{r_q}} = \mathbf{x}_{2^{r_q}}$ has been computed already in the step $r = r_q$. With $\mathbf{x}_0 := \mathrm{o}$, we then obtain

$$r = r_q - 1(-1)0, \; j := 2^r(2^{r+1})j_r - 2^r \tag{29}$$

$$
\begin{aligned}
\mathbf{x}_j &= \mathbf{p}_j^r + L\tilde{E}S\big(\mathbf{A}^{(r)}, \mathrm{S}^{(r)}\mathbf{x}_{j-2^r} + \mathrm{T}^{(r)}\mathbf{x}_{j+2^r} + \mathbf{q}_j^r\big) \\
&\subseteq \tilde{\mathbf{p}}_j^r + L\tilde{E}S\big(\mathbf{A}^{(r)}, \mathrm{S}^{(r)}(\tilde{\mathbf{x}}_{j-2^r} + \hat{\mathbf{x}}_{j-2^r}) + \mathrm{T}^{(r)}(\tilde{\mathbf{x}}_{j+2^r} + \hat{\mathbf{x}}_{j+2^r}) + \tilde{\mathbf{q}}_j^r\big) \\
&\subseteq \tilde{\mathbf{x}}_j + \hat{\mathbf{x}}_j \\
\tilde{\mathbf{x}}_j &= \tilde{\mathbf{p}}_j^r + L\tilde{E}S\big(\mathbf{A}^{(r)}, \mathrm{S}^{(r)}\tilde{\mathbf{x}}_{j-2^r} + \mathrm{T}^{(r)}\tilde{\mathbf{x}}_{j+2^r} + \tilde{\mathbf{q}}_j^r\big) \\
\hat{\mathbf{x}}_j &= L\tilde{E}S\big(\mathbf{A}^{(r)}, \mathrm{S}^{(r)}\hat{\mathbf{x}}_{j-2^r} + \mathrm{T}^{(r)}\hat{\mathbf{x}}_{j+2^r}\big) = \big(\mathbf{A}^{(r)}\big)^{-1}\big(\mathrm{S}^{(r)}\hat{\mathbf{x}}_{j-2^r} + \mathrm{T}^{(r)}\hat{\mathbf{x}}_{j+2^r}\big).
\end{aligned}
$$

Similarly, we get

$$r = r_q - 1(-1)\min\{r_s, r_q - 1\} + 1, \; j = j_r \tag{30}$$

$$
\begin{aligned}
\mathbf{x}_j &= \mathbf{p}_j^r + L\tilde{E}S\big(\mathbf{B}^{(r)}, \mathrm{C}^{(r)}\big\{\mathrm{S}^{(r)}\mathbf{x}_{j-2^r} + \mathbf{q}_j^r\big\}\big) \subseteq \tilde{\mathbf{x}}_j + \hat{\mathbf{x}}_j \\
\tilde{\mathbf{x}}_j &= \tilde{\mathbf{p}}_j^r + L\tilde{E}S\big(\mathbf{B}^{(r)}, \mathrm{C}^{(r)}\big\{\mathrm{S}^{(r)}\tilde{\mathbf{x}}_{j-2^r} + \tilde{\mathbf{q}}_j^r\big\}\big) \\
\hat{\mathbf{x}}_j &:= L\tilde{E}S\big(\mathbf{B}^{(r)}, \mathrm{C}^{(r)}\mathrm{S}^{(r)}\hat{\mathbf{x}}_{j-2^r}\big) = \big(\mathbf{B}^{(r)}\big)^{-1}\mathrm{C}^{(r)}\mathrm{S}^{(r)}\hat{\mathbf{x}}_{j-2^r}
\end{aligned}
$$

with odd $q_r$ $\big($see (3.5)$\big)$. We now repeatedly insert (4.25) in (4.26), (4.27) and vice versa in order to finally obtain from (4.28):

$$\hat{x}_{2^{r_q}} = \prod_{r=r_t+2}^{r_q} \Big(\mathrm{I} + \big(\mathbf{B}^{(r)}\big)^{-1}\mathrm{C}^{(r)}\mathrm{U}^{(r)}\Big)\big(\mathbf{B}^{(r_t+1)}\big)^{-1}\mathrm{C}^{(r_t+1)}\hat{\mathbf{q}}_{j_{r_t+1}}^{r_t+1}. \tag{31}$$

Before inserting (4.23) into (4.31), we first estimate the term $\mathrm{U}^{(r+1)}\mathrm{d}\big(\tilde{\mathbf{p}}_{j_r}^{(r+1)}\big) + \mathrm{T}^{(r)}\mathrm{d}\big(\tilde{\mathbf{q}}_{j_r+2^r}^r\big)$ for $r = r_t$ in (4.23). In case (a), the optimal enclosure is given by $\tilde{\mathbf{x}} = \lfloor \mathrm{M}^{-1}\mathrm{i}(\mathbf{b}), \mathrm{M}^{-1}\mathrm{s}(\mathbf{b})\rfloor$, i.e. $\tilde{\mathbf{x}} = \mathrm{M}^{-1}\tilde{\mathbf{b}}$ with $\tilde{\mathrm{b}} := \mathrm{d}(\mathbf{b})$, $\tilde{\mathrm{x}} := \mathrm{d}(\tilde{\mathbf{x}})$. we then formally apply the interval algorithm (3.5) to the point system $\mathrm{M}\tilde{\mathrm{x}} = \tilde{\mathrm{b}}$ and we get auxiliary vectors $\tilde{\mathrm{p}}_j^r := \mathrm{d}(\tilde{\mathbf{p}}_j^r)$, $\tilde{\mathrm{q}}_j^r := \mathrm{d}(\tilde{\mathbf{q}}_j^r)$. For convenience, we define the abbreviations:

$$\tilde{\mathrm{w}}^{s,j} := \sum_{k=1}^{2^s} \big\{\mathrm{S}^{2^s-1+k}\mathrm{T}^{2^s-k}\mathrm{d}\big(\tilde{\mathbf{x}}_{j-(2k-1)}\big) + \mathrm{S}^{2^s-k}\mathrm{T}^{2^s-1+k}\mathrm{d}\big(\tilde{\mathbf{x}}_{j+(2k-1)}\big)\big\}$$

$$
\begin{aligned}
\mathrm{M}^{(r)} &= \prod_{i=r_s+1}^{r} \left(\mathrm{I} + \left(\mathrm{B}^{(i)}\right)^{-1}\mathrm{C}^{(i)}\mathrm{U}^{(i)}\right)\left(\mathrm{A}^{(r_s)}\right)^{-1}\mathrm{T}^{(r_s)} \\
\mathrm{X}^{(r)} &= \prod_{i=0}^{r}\left(\mathrm{A}^{(i)}\right)^{-1} \\
\tilde{\mathrm{z}}^{r,j} &:= \left(\mathrm{B}^{(r)}\right)^{-1}\mathrm{C}^{(r)}\mathrm{X}^{(r-1)}\mathrm{S}^{(r)}\tilde{\mathrm{w}}^{r-1,j-2^r}.
\end{aligned}
\tag{32}
$$

Using the relations

$$
\begin{aligned}
\mathrm{B}^{(r+1)}\left(\mathrm{C}^{(r+1)}\right)^{-1} &= \mathrm{A}^{(r)}\mathrm{B}^{(r)}\left(\mathrm{C}^{(r)}\right)^{-1} - \mathrm{U}^{(r+1)} \\
\mathrm{B}^{(r)}\left(\mathrm{C}^{(r)}\right)^{-1} &= \mathrm{A}^{(r)} \quad (0 \le r \le r_s)
\end{aligned}
\tag{4.33}
$$

the equations (3.2) and $\tilde{\mathrm{y}}_j^r = \mathrm{A}^{(r)}\tilde{\mathrm{p}}_j^r + \tilde{\mathrm{q}}_j^r$ $(j \ne j_r)$ and $\tilde{\mathrm{y}}_j^r = \mathrm{B}^{(r)}\left(\mathrm{C}^{(r)}\right)^{-1}\tilde{\mathrm{p}}_j^r + \tilde{\mathrm{q}}_j^r$ $(j = j_r)$, we can show by induction that

for $0 \le r \le r_s - 1$, $j := 2^{r+1}(2^{r+1})j_{r+1}$ $\qquad (34)$

$$
\begin{aligned}
\mathrm{d}\left(\tilde{\mathrm{p}}_j^{r+1}\right) &= \mathrm{d}\left(\tilde{\mathrm{x}}_j\right) - \mathrm{X}^{(r)}\tilde{\mathrm{w}}^{r,j} \\
\mathrm{d}\left(\tilde{\mathrm{q}}_j^{r+1}\right) &= -\left\{\mathrm{S}^{(r+1)}\mathrm{d}\left(\tilde{\mathrm{x}}_{j-2^{r+1}}\right) + \mathrm{T}^{(r+1)}\mathrm{d}\left(\tilde{\mathrm{x}}_{j+2^{r+1}}\right)\right\} + \mathrm{A}^{(r+1)}\mathrm{X}^{(r)}\tilde{\mathrm{w}}^{r,j}
\end{aligned}
$$

for $r_s \le r \le r_t - 1$, $j := j_{r+1}$ $\qquad (35)$

$$
\begin{aligned}
\mathrm{d}\left(\tilde{\mathrm{p}}_j^{r+1}\right) &= \mathrm{d}\left(\tilde{\mathrm{x}}_j\right) - \mathrm{M}^{(r)}\mathrm{d}\left(\tilde{\mathrm{x}}_{j_{r_s}+2^{r_s}}\right) - \tilde{\mathrm{z}}^{r,j} \\
\mathrm{d}\left(\tilde{\mathrm{q}}_j^{r+1}\right) &= -\left\{\mathrm{S}^{(r+1)}\mathrm{d}\left(\tilde{\mathrm{x}}_{j-2^{r+1}}\right) + \mathrm{U}^{(r+1)}\mathrm{M}^{(r)}\mathrm{d}\left(\tilde{\mathrm{x}}_{j_{r_s}+2^{r_s}}\right)\right\} \\
&\quad + \mathrm{B}^{(r+1)}\left(\mathrm{C}^{(r+1)}\right)^{-1}\tilde{\mathrm{z}}^{r,j}
\end{aligned}
$$

and finally

for $r = r_t$, $j := j_{r+1}$ $\qquad (36)$

$$
\begin{aligned}
\mathrm{d}\left(\tilde{\mathrm{p}}_j^{r+1}\right) := \mathrm{d}\left(\tilde{\mathrm{x}}_j\right) - \left(\mathrm{A}^{(r)}\right)^{-1}&\Big(\mathrm{T}^{(r)}\mathrm{M}^{(r-1)}\mathrm{d}\left(\tilde{\mathrm{x}}_{j_{r_s}+2^{r_s}}\right) \\
&+ \mathrm{S}^{(r)}\mathrm{X}^{(r-1)}\tilde{\mathrm{w}}^{r-1,j-2^r} + \mathrm{T}^{(r)}\tilde{\mathrm{z}}^{r-1,j+2^r}\Big).
\end{aligned}
$$

With (4.35) and (4.36), we get for $r = r_t + 1$, $j := j_{r+1}$

$$
\begin{aligned}
&\mathrm{U}^{(r_t+1)}\mathrm{d}\left(\tilde{\mathrm{p}}_{j_{r_t}}^{r_t+1}\right) + \mathrm{T}^{(r_t)}\mathrm{d}\left(\tilde{\mathrm{q}}_{j_{r_t}+2^{r_t}}^{r_t}\right) \\
&= \mathrm{U}^{(r_t+1)}\mathrm{d}\left(\tilde{\mathrm{x}}_{j_{r_t}}\right) - \mathrm{U}^{(r_t+1)}\left(\mathrm{A}^{(r_t)}\right)^{-1} \\
&\quad *\left(\mathrm{T}^{(r_t)}\mathrm{M}^{(r_t-1)}\mathrm{d}\left(\tilde{\mathrm{x}}_{j_{r_s}+2^{r_s}}\right) + \mathrm{S}^{(r_t)}\mathrm{X}^{(r_t-1)}\tilde{\mathrm{w}}^{r_t-1,j_{r_t}-2^{r_t}} + \mathrm{T}^{(r_t)}\tilde{\mathrm{z}}^{r_t-1,j_{r_t}+2^{r_t}}\right)
\end{aligned}
\tag{37}
$$

$$- \left\{ S^{(r_t)} T^{(r_t)} d(\tilde{\mathbf{x}}_j) + T^{(r_t)} U^{(r_t)} M^{(r_t-1)} d(\tilde{\mathbf{x}}_{j_{r_s}+2^{r_s}}) \right\}$$
$$+ B^{(r_t)} (C^{(r_t)})^{-1} T^{(r_t)} \tilde{z}^{r_t-1, j_{r_t}+2^{r_t}}$$
$$= -U^{(r_t+1)} S^{(r_t)} X^{(r_t)} \tilde{w}^{r_t-1, j_{r_t}-2^{r_t}}$$
$$- \left( I + (A^{(r_t)})^{-1} U^{(r_t)} \right) T^{(r_t)} U^{(r_t)} M^{(r_t-1)} d(\tilde{\mathbf{x}}_{j_{r_s}+2^{r_s}})$$
$$+ (A^{(r_t)})^{-1} B^{(r_t+1)} (C^{(r_t+1)})^{-1} T^{(r_t)} \tilde{z}^{r_t-1, j+2^{r_t}}.$$

We insert (4.32), (4.37) in (4.23) and then (4.23) in (4.31). As $\hat{x}_{2^{r_q}}$ is a symmetric interval vector, it is sufficient to estimate its width. We use (2.2) and the nonnegativity of all matrices whose spectral radius has been estimated in Lemma 4.8. According to this Lemma, all these spectral radii $\rho(.)$ are positive eigenvalues of the respective matrices and they are all functions of $\lambda_{min}(A) = 1/\rho(A^{-1})$ and $\omega_{max}(U) = \rho(U)$. Then (4.5, i) shows the existence of a vector $c > 0$, which is an eigenvector for all matrices in Lemma 4.8 (f)–(j). From Lemma 2.1, we deduce $\|.\|_{c,c} \leq \rho(.)$ in the **same** matrix norm for all these matrices. This results in the inequality

$$\|d(\hat{x}_{2^{r_q}})\|_c \leq c_4 \max_{1 \leq j \leq q} \left\{ \|d(\tilde{x}_j)\|_c \right\}$$

where
$$
\begin{aligned}
c_4 &= 2 \prod_{r=r_t+2}^{r_q} \left( 1 + \left\| (B^{(r)})^{-1} C^{(r)} U^{(r)} \right\|_{c,c} \right) |\lambda_{k_{r_t}}^{(r_t)} - \gamma_1^{(r_t)}| \\
&\quad * \left( \left\| (A - \gamma_1^{(r_t)} U)^{-1} U \right\|_{c,c} (c_1 + c_2) + c_3 \right) \\
c_1 &:= \left\| (B^{(r_t+1)})^{-1} C^{(r_t+1)} U^{(r_t+1)} \right\|_{c,c} \left\| (B^{(r_t)})^{-1} C^{(r_t)} S^{(r_t)} \right\|_{c,c} \\
&\quad * \left( \left\| X^{(r_t-1)} \right\|_{c,c} v^{(r_t-1)} \right) \\
c_2 &:= \left\| (B^{(r_t)})^{-1} C^{(r_t)} T^{(r_t)} \right\|_{c,c} \left\| (B^{(r_t-1)})^{-1} C^{(r_t-1)} S^{(r_t-1)} \right\|_{c,c} \\
&\quad * \left( \left\| X^{(r_t-2)} \right\|_{c,c} v^{(r_t-2)} \right) \\
c_3 &:= \left\| (B^{(r_t+1)})^{-1} C^{(r_t+1)} T^{(r_t)} U^{(r_t)} \right\|_{c,c} \\
&\quad * \left\| (B^{(r_t)})^{-1} C^{(r_t)} A^{(r_t)} (A - \gamma_1^{(r_t)} U)^{-1} U \right\|_{c,c} \left( 1 + \left\| (A^{(r_t)})^{-1} U^{(r_t)} \right\|_{c,c} \right) \\
&\quad * \prod_{r=r_s+1}^{r_t-1} \left( 1 + \left\| (B^{(r)})^{-1} C^{(r)} U^{(r)} \right\|_{c,c} \right) \left\| (A^{(r_s)})^{-1} T^{(r_s)} \right\|_{c,c} \\
v^{(r)} &:= \sum_{k=1}^{2^r} \left\{ \|S\|_{c,c}^{2^r-1+k} \|T\|_{c,c}^{2^r-k} + \|S\|_{c,c}^{2^r-k} \|T\|_{c,c}^{2^r-1+k} \right\}.
\end{aligned}
$$

(4.38)

Before terminating the proof of assertion 4), i.e. of assertion 1) under restriction (j), we note

for $r := r_q - 1(-1)0$: $\qquad\qquad\qquad\qquad\qquad\qquad$ (4.39a)

$$\max_{j=2^r(2^{r+1})j_r-2^r} \left\{ \|\mathrm{d}(\hat{\mathbf{x}}_j)\|_{\mathrm{c}} \right\} \leq \left( \left\|(\mathrm{A}^{(r)})^{-1}\mathrm{S}^{(r)}\right\|_{\mathrm{c,c}} + \left\|(\mathrm{A}^{(r)})^{-1}\mathrm{T}^{(r)}\right\|_{\mathrm{c,c}} \right)$$

$$* \max_{j=2^{r+1}(2^{r+1})j_r} \left\{ \|\mathrm{d}(\hat{\mathbf{x}}_j)\|_{\mathrm{c}} \right\}$$

$$\leq \left( a_{\mathrm{S}}^{(r)} + a_{\mathrm{T}}^{(r)} \right) \max_{j=2^{r+1}(2^{r+1})j_r} \left\{ \|\mathrm{d}(\hat{\mathbf{x}}_j)\|_{\mathrm{c}} \right\}$$

and, again with Lemma 4.8,

for $r := r_q - 1(-1)0$: $\qquad\qquad\qquad\qquad\qquad\qquad$ (4.39b)

$$\|\mathrm{d}(\hat{\mathbf{x}}_{j_r})\|_{\mathrm{c}} \leq \begin{cases} \left\|(\mathrm{B}^{(r)})^{-1}\mathrm{C}^{(r)}\mathrm{S}^{(r)}\right\|_{\mathrm{c,c}} \|\mathrm{d}(\hat{\mathbf{x}}_{j-2^r})\|_{\mathrm{c}} \\ \quad \big( r := r_q - 1(-1)\min\{r_s, r_q-1\}+1 \big) \\ \left\|(\mathrm{A}^{(r)})^{-1}\mathrm{S}^{(r)}\right\|_{\mathrm{c,c}} \|\mathrm{d}(\hat{\mathbf{x}}_{j-2^r})\|_{\mathrm{c}} \\ \quad \big( r := \min\{r_s, r_q-1\}(-1)0 \big) \end{cases}$$

$$\leq \|\mathrm{d}(\hat{\mathbf{x}}_{j-2^r})\|_{\mathrm{c}}.$$

Note also that according to our assumption (4.24), $q$ is defined by $r_s$, $r_t$, i.e. $c(q) \equiv c(r_s, r_t)$. With

$$c'(q) := c_4 \max \left\{ 1, \max \left\{ \prod_{r=i}^{r_q-1} \left( \left\|(\mathrm{A}^{(r)})^{-1}\mathrm{S}^{(r)}\right\|_{\mathrm{c,c}} + \left\|(\mathrm{A}^{(r)})^{-1}\mathrm{T}^{(r)}\right\|_{\mathrm{c,c}} \right) \right. \right. \tag{4.40}$$

$$\left. \left. \Big| \ 0 \leq i \leq r_q - 1 \right\} \right\}$$

and using Lemma 4.8, we get $c'(q) \leq c(q)$ with $c(q)$ defined in assertion 4). Under restriction (j) the monotone convergence $c(q) \downarrow 0 \ (d \uparrow \infty)$, i.e. assertion 2), immediately follows with Lemma 4.8 using $d \uparrow \infty \Leftrightarrow \theta \uparrow \infty \Leftrightarrow \lambda' := \lambda + d \uparrow \infty$ and the monotonicity in $d$ from the right-hand side in the relation $\left\|(\mathrm{B}^{(r_t)})^{-1}\mathrm{C}^{(r_t)}\mathrm{A}^{(r_t)}(\mathrm{A}-\gamma_1^{(r_t)}\mathrm{U})^{-1}\mathrm{U}\right\|_{\mathrm{c,c}} \leq b_{\mathrm{A}}^{(r_t)} \left\|(\mathrm{A}-\gamma_1^{(r_t)}\mathrm{U})^{-1}\mathrm{U}\right\|_{\mathrm{c,c}}$. With (4.38), we have completed the proof of assertion 4). From (4.39) we immediately deduce assertion 5).

If we remove the restriction (j) in order to prove assertions 1), 2) for general $q$, we are faced with considerably more complicated enclosures. First, (4.23) has to be replaced by

$$\hat{\mathbf{q}}_j^{r+1} := \mathrm{U}^{(r+1)}\hat{\mathbf{p}}_j^{r+1} + |\lambda_{k_r}^{(r)} - \gamma_1^{(r)}|\mathrm{U}\mathit{LES}\Big( (\mathrm{A}-\lambda_{k_r}^{(r)}\mathrm{U})\mathrm{B}^{(r)},$$

$$\mathrm{C}^{(r)}\mathrm{A}^{(r)}\{ \mathrm{U}^{(r+1)}\big(\mathrm{d}(\tilde{\mathbf{p}}_j^{r+1}) + \mathrm{d}(\hat{\mathbf{p}}_j^{r+1})\big) \tag{4.23'}$$

$$+ \mathrm{T}^{(r)}\big(\mathrm{d}(\tilde{\mathbf{q}}_{j+2^r}^r) + \mathrm{d}(\hat{\mathbf{q}}_{j+2^r}^r)\big) \}[-1,1] \Big)$$

also for possibly several $r \geq r_t + 1$ and $j = j_r$. In the subsequent steps of both the reduction and the solution phase, additional terms depending on $\hat{\mathbf{p}}_{j_r}^r$, $\hat{\mathbf{q}}_{j_r}^r$ have to be added. The assertions 1), 2), however, remain valid, as the proof is based on the behaviour of the spectral radii listed in Lemma 4.8 for $d \to \infty$. Because of the unpredictable number of case decision, a useful closed expression like (4.38) cannot be derived and in addition we cannot guarantee anymore that the maximal error in the chosen norm occurs for $\hat{x}_{2^{r_q}}$.

The whole proof can now be repeated for case (b) by applying (2.7) and (2.14, b) for *LES*. In case (c), (d), we only get a weaker estimate. We simplify (c), (d) from Lemma 4.9 with (2.5) to

$$\mathbf{z} \subseteq \tilde{\mathbf{z}} + \left|\lambda_{k_r}^{(r)} - \gamma_1^{(r)}\right| \left(\mathrm{i}(\mathbf{A}) - \gamma_i^{(r)}\mathrm{U}\right)^{-1} \mathrm{U}\left(\mathrm{B}_{inf}^{(r)}\right)^{-1} \mathrm{C}_{inf}^{(r)} \mathrm{A}_{inf}^{(r)} |\mathrm{s}(\mathbf{y})|[-1,1] \quad (4.41)$$

and replace (4.23) and all subsequent estimates appropriately by applying (2.8).                                                                                      $\square$

Condition (k), i.e. assertion 5) is satisfied in many cases (compare Lemma 4.8), as $\cosh(2^r \theta_{min}) \gg 1$ often holds. But for $\cosh(2^r \theta_{min}) \approx 1$, we possibly get $a_\mathrm{S}^{(r)} + a_\mathrm{T}^{(r)} \geq 1$, taking into account that $\frac{\mathrm{S}}{\omega_{max}} = \sqrt{\frac{S}{t}}$ and $\frac{\mathrm{t}}{\omega_{max}} = \sqrt{\frac{t}{S}}$ for $\mathrm{S} \neq \mathrm{T}$, hence $\frac{\mathrm{t}}{\omega_{max}} + \frac{\mathrm{S}}{\omega_{max}} \geq 1$. For $\mathrm{S} = \mathrm{T}$, Theorem 4.10 can be significantly simplified:

**Corollary 4.11.** *If, in addition,*

  *k')* $\mathrm{S} = \mathrm{T}$

*then*

  *4')  the assertions 1), 2) hold for*

$$c(q) := 2 \prod_{r=r_t+2}^{r_q} (1 + b^{(r)}) \left|\lambda_{k_{r_t}}^{(r_t)} - \gamma_1^{(r_t)}\right| \rho\left((\mathrm{A} - \gamma_1^{(r_t)}\mathrm{U})^{-1}\mathrm{U}\right)$$
$$* \left\{ b^{(r_t+1)} b^{(r_t)} 2^{r_t-1} \prod_{i=0}^{r_t-1} a^{(i)} + b^{(r_t)} b^{(r_t-1)} 2^{r_t-2} \prod_{i=0}^{r_t-2} a^{(i)} \right.$$
$$\left. + b^{(r_t+1)} b_\mathrm{A}^{(r_t)} (1 + a^{(r_t)}) \prod_{r=r_s+1}^{r_t-1} (1 + b^{(r)}) a^{(r_s)} \right\}$$

where $b^{(r)} \equiv b_U^{(r)}$, $a^{(r)} \equiv a_U^{(r)}$

5') $\displaystyle\max_{1\le j\le q}\left\{\|\mathrm{d}(\hat{\mathbf{x}}_j)\|_c\right\} = \|\mathrm{d}(\tilde{\mathbf{x}}_{2^{r_q}})\|_c$

*Proof.* In the special case S = T, i.e. $\forall r \in \{0,\ldots,r_q\} : S^{(r)} = T^{(r)} = U^{(r)}$, $c(q)$ can be simplified with $a^{(r)} \equiv a_U^{(r)}$, $b^{(r)} \equiv b_U^{(r)}$, and $x^{(r)}w^{(r)} = 2^r \prod_{i=0}^r a^{(i)}$. We further note (compare Lemma 4.8) $a_S^{(r)} + a_T^{(r)} = 2a_U^{(r)} = \frac{1}{\cosh(2^r\theta_{min})} \le 1$. $\square$

The conditions (4.5, g–i) ensure that estimates for nonsymmetric coefficients matrices like that of the five point discretization of the elliptic operator $a(x)u_{xx} + bu_{yy} + c(x)u_x + du_y$ can be treated. In the symmetric case, the Euclidean norm can be used.

**Corollary 4.12.** *For symmetric A, S, T, the assertions 2), 4) of Theorem 4.10 can be improved:*

2')
$$\|\mathrm{d}(\hat{\mathbf{x}})\|_\infty \le \begin{cases} \begin{aligned} \max_{1\le j\le q}\{\|\mathrm{d}(\hat{\mathbf{x}}_j)\|_2\} &\le c(q)\max_{1\le j\le q}\{\|\mathrm{d}(\tilde{\mathbf{x}}_j)\|_2\} \\ &\le c(q)\sqrt{p}\,\|\mathrm{d}(\tilde{\mathbf{x}})\|_\infty \end{aligned} & a),\ b) \\[2ex] \begin{aligned} \max_{1\le j\le q}\{\|\mathrm{s}(\hat{\mathbf{x}}_j)\|_2\} &\le c(q)\max_{1\le j\le q}\{\||\mathrm{s}(\tilde{\mathbf{x}}_j)|\|_2\} \\ &\le c(q)\sqrt{p}\||\mathrm{s}(\tilde{\mathbf{x}})|\|_\infty \end{aligned} & c),\ d) \end{cases}$$

4') $\displaystyle\max_{1\le j\le q}\{\|\mathrm{d}(\hat{\mathbf{x}}_j)\|_2\} = \|\mathrm{d}(\hat{\mathbf{x}}_{2^{r_q}})\|_2.$

*Proof.* If A, S, and T are symmetric, then all matrices needed in the definition of $c(q)$ commute. This implies their symmetry and $\rho(.) = \|.\|_{2,2}$. $\square$

In the cases (a), (b), we have derived the desired estimate for $\mathrm{d}(\hat{\mathbf{x}})$ depending on $\mathrm{d}(\tilde{\mathbf{x}})$. In the cases (c), (d), we have to accept instead a dependence of $\mathrm{d}(\hat{\mathbf{x}})$ on $|\mathrm{s}(\tilde{\mathbf{x}})|$ which is due to (c), (d) of Lemma 4.9, i.e. basically (2.14, c) and (2.14, d).

**Corollary 4.13.** *With* $C = (\delta_{i,j} c_i)_{i,j=1}^p$, *we get*

$$
\|d(\hat{\mathbf{x}})\|_\infty \leq
\begin{cases}
c(q)\|C\|_\infty \max\limits_{1 \leq j \leq q}\{\|d(\tilde{\mathbf{x}}_j)\|_c\} & \\
\quad \leq c(q)\|C\|_\infty \|C^{-1}\|_\infty \|d(\tilde{\mathbf{x}})\|_\infty & a), b) \\[2ex]
c(q)\|C\|_\infty \max\limits_{1 \leq j \leq q}\{\||s(\tilde{\mathbf{x}}_j)|\|_c\} & \\
\quad \leq c(q)\|C\|_\infty \|C^{-1}\|_\infty \||s(\tilde{\mathbf{x}})|\|_\infty & c), d).
\end{cases}
$$

The corollaries contain simplifications, mainly estimates in the Euclidean and the $\infty$-norm which can be significant for computational purposes. The eigenvector c defining the $\|.\|_c$-norm should also be accessible in many cases. The practical significance of the estimates containing $\|C\|_\infty$, $\|C^{-1}\|_\infty$ depends on the size of these norms in a specific context.

The assertions of Theorem 4.10 are the key for the application of *IBUD* for values of the block dimension $q$ for which optimal enclosures cannot be expected. We have derived an estimate for the additional width of the enclosure vector *IBUD*$(\mathbf{M}, \mathbf{b})$ when compared to an optimal enclosure in the "usual" four cases, in which such an enclosure is principally possible. Assertion 2) shows that *IBUD* can yield satisfactory enclosures for any value of $q$ under the condition that the coefficient matrix $\mathbf{M}$ is sufficiently diagonally dominant.

For a selected set of values of $q$, i.e. those which cause only one occurrence of (3.6c)/(3.9), the width can be explicitly estimated by a rather ugly, but very easily computable constant $c(q)$. It is principally possible to derive an estimate for general $q$. But on the one hand, such an estimate would have to be computed, due to the case decisions which are necessary for general $q$, by an algorithm rather than by a simple constant like $c(q)$ from assertion 5). On the other hand, more general $q$ than those defined by condition (j) are characterized by more than one application of (3.6c). Numerical results will show that for these values of $q$ acceptable enclosures can be observed only for a very (in practice usually unrealistically) strong diagonal dominance of $\mathbf{M}$. Unfortunately, the constants defining $c(q)$ are not equally monotone in $r$: while $a_\gamma^{(r)}$, $a_{\mathrm{T}}^{(r)}$, $a_{\mathrm{U}}^{(r)}$, $b_{\gamma\mathrm{A}}^{(r)}$, $b_{\mathrm{S}}^{(r)}$, $b_{\mathrm{T}}^{(r)}$, $b_{\mathrm{U}}^{(r)}$, $b_{\mathrm{TU}}^{(r)}$, $x^{(r)}$ decrease with increasing $r$, $\left|\lambda_{k_{r_t}}^{(r_t)} - \gamma_1^{(r_t)}\right|$, $\rho\left((A - \gamma_1^{(r_t)}U)^{-1}U\right)$ or the products involving $1 + a_{\mathrm{U}}^{(r)}$, $1 + b_{\mathrm{U}}^{(r)}$ may increase. Therefore, it is not possible to predict suitable $q$ from monotonicity arguments on $r$.

Assertion 4) is only valid for the $q$ defined by (j) because for more general $q$, $\mathbf{x}_j$ in (4.30) can depend also from $\hat{\mathbf{p}}_j^{(r)}$ and $\hat{\mathbf{q}}_j^{(r)}$. This dependence has to be mainly considered in (4.37).

# 5   Numerical results

The numerical examples have been computed on a workstation IBM RS 6000/560. We have used a simulation of an interval arithmetic which follows the principles described in [1]. All bounds are computed by the near-IEEE floating point arithmetic on the IBM workstation. All rounding errors are included in the computed interval vector. Lower and upper bounds are multiplied by factors $1 - 2^{-53}$ and $1 + 2^{-52}$. Underflow results are avoided by substituting them by one of the constants $0$ or $\pm pmin$, where $pmin = 2^{-1022}$ denotes the smallest positive normalized machine number. This simulation represents a compromise between precision and speed. The usually slight overestimates caused by this simulation and indirectly by the underlying floating point arithmetic can be reduced by using an accurate arithmetic like that proposed in [6, 7] and related publications. This has to be paid for by a reduction of speed by roughly at least one order of magnitude. In the following examples the quality of enclosures is indicated by the following relative width

$$rd(x) := \max_{1 \leq i \leq N} \left\{ \frac{\mathrm{d}(\hat{X}_i)}{\mathrm{d}(\tilde{X}_i)} \right\} \tag{5.1}$$

where $\hat{\mathbf{x}}$ denotes an estimate for the difference of $\mathbf{x} = IBUD(\mathbf{M}, \mathbf{b})$ to the optimal enclosure $\tilde{\mathbf{x}} = \mathbf{x}_{opt}$ in case (a) according to $\mathbf{x} \subseteq \tilde{\mathbf{x}} + \hat{\mathbf{x}}$.

Figure 1 illustrates the dependence of the width on the diagonal dominance of the coefficient matrix for the example $\mathbf{M} \equiv \mathrm{M} = (-\mathrm{I}, \mathrm{A} + d\mathrm{I}, -\mathrm{I})$, $\mathrm{A} = (-1, 4, -1)$, $d = 10^k$, $-7 \leq k \leq 3$ or $k = 0$, $\mathbf{x}_{opt} = ([1, 2])_{i=1}^N$, $\mathbf{b} := [\mathrm{M}\,\mathrm{i}(\mathbf{x}_{opt}), \mathrm{M}\,\mathrm{s}(\mathbf{x}_{opt})]$, $127 \leq q \leq 146$, $p = 255$.

As predicted, we observe widths which strongly differ depending on $q$ for small $d$, but which decrease to 0 with increasing diagonal dominance. Table 1 confirms this fact for the same example for some values of $d$ and values $q = 2^n \big( 2^m (2^l + 1) + 1 \big) - 1$ for which optimal enclosures cannot be expected. This table illustrates the fact that $c(q)$ is a good measure
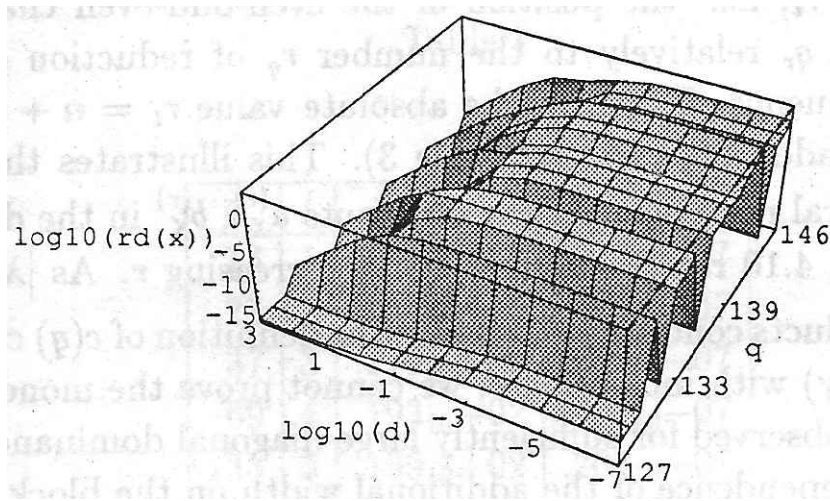
Figure 1

for the relative additional width (5.1). Note that this is not necessarily an upper bound for (5.1) according to its definition in Theorem 4.10 on the basis of maxima over the blocks $1 \leq j \leq q$ in the Euclidean norm. The last column reveals that the computed relative width remains roughly in the order of the relative machine precision while values for $c(q)$ indicate that the theoretical additional width is smaller by orders of magnitude. This is due to unavoidable overestimates of the rounding errors by the simulated interval arithmetic. These overestimates are mainly caused by the fact that in the treatment of the sequences of subsystems with matrices $\mathbf{A} - \alpha \mathbf{U}$, in particular for large $r$, the absolute values of right-hand sides and intermediate enclosure vectors extremely vary over almost all orders of magnitude while the initial r.h.s. and the final enclosures have a reasonable size. It is typical for all variants of the Buneman algorithm, that, depending on the used floating point arithmetic, underflow frequently occurs already for moderate system sizes. As underflow results are replaced by the next admissible normalized machine number, they can be subject to a significant overestimate of the relative error in subsequent steps.

The Tables 2, 3 show the behaviour of $c(q)$, i.e. the additional width, depending on the constants $m$, $n$, $l$ in the definition of $q$ in Theorem 4.10 (j). While $l = r_q - r_t$, i.e. the position of the even-odd-even change in the sequence of the $q_r$ relatively to the number $r_q$ of reduction steps, only has a minor influence (Table 2), the absolute value $r_t = n + m$ strongly determines the additional width (Table 3). This illustrates the fact that almost all spectral radii defining the constants $a_X^{(r)}$, $b_X^{(r)}$ in the definition of $c(q)$ in

| $q$ | $d = 0$ | | $d = 0.1$ | |
|---|---|---|---|---|
| | $c(q)$ | $rd(x)$ | $c(q)$ | $rd(x)$ |
| 130 | $.363_{10}+01$ | $.150_{10}+01$ | $.122_{10}+01$ | $.672_{10}+00$ |
| 132 | $.103_{10}+02$ | $.543_{10}+01$ | $.109_{10}+01$ | $.975_{10}+00$ |
| 133 | $.199_{10}+02$ | $.834_{10}+01$ | $.182_{10}+01$ | $.141_{10}+01$ |
| 136 | $.218_{10}+02$ | $.138_{10}+02$ | $.181_{10}+00$ | $.242_{10}+00$ |
| 137 | $.423_{10}+02$ | $.229_{10}+02$ | $.339_{10}+00$ | $.365_{10}+00$ |
| 139 | $.789_{10}+02$ | $.340_{10}+02$ | $.484_{10}+00$ | $.464_{10}+00$ |
| 144 | $.397_{10}+02$ | $.292_{10}+02$ | $.222_{10}-02$ | $.367_{10}-02$ |
| 145 | $.780_{10}+02$ | $.521_{10}+02$ | $.457_{10}-02$ | $.561_{10}-02$ |
| $q$ | $d = 1$ | | $d = 10$ | |
| | $c(q)$ | $rd(x)$ | $c(q)$ | $rd(x)$ |
| 130 | $.769_{10}-01$ | $.714_{10}-01$ | $.134_{10}-03$ | $.134_{10}-03$ |
| 132 | $.501_{10}-02$ | $.533_{10}-02$ | $.197_{10}-07$ | $.197_{10}-07$ |
| 133 | $.637_{10}-02$ | $.629_{10}-02$ | $.205_{10}-07$ | $.205_{10}-07$ |
| 136 | $.520_{10}-05$ | $.580_{10}-05$ | $.108_{10}-15$ | $.355_{10}-14$ |
| 137 | $.665_{10}-05$ | $.669_{10}-05$ | $.110_{10}-15$ | $.355_{10}-14$ |
| 139 | $.689_{10}-05$ | $.688_{10}-05$ | $.111_{10}-15$ | $.355_{10}-14$ |
| 144 | $.220_{10}-11$ | $.253_{10}-11$ | $.135_{10}-32$ | $.333_{10}-14$ |
| 145 | $.285_{10}-11$ | $.290_{10}-11$ | $.137_{10}-32$ | $.333_{10}-14$ |

Table 1

| $q$ | $l$ | $d = 0$ | $d = 10$ |
|---|---|---|---|
| 13 | 1 | $.123_{10}+02$ | $.205_{10}-07$ |
| 21 | 2 | $.156_{10}+02$ | $.205_{10}-07$ |
| 37 | 3 | $.179_{10}+02$ | $.205_{10}-07$ |
| 69 | 4 | $.193_{10}+02$ | $.205_{10}-07$ |
| 133 | 5 | $.199_{10}+02$ | $.205_{10}-07$ |
| 261 | 6 | $.200_{10}+02$ | $.205_{10}-07$ |
| 517 | 7 | $.200_{10}+02$ | $.205_{10}-07$ |

Table 2: $q = 2^{n+m+l} + 2^{n+m} + 2^n - 1 = 4 * 2^l + 5 \quad (n = m = 1)$

| $q$ | $n$ | $m$ | $n+m$ | $d=0$ | $d=0.1$ |
|---:|:---:|:---:|---:|:---:|:---:|
| 6 | 0 | 1 | 1 | $.216_{10}+01$ | $.114_{10}+01$ |
| 12 | 0 | 2 | 2 | $.674_{10}+01$ | $.108_{10}+01$ |
| 13 | 1 | 1 | 2 | $.123_{10}+02$ | $.181_{10}+01$ |
| 24 | 0 | 3 | 3 | $.156_{10}+02$ | $.181_{10}+00$ |
| 25 | 1 | 2 | 3 | $.294_{10}+02$ | $.339_{10}+00$ |
| 27 | 2 | 1 | 3 | $.521_{10}+02$ | $.484_{10}+00$ |
| 48 | 0 | 4 | 4 | $.313_{10}+02$ | $.222_{10}-02$ |
| 49 | 1 | 3 | 4 | $.610_{10}+02$ | $.457_{10}-02$ |
| 51 | 2 | 2 | 4 | $.115_{10}+03$ | $.697_{10}-02$ |
| 55 | 3 | 1 | 4 | $.201_{10}+03$ | $.781_{10}-02$ |
| 96 | 0 | 5 | 5 | $.532_{10}+02$ | $.175_{10}-06$ |
| 97 | 1 | 4 | 5 | $.105_{10}+03$ | $.382_{10}-06$ |
| 99 | 2 | 3 | 5 | $.204_{10}+03$ | $.604_{10}-06$ |
| 103 | 3 | 2 | 5 | $.382_{10}+03$ | $.690_{10}-06$ |
| 111 | 4 | 1 | 5 | $.661_{10}+03$ | $.694_{10}-06$ |
| 192 | 0 | 6 | 6 | $.567_{10}+02$ | $.578_{10}-15$ |
| 193 | 1 | 5 | 6 | $.113_{10}+03$ | $.131_{10}-14$ |
| 195 | 2 | 4 | 6 | $.225_{10}+03$ | $.211_{10}-14$ |
| 199 | 3 | 3 | 6 | $.439_{10}+03$ | $.244_{10}-14$ |
| 207 | 4 | 2 | 6 | $.819_{10}+03$ | $.246_{10}-14$ |
| 223 | 5 | 1 | 6 | $.137_{10}+04$ | $.246_{10}-14$ |
| 384 | 0 | 7 | 7 | $.199_{10}+02$ | $.339_{10}-32$ |
| 385 | 1 | 6 | 7 | $.404_{10}+02$ | $.781_{10}-32$ |
| 387 | 2 | 5 | 7 | $.828_{10}+02$ | $.127_{10}-31$ |
| 391 | 3 | 4 | 7 | $.171_{10}+03$ | $.148_{10}-31$ |
| 399 | 4 | 3 | 7 | $.347_{10}+03$ | $.149_{10}-31$ |
| 415 | 5 | 2 | 7 | $.655_{10}+03$ | $.149_{10}-31$ |
| 447 | 6 | 1 | 7 | $.103_{10}+04$ | $.149_{10}-31$ |
| 768 | 0 | 8 | 8 | $.108_{10}+01$ | $.609_{10}-67$ |
| 769 | 1 | 7 | 8 | $.227_{10}+01$ | $.142_{10}-66$ |

Table 3: $q = 2^{n+m+1} + 2^{n+m} + 2^n - 1 = 3 * 2^{n+m} + 2^n - 1$   $(l = 1)$

Theorem 4.10 rapidly decrease with increasing $r$. As $\left|\lambda_{k_{r_t}}^{(r_t)} - \gamma_1^{(r_t)}\right|$, $a_{\gamma^A}^{(r)}$ and the products concerning $1 + b_U^{(r)}$ in the definition of $c(q)$ can increase (relatively slowly) with increasing $r$, we cannot prove the monotonicity in $r$ which can be observed for sufficiently large diagonal dominance. Table 4 illustrates the dependence of the additional width on the block size $p$.

| $p$ | $c(q)$ |
|---:|:---|
| 7 | $.218_{10}-06$ |
| 15 | $.499_{10}-05$ |
| 31 | $.122_{10}-04$ |
| 63 | $.155_{10}-04$ |
| 127 | $.164_{10}-04$ |
| 255 | $.167_{10}-04$ |
| 511 | $.167_{10}-04$ |

Table 4: $q = 151 \quad p = 2^{k+1} - 1 \quad (2 \le k \le 8)$

| $q$ | $k = -3$ | | $k = -1$ | |
|---|---|---|---|---|
| | $rd(x)$ | $c(q)$ | $rd(x)$ | $c(q)$ |
| 130 | $.144_{10}+00$ | $.167_{10}+00$ | $.148_{10}+00$ | $.170_{10}+00$ |
| 132 | $.225_{10}-01$ | $.220_{10}-01$ | $.264_{10}-01$ | $.232_{10}-01$ |
| 133 | $.329_{10}-01$ | $.301_{10}-01$ | $.340_{10}-01$ | $.318_{10}-01$ |
| 136 | $.124_{10}-03$ | $.957_{10}-04$ | $.133_{10}-03$ | $.108_{10}-03$ |
| 137 | $.157_{10}-03$ | $.134_{10}-03$ | $.168_{10}-03$ | $.151_{10}-03$ |
| 139 | $.167_{10}-03$ | $.143_{10}-03$ | $.179_{10}-03$ | $.161_{10}-03$ |
| 144 | $.995_{10}-09$ | $.714_{10}-09$ | $.115_{10}-08$ | $.921_{10}-09$ |
| 145 | $.125_{10}-08$ | $.102_{10}-08$ | $.144_{10}-08$ | $.131_{10}-08$ |

Table 5

Table 5 shows similar results for a non-symmetric example where the non-symmetry is determined by the "mesh size" $h$ in accordance with typical discretizations of general elliptic problems which are not given in self-adjoint

form. We consider the example

$$
\begin{aligned}
\mathbf{M} &= (-\mathrm{S}, \mathbf{A}, -\mathrm{T}), \\
\mathrm{S} &= -\big(b - (h/2)e\big)\mathrm{I}, \\
\mathrm{T} &= -\big(b + (h/2)e\big)\mathrm{I}, \\
\mathbf{A} &= \Big(-\big(a - (h/2)c\big), 2(a+b), -\big(a + (h/2)c\big)\Big) + [0,1]\mathrm{I} + d\mathrm{I}, \\
a &= 4.1, \quad b = 3.1, \quad c = 1.0, \quad e = 1.1, \quad d = 2, \\
\mathbf{x} &= \big([i+j, 2(i+j)]\big)_{i=1\,j=1}^{p\quad q}, \quad p = 255, \quad h = 10^{-k}, \quad k = -3, -1.
\end{aligned}
$$

# 6   Conclusion

The interval Buneman algorithm can be defined, like its noninterval counterpart, for arbitrary values of the block dimension. This increase of flexibility is useful in applications where the system size resulting from the requirement $q = 2^{k+1} - 1$ in the variants of the original Buneman algorithm is not adequate for various reasons. Under appropriate conditions *IBUD* yields optimal enclosures for selected values of the block dimension $q$. In the present paper, we have extended the range of suitable values by admitting also nonoptimal enclosures. For the special choice $q = 2^n\big(2^m(2^l+1)+1\big) - 1$, we can prove suboptimal enclosures already in the presence of a moderate degree of diagonal dominance of the coefficient matrix. The enclosure quality can be controlled by fairly easily computable estimates.

**Acknowledgment**. I am indebted to the referees for their helpful comments and suggestions.

# References

[1] Alefeld, G. and Herzberger, J. *Introduction to interval computations.* Academic Press, New York, 1983.

[2] Alefeld, G. and Mayer, G. *The Cholesky method for interval data* (in preparation).

[3] Barth, W. and Nuding, E. *Optimale Lösung von Intervallgleichungssystemen.* Computing **12** (1974), pp. 117–125.

[4] Buneman, O. *A compact noniterative Poisson solver.* Institute for Plasma Research Report 294, Stanford University, 1969.

[5] Buzbee, B., Golub, G., and Nielson, C. *On direct methods for solving Poisson's equation.* SIAM J. Num. Anal. **7** (1970), pp. 627–656.

[6] Kulisch, U. and Miranker, W. *Computer arithmetic in theory and practice.* Academic Press, New York, 1981.

[7] Kulisch, U. and Miranker, W. *A new approach to scientific computation.* Academic Press, New York, 1983.

[8] Schwandt, H. *An interval arithmetic approach for the construction of an almost globally convergent method for the solution of the nonlinear Poisson equation on the unit square.* SIAM J. Sci. Stat. Comp. **5** (1984), pp. 427–452.

[9] Schwandt, H. *Interval arithmetic methods for systems of nonlinear equations arising from discretizations of quazilinear elliptic and parabolic partial differential equations.* Appl. Num. Math. **3** (1987), pp. 257–287.

[10] Schwandt, H. *Cyclic reduction for tridiagonal systems of equations with interval coefficients on vector computers.* SIAM J. Num. Anal. **26** (1989), pp. 661–680.

[11] Schwandt, H. *The interval Buneman algorithm for arbitrary block dimension.* To appear in Computing Suppl. **9** (1993).

[12] Schwandt, H. *An interval arithmetic domain decomposition method for a class of elliptic PDE on non-rectangular domains.* To appear in J. Comp. Appl. Math. (1993).

[13] Sweet, R. *A cyclic reduction algorithm for solving block tridiagonal systems of arbitrary dimension.* SIAM J. Num. Anal. **14** (1977), pp. 706–720.

[14] Varga, R. *Matrix iterative analysis.* Prentice Hall, Englewood Cliffs, New Jersey, 1962.

Technische Universität Berlin
Fachbereich Mathematik, MA 6–4
Straße des 17. Juni 136
D–10623 Berlin
Germany
E-mail: `schwandt@math.tu-berlin.de`