

NUMERICAL STABILITY OF STRUCTURE-PRESERVING DIRECT SOLVERS FOR CENTROSYMMETRIC LINEAR SYSTEMS *

SARAH NATAJ* CHEN GREIF† MANFRED TRUMMER‡

Abstract. This paper analyzes direct solvers for centrosymmetric linear systems, applying structure-preserving factorizations with a particular focus on assessing their stability. We build on existing algorithms and complement the factorizations with equilibration and mixed-precision computations. The solvers are applied to linear systems arising from spectral discretizations of partial differential equations and the results demonstrate their effectiveness. We evaluate the accumulation of roundoff errors during the computational process and their impact on the numerical solution. The study demonstrates that errors originating from the factorization of the matrix and a modified substitution propagate in a stable manner, establishing the direct solver’s robustness. Additionally, we provide insights into the solver’s stability by proving a bound on the relative error.

Key words. direct solution of linear systems, centrosymmetric and skew-centrosymmetric matrices, numerical stability, mixed precision

AMS subject classifications. 65F05, 65G50, 65N35

1. Introduction. A matrix $\mathcal{A} \in \mathbb{R}^{n \times n}$ is centrosymmetric if it is symmetric with respect to its center, i.e. it satisfies the condition $\mathcal{A}_{ij} = \mathcal{A}_{n-i+1, n-j+1}$ for all i, j . These matrices are commonly observed in a variety of applications in computational science and engineering, such as the numerical solution of partial differential equations (PDEs), signal processing, Markov processes, and more [8, 11, 16]. Important cases of centrosymmetric matrices are symmetric Toeplitz matrices and the discretization of the Laplacian by spectral collocation methods. Numerical methods for solving problems involving such matrices have been explored in the literature; see for instance [7, 14]. Centrosymmetric matrices can be considered in some ways as related to a large family of structured matrices that include Hamiltonian, J -symmetric and persymmetric matrices [19, 20, 26].

Given the special structure of centrosymmetric matrices, substantial speedup and reduction in storage requirements can be accomplished, in comparison with general linear solvers. Andrew [3] proposed to solve two linear systems of half the size each, cutting the storage by half and the computational time to a quarter of the standard method. This important early work of Andrew has been expanded in subsequent research, as detailed in [1, 10].

Structure-preserving factorizations aim to maintain properties of the original matrix. In the case of centrosymmetric matrices, these factorizations ensure that the symmetry with respect to the center is retained in the factors. Utilizing such factorizations may lead to more efficient and robust algorithms. In applications such as signal processing, control theory, and physics, preserving the structure in factorizations is potentially beneficial in ensuring that the algorithms accurately reflect the underlying physical or theoretical models.

Burnik [4] introduced a new factorization for a centrosymmetric matrix as a product of an orthogonal matrix Q and a matrix X of a so-called double-cone structure, both centrosymmetric. An analogous QX factorization was also proposed by Steele et al. [24] based on a similarity transformation of the centrosymmetric matrix. In [14] we developed a structure-preserving

*Received... Accepted... Published online on... Recommended by...

†Department of Mathematics, Simon Fraser University, Burnaby, BC, Canada and Department of Computer Science, The University of British Columbia, Vancouver, BC, Canada
Current address: Department of Applied Mathematics, University of Waterloo, ON, Canada, Email: sarah.nataj@uwaterloo.ca

‡Department of Computer Science, The University of British Columbia, Vancouver, BC, Canada, Email: greif@cs.ubc.ca

§Department of Mathematics, Simon Fraser University, BC, Canada, Email: trummer@sfu.ca

LU-type factorization for centrosymmetric matrices, which we called XY factorization or double-cone factorization, where X and Y are centrosymmetric double-cone matrices. This factorization is computed by using a similarity transformation of the centrosymmetric matrix. We applied incomplete XY factorizations as preconditioners for iterative methods in solving sparse linear systems arising in spectral discretization of some linear PDEs [14].

In this work we investigate the numerical stability of a direct solver based on the XY factorization of centrosymmetric matrices originally presented in [14]. There are relatively few papers addressing the error analysis of numerical solutions for linear systems involving centrosymmetric matrices. Lv and Zheng in [18] performed a perturbation analysis of the QX factorization proposed by Burnik in [4].

In an effort to construct an efficient and robust solver, we apply equilibration and iterative refinement with mixed-precision arithmetic, and explore the suitability of these techniques to a set of centrosymmetric and skew-centrosymmetric matrices arising from the important class of spectral methods for the numerical solution of PDEs.

Our analysis is primarily focused on performing an error assessment to determine if errors propagate stably through each step of the given algorithm, including matrix multiplications, the factorization process and a (modified) backward substitution which will be discussed later.

An outline of the remainder of this paper follows. In Section 2, we discuss some known properties of centrosymmetric matrices, review the structure-preserving LU-type factorization for centrosymmetric matrices, and a direct solver based on this factorization. In Section 3, equilibration and iterative refinement algorithms with mixed precision are presented. In Section 4 we investigate the stability of the proposed direct solver based on structure-preserving factorization for centrosymmetric linear systems. Section 5 offers some numerical experiments that validate our claims. Finally, in Section 6, we draw conclusions.

2. Direct solver for centrosymmetric and skew-centrosymmetric linear systems. A matrix $\mathcal{A} \in \mathbb{R}^{n \times n}$ is centrosymmetric iff $J\mathcal{A}J = \mathcal{A}$, where $J \in \mathbb{R}^{n \times n}$ is the flip matrix, a matrix that has ones along the anti-diagonal and zero elsewhere. If $J\mathcal{A}J = -\mathcal{A}$, then \mathcal{A} is called skew-centrosymmetric.

A key finding regarding centrosymmetric matrices is a special similarity transformation, as demonstrated by Weaver [27]. Let $\mathcal{A} \in \mathbb{R}^{n \times n}$ be centrosymmetric and $n = 2k$. There exist matrices $A, C \in \mathbb{R}^{k \times k}$ such that

$$(2.1) \quad \mathcal{A} = \begin{bmatrix} A & JCJ \\ C & JAJ \end{bmatrix}.$$

The matrix can be expressed by the similarity transformation

$$(2.2) \quad \mathcal{A} = U \begin{bmatrix} A + JC & 0 \\ 0 & A - JC \end{bmatrix} U^T, \quad U = \frac{1}{\sqrt{2}} \begin{bmatrix} I & I \\ J & -J \end{bmatrix},$$

where $I \in \mathbb{R}^{k \times k}$ represents the identity matrix. Here U is orthogonal, which means the spectrum of \mathcal{A} consists of the union of the spectra of $A \pm JC$. A similar result is applicable when n is odd: a square centrosymmetric matrix \mathcal{A} of order $n = 2k + 1$ has the form

$$\mathcal{A} = \begin{bmatrix} A & z & JCJ \\ y^T & q & y^T J \\ C & Jz & JAJ \end{bmatrix},$$

where $A, C, J \in \mathbb{R}^{k \times k}$, $z, y \in \mathbb{R}^k$ and q is a constant. Then

$$(2.3) \quad \mathcal{A} = U \begin{bmatrix} A + JC & \sqrt{2}z & 0 \\ \sqrt{2}y^T & q & 0 \\ 0 & 0 & A - JC \end{bmatrix} U^T, \quad U = \frac{1}{\sqrt{2}} \begin{bmatrix} I & 0 & I \\ 0 & \sqrt{2} & 0 \\ J & 0 & -J \end{bmatrix}.$$

where $I \in \mathbb{R}^{k \times k}$ is the identity matrix.

A skew-centrosymmetric matrix \mathcal{A} of even or odd dimensions has, respectively, the form:

$$\mathcal{A} = \begin{bmatrix} A & -JCJ \\ C & -JAJ \end{bmatrix}, \quad \text{or} \quad \mathcal{A} = \begin{bmatrix} A & z & -JCJ \\ y^T & 0 & -y^T J \\ C & -Jz & -JAJ \end{bmatrix}.$$

If \mathcal{A} is skew-centrosymmetric and of even dimension, then there is an orthogonal skew-centrosymmetric matrix E such that $E\mathcal{A}$ is centrosymmetric, where

$$(2.4) \quad E = \begin{bmatrix} -I & 0 \\ 0 & I \end{bmatrix}.$$

This result cannot be extended to skew-centrosymmetric matrices of odd order [1].

The set \mathcal{C}_n of $n \times n$ centrosymmetric matrices is an algebra: If $\mathcal{A}, \mathcal{B} \in \mathcal{C}_n$ and $a \in \mathbb{R}$, then $\mathcal{A} + \mathcal{B}$, $\mathcal{A}\mathcal{B}$, $a\mathcal{A} \in \mathcal{C}_n$. If $\mathcal{A} \in \mathcal{C}_n$, so is \mathcal{A}^T . If $\mathcal{A} \in \mathcal{C}_n$ is invertible, then $\mathcal{A}^{-1} \in \mathcal{C}_n$, and each diagonal block in (2.2) is invertible [13]. For various other properties of centrosymmetric matrices refer to [2, 25, 27].

Let $\lfloor a \rfloor$ and $\lceil a \rceil$ denote the floor and ceiling of a real number a , respectively.

DEFINITION 2.1. Let $n \geq 3$. For a given k , $1 \leq k \leq \lceil n/2 \rceil - 1$, consider the following two-column sub-matrix of $A = (a_{ij}) \in \mathbb{R}^{n \times n}$,

$$\begin{bmatrix} a_{k+1,k} & a_{k+1,n-k+1} \\ \vdots & \vdots \\ a_{n-k,k} & a_{n-k,n-k+1} \end{bmatrix}.$$

The matrix A is called a vertical double-cone, or v-double-cone, if every two-column sub-matrix of A has all zero entries for each $1 \leq k \leq \lceil n/2 \rceil - 1$.

DEFINITION 2.2. Let $n \geq 3$. For a given k , $1 \leq k \leq \lceil n/2 \rceil - 1$, consider the following two-row sub-matrix of $A = (a_{ij}) \in \mathbb{R}^{n \times n}$,

$$\begin{bmatrix} a_{k,k+1} & \cdots & a_{k,n-k} \\ a_{n-k+1,k+1} & \cdots & a_{n-k+1,n-k} \end{bmatrix}.$$

The matrix A is called a horizontal double-cone, or h-double-cone, if every two-row sub-matrix of A has all zero entries for each $1 \leq k \leq \lceil n/2 \rceil - 1$.

We call a matrix “double-cone” if it is either v-double-cone or h-double-cone. Note that all 1×1 and 2×2 matrices are double-cone.

Example: A_1 represents a vertical double-cone (v-double cone) centrosymmetric matrix, and A_2 represents a horizontal double-cone (h-double cone) centrosymmetric matrix,

$$A_1 = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} & a_{15} \\ 0 & a_{22} & a_{23} & a_{24} & 0 \\ 0 & 0 & a_{33} & 0 & 0 \\ 0 & a_{24} & a_{23} & a_{22} & 0 \\ a_{15} & a_{14} & a_{13} & a_{12} & a_{11} \end{bmatrix}, \quad A_2 = \begin{bmatrix} a_{11} & 0 & 0 & 0 & a_{51} \\ a_{21} & a_{22} & 0 & a_{42} & a_{41} \\ a_{31} & a_{32} & a_{33} & a_{32} & a_{31} \\ a_{41} & a_{42} & 0 & a_{22} & a_{21} \\ a_{51} & 0 & 0 & 0 & a_{11} \end{bmatrix}.$$

For a nonsingular $\mathcal{A} \in \mathcal{C}_n$ with n even, each diagonal block in the similarity transformation (2.2) of \mathcal{A} is nonsingular. Then there are permutation matrices P_1 and P_2 , unit lower triangular matrices L_1 and L_2 and nonsingular upper triangular matrices U_1 and U_2 such that $P_1(A + JC) = L_1U_1$ and $P_2(A - JC) = L_2U_2$. Applying the LU factorization to

each diagonal block of (2.2) leads to a factorization of the form $QA = XY$, where Q is centrosymmetric orthogonal given by

$$(2.5) \quad Q = \mathcal{U} \begin{bmatrix} P_1 & 0 \\ 0 & P_2 \end{bmatrix} \mathcal{U}^T,$$

and X and Y are centrosymmetric double-cone matrices,

$$(2.6) \quad X = \mathcal{U} \begin{bmatrix} L_1 & 0 \\ 0 & L_2 \end{bmatrix} \mathcal{U}^T, \quad Y = \mathcal{U} \begin{bmatrix} U_1 & 0 \\ 0 & U_2 \end{bmatrix} \mathcal{U}^T.$$

We call this factorization “double-cone factorization” or “XY factorization” of a centrosymmetric matrix [14]. When dealing with an odd-sized matrix, a similar factorization applies.

Consider a nonsingular $\mathcal{A} \in \mathcal{C}_n$ with n odd. Define $m = \lfloor n/2 \rfloor$. Then, in (2.3) both diagonal blocks are nonsingular and there exist unit lower triangular matrices L_1 and L_2 and nonsingular upper triangular matrices U_1 and U_2 such that

$$P_1 \begin{bmatrix} A + JC & \sqrt{2}z \\ \sqrt{2}y^T & q \end{bmatrix} = L_1 U_1 = \begin{bmatrix} \hat{L}_1 & 0 \\ \ell^T & 1 \end{bmatrix} \begin{bmatrix} \hat{U}_1 & u \\ 0 & \rho \end{bmatrix}, \quad P_2(A - JC) = L_2 U_2,$$

where P_1 and P_2 are permutation matrices, \hat{L}_1 is unit lower triangular, \hat{U}_1 is upper triangular, ℓ and u are vectors of length m , and ρ is a nonzero real number. Also P_1 can be written as $P_1 = \begin{bmatrix} \hat{P}_1 & s \\ t^T & \gamma \end{bmatrix}$, where $\hat{P}_1 \in \mathbb{R}^{m \times m}$, s and t are vectors of size m , and γ is zero or one. By similar calculations as for the even case, we obtain $QA = XY$, where Q is an orthogonal centrosymmetric matrix given by

$$(2.7) \quad Q = \mathcal{U} \begin{bmatrix} \begin{bmatrix} \hat{P}_1 & s \\ t^T & \gamma \end{bmatrix} & \\ & P_2 \end{bmatrix} \mathcal{U}^T,$$

and X and Y are centrosymmetric double-cone matrices,

$$(2.8) \quad X = \mathcal{U} \begin{bmatrix} \begin{bmatrix} \hat{L}_1 & 0 \\ \ell^T & 1 \end{bmatrix} & \\ & L_2 \end{bmatrix} \mathcal{U}^T, \quad Y = \mathcal{U} \begin{bmatrix} \begin{bmatrix} \hat{U}_1 & u \\ 0 & \rho \end{bmatrix} & \\ & U_2 \end{bmatrix} \mathcal{U}^T.$$

To solve the linear system $\mathcal{A}z = b$, the first step is to solve $Xw = \tilde{b}$, with $\tilde{b} = Qb$, followed by solving $Yz = w$. These linear systems involve double-cone matrices and can be efficiently solved using a modified backward substitution method as detailed in [4], and outlined in Algorithms 1 and 2. Algorithm 3 presents a direct solver for $n \times n$ centrosymmetric linear system $\mathcal{A}z = b$ based on the XY factorization.

Notice that Q is not a permutation matrix, however, the action of Q on the vector b can be executed using permutation matrices allowing the formation of \tilde{b} without the need to store Q ,

$$(2.9) \quad \begin{aligned} \tilde{b} = Qb &= \frac{1}{2} \begin{bmatrix} P_1 + P_2 & (P_1 - P_2)J \\ J(P_1 - P_2) & J(P_1 + P_2)J \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} P_1 b_1 + P_2 b_2 + P_1 J b_1 - P_2 J b_2 \\ J P_1 b_1 - J P_2 b_2 + J P_1 J b_1 + J P_2 J b_2 \end{bmatrix}. \end{aligned}$$

In order to solve $\mathcal{A}z = b$, where \mathcal{A} is skew-centrosymmetric, we multiply both side of the equation by E given by (2.4) and solve $\mathcal{B}z = \bar{b}$ using the double-cone factorization, where $\mathcal{B} = E\mathcal{A}$ is centrosymmetric and $\bar{b} = Eb$.

Algorithm 1 Modified backward substitution for solving an h-double cone linear system $Xw = b$, where $X \in \mathcal{C}_n$.

- 1: Let w be a zero vector of size n and $k = \lfloor n/2 \rfloor$
- 2: For $j = 1, \dots, k$, $p = j$ and $q = n - j + 1$ solve

$$\begin{bmatrix} x_{pp} & x_{pq} \\ x_{qp} & x_{qq} \end{bmatrix} \begin{bmatrix} w_p \\ w_q \end{bmatrix} = \begin{bmatrix} b_p - \sum_{i=1}^{p-1} x_{pi} w_i - \sum_{i=q+1}^n x_{pi} w_i \\ b_q - \sum_{i=1}^{p-1} x_{qi} w_i - \sum_{i=q+1}^n x_{qi} w_i \end{bmatrix}$$

- 3: If n is odd set $w_{k+1} = (\hat{b}_{k+1} - \sum_{i=1}^{k-1} x_{pi} w_i - \sum_{i=k+1}^n x_{pi} w_i) / x_{k+1, k+1}$
-

Algorithm 2 Modified backward substitution for solving an v-double cone linear system $Yz = w$, where $Y \in \mathcal{C}_n$.

- 1: Let z be a zero vector of size n and $k = \lfloor n/2 \rfloor$
- 2: If n is odd set $z_{k+1} = w_{k+1} / y_{k+1, k+1}$

- 3: For $j = 1, \dots, k$, $p = k - j + 1$ and $q = k + j$ solve $\begin{bmatrix} y_{pp} & y_{pq} \\ y_{qp} & y_{qq} \end{bmatrix} \begin{bmatrix} z_p \\ z_q \end{bmatrix} = \begin{bmatrix} w_p - \sum_{i=p+1}^{q-1} y_{pi} z_i \\ w_q - \sum_{i=p+1}^{q-1} y_{qi} z_i \end{bmatrix}$
-

Using the double-cone factorization to solve a centrosymmetric system is asymptotically four times faster than solving by a standard LU factorization. This is the same speed-up as the approach suggested by Andrew in [2]. When the given matrix is symmetric positive definite (SPD), Cholesky factorizations are used for the diagonal blocks of the similarity transformation of the matrix in a manner similar to the nonsymmetric case, obtaining similar gains. We call it the XX^T factorization in this case.

3. Equilibration and iterative refinement with mixed precision. Equilibration [12] is an effective way to improve the conditioning of linear systems. We use the row and column equilibration algorithm by Knight et al. [17], which is known to preserve symmetry. Let $A \in \mathbb{R}^{n \times n}$ be a nonsingular matrix. Algorithm 4 computes nonsingular diagonal matrices R and S such that $B = RAS$ has the property that $\max_k |b_{ik}| = \max_k |b_{ki}| = 1$ for all $1 \leq i \leq n$. This algorithm is linearly convergent and permutation invariant. Note that matrices appearing in spectral methods often are extremely ill-conditioned. Our experiments in Section 5.3 demonstrate significant reduction in condition number when we apply equilibration to spectral differentiation matrices.

PROPOSITION 1. *Algorithm 4 (equilibration) preserves centrosymmetry.*

Proof. Assume A is centrosymmetric. In each step of the algorithm the diagonal matrices \tilde{R} and \tilde{S} are centrosymmetric, hence $B = RAS$ is a product of centrosymmetric matrices. Recall that the set of $n \times n$ centrosymmetric matrices form an algebra [27]. Therefore $B = RAS$ is indeed centrosymmetric. \square

Iterative refinement algorithms are widely used to improve the accuracy of the numerical solution of linear systems. Carson and Higham [6] introduced a GMRES iterative refinement algorithm (GMRES-IR) with mixed precision. The factors of the matrix are computed in low precision. The algorithm then solves the correction equation with GMRES in high precision,

Algorithm 3 Direct solver based on the XY factorization for solving an $n \times n$ centrosymmetric linear system $\mathcal{A}z = b$

- 1: Find a similarity transformation of $\mathcal{A} \in \mathcal{C}_n$ of the form of (2.2) and (2.3)
 - 2: Compute the LU factorization of diagonal blocks in (2.2) and (2.3)
 - 3: Form X and Y so that $Q\mathcal{A} = XY$ using (2.5)–(2.6) and (2.7)–(2.8)
 - 4: Set $\tilde{b} = Qb$ using (2.9) and solve $Xw = \tilde{b}$ using modified backward substitution (Algorithm 1)
 - 5: Solve $Yz = w$ using modified backward substitution (Algorithm 2)
-

Algorithm 4 Row and column equilibration of a matrix $\mathcal{A} \in \mathbb{R}^{n \times n}$

```

1: Let  $R = \text{eye}(n)$ ;  $S = \text{eye}(n)$ ;  $r = \text{zeros}(n, 1)$ ;  $s = \text{zeros}(n, 1)$ 
2: while  $\max_i |r(i) - 1| > \text{tol}$  or  $\max_i |s(i) - 1| > \text{tol}$ 
3:   for  $i = 1 : n$  do
4:      $r(i) = \|\mathcal{A}(i, :)\|_{\infty}^{-1/2}$ 
5:      $s(i) = \|\mathcal{A}(:, i)\|_{\infty}^{-1/2}$ 
6:   end for
7:    $\tilde{R} = \text{diag}(r)$ ;  $\tilde{S} = \text{diag}(s)$ 
8:    $\mathcal{A} = \tilde{R} \mathcal{A} \tilde{S}$ 
9:    $R = \tilde{R} R$ 
10:   $S = S \tilde{S}$ 
11: end while

```

using the product of the factors as a preconditioner. Their analysis shows that GMRES-IR with mixed precision can provide accurate solutions to systems with condition numbers of magnitude u^{-1} and larger, where u is the unit machine roundoff. Their algorithm uses three precisions plus the working precision. Let u be the working precision in which the matrix A and vector b are stored, u_f the precision in which the factorization of A is computed, u_r the precision in which the residual is calculated, and u_s the precision in which the correction equation is solved. Usually $u_r \leq u \leq u_s \leq u_f$. In Algorithm 5, we adapt GMRES-IR for solving centrosymmetric systems by using the double-cone factorization and double-cone solvers that were introduced in Section 2. Numerical results applying Algorithm 5 are given in Section 5.3.

Algorithm 5 GMRES- IR adopted for double-cone factorization for a centrosymmetric system $\mathcal{A}x = b$ with $\mathcal{A} \in \mathcal{C}_n$ and $b \in \mathbb{R}^n$ given in precision u

```

1: Obtain  $S$  and  $R$  from Algorithm 4. Calculate  $\hat{\mathcal{A}} = RAS$  and  $\hat{b} = Rb$  in precision  $u$ 
2: Compute a double-cone factorization  $\hat{Q}\hat{\mathcal{A}} = \hat{X}\hat{Y}$  in precision  $u_f$ 
3: Solve  $\hat{\mathcal{A}}y_0 = \hat{b}$  in precision  $u_f$  using the  $\hat{X}\hat{Y}$  factors and double-cone solvers, store the approximate solution  $x_0 = Sy_0$  in precision  $u$ 
4: for  $i = 1 : i_{\max}$  do
5:   Compute  $r_i = b - \mathcal{A}x_i$  at precision  $u_r$  and round  $r_i$  to precision  $u$ 
6:   Solve  $MAd_i = Mr_i$  by GMRES in precision  $u$ , with  $M = S\hat{Y}^{-1}\hat{X}^{-1}\hat{Q}R$ , where the matrix vector product in GMRES is computed in precision  $u_r$ , store  $d_i$  in precision  $u$ 
7:   Compute  $x_{i+1} = x_i + d_i$  at precision  $u$ 
8:   if converged then
9:     return  $x_{i+1}$ , quit
10:  end if
11: end for

```

4. Error analysis for the direct solver. We now investigate the stability of the direct solver based on XY factorization for centrosymmetric matrices.

4.1. Preliminaries. Throughout this section we use the following notation. For $A = (a_{ij}) \in \mathbb{R}^{m \times n}$, define $|A| = (|a_{ij}|)$. The matrix Euclidean, infinity and Frobenius norms are denoted by $\|A\|_2$, $\|A\|_{\infty}$ and $\|A\|_F$, respectively. Additionally, define

$$\|A\|_M := \max_{i,j} |a_{ij}|.$$

This norm is not consistent; the best bound that holds for all $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$, is $\|AB\|_M \leq n\|A\|_M\|B\|_M$. Also we have $\|A\|_M \leq \|A\|_{\infty} \leq n\|A\|_M$, for all $A \in \mathbb{R}^{m \times n}$. Let u be the unit roundoff. We consider a model of floating point arithmetic where a single basic operation yields a relative error δ bounded by u ,

$$\text{fl}(x \text{ op } y) = (1 + \delta)(x \text{ op } y), \quad |\delta| \leq u, \quad \text{op} \in \{+, -, \times, \div\}.$$

We follow Higham's standard model of arithmetic [15, Section 2.4], which may not capture all features that we expect in a floating point system.

LEMMA 2. [15, Lemma 3. 1] If $|\delta_i| \leq u$, for $i = 1, \dots, k$, where k is any positive integer and $\rho = \pm 1$, and

$$(4.1) \quad ku < 1,$$

then $\prod_{i=1}^k (1 + \delta_i)^{\rho_i} = 1 + \theta_k$, where

$$(4.2) \quad |\theta_k| \leq \frac{ku}{1 - ku} =: \gamma_k.$$

Two useful properties of γ_k are given in the following lemma.

LEMMA 3. [15, Lemma 3. 3] For any positive integers n, m , the following relations hold:

$$\begin{aligned} \gamma_m + \gamma_n + \gamma_m \gamma_n &\leq \gamma_{m+n}, \\ c\gamma_n &\leq \gamma_{cn}, \quad c \geq 1. \end{aligned}$$

We make the following assumptions about matrix operations [15]. If $A, B \in \mathbb{R}^{n \times n}$, $\alpha \in \mathbb{R}$,

$$\begin{aligned} \text{fl}(\alpha A) &= \alpha A + E, & |E| &\leq u|\alpha A|, \\ \text{fl}(A + B) &= (A + B) + E, & |E| &\leq u|A + B|, \\ \text{fl}(AB) &= AB + E, & |E| &\leq \gamma_n |A| |B|. \end{aligned}$$

For a product of matrices, the corresponding normwise bounds are given by [15, Section 3.5]

$$(4.3) \quad \|E\|_p \leq c_1(n)u\|A\|_p\|B\|_p + O(u^2), \quad p = 1, \infty, F, M,$$

where $c_1(n)$ denotes a constant depending polynomially on n . We also assume the LU factorization is done in such a way that the computed L and U factors of $A \in \mathbb{R}^{n \times n}$, which we denote by \hat{L} and \hat{U} , satisfy

$$(4.4) \quad \hat{L}\hat{U} = A + \Delta A, \quad |\Delta A| \leq \gamma_n |\hat{L}| |\hat{U}|.$$

Algorithms for computing the factorizations that satisfy these bounds can be found in [15, Algorithm 9. 2] and in [12, Algorithm 3. 2. 1]. The corresponding normwise bounds are

$$(4.5) \quad \|\Delta A\|_p \leq c_2(n)u\|\hat{L}\|_p\|\hat{U}\|_p + O(u^2), \quad p = 1, \infty, F, M,$$

where $c_2(n)$ denotes a constant polynomially dependent on n .

The following theorem gives a bound on the backward error for the solution of general linear systems using Gaussian Elimination with partial pivoting (GEPP). Define the growth factor of a matrix as

$$\rho_n = \frac{\max_{i,j,k} |a_{ij}^{(k)}|}{\max_{i,j} |a_{ij}|},$$

where $|a_{ij}^{(k)}|$, $k = 1, \dots, n - 1$ are the elements that occur during the elimination. It can be shown that $\rho_n \leq 2^{n-1}$ for partial pivoting. Also if A is diagonally dominant by rows or columns, then $\rho_n \leq 2$, see [15, Section 9. 5].

THEOREM 4. ([28, Chapter 3, Section 25], [15, Section 9. 3]) If GEPP is used to produce a computed solution \hat{x} to $Ax = b$, then

$$(4.6) \quad (A + \Delta A)\hat{x} = b, \quad \|\Delta A\|_\infty \leq n^2 \gamma_{3n} \rho_n \|A\|_\infty.$$

4.2. Roundoff error analysis for the modified backward substitution. Let us consider Algorithms 1 and 2.

THEOREM 5. *Let X be an $n \times n$ centrosymmetric h-double-cone matrix and $b \in \mathbb{R}^n$. Assume that (4.1), (4.4), and (4.6) hold, and $k = \lfloor n/2 \rfloor \geq 112$. If the solution \hat{w} to the linear system $Xw = b$ is computed using the modified backward substitution of Algorithm 1, then*

$$(4.7) \quad (X + \Delta X)\hat{w} = b, \quad \|\Delta X\|_M \leq \gamma_{2k}\|X\|_M.$$

Proof. We consider Algorithm 1 for an h-double cone matrix of even order; the odd case can be carried out in the same way. In step 2 of Algorithm 1, we need to solve a 2×2 linear symmetric centrosymmetric system. Let $E = \begin{bmatrix} x_{pp} & x_{pq} \\ x_{qp} & x_{qq} \end{bmatrix}$. Note that in fact $x_{pp} = x_{qq}$ and $x_{pq} = x_{qp}$, given that E is symmetric centrosymmetric, but we will not take advantage of these equalities here, and use instead Theorem 4, which applies to general matrices. To that end, using Lemma 2, the accumulated roundoff error in solving this linear system is

$$(E + \Delta E) \begin{bmatrix} w_p \\ w_q \end{bmatrix} = \begin{bmatrix} b_p(1 + \epsilon_1) - \sum_{i=1}^{p-1} x_{pi}w_i(1 + \theta_{pi}) - \sum_{i=q+1}^n x_{pi}w_i(1 + \theta'_{pi}) \\ b_q(1 + \epsilon_2) - \sum_{i=1}^{p-1} x_{qi}w_i(1 + \theta_{qi}) - \sum_{i=q+1}^n x_{qi}w_i(1 + \theta'_{qi}) \end{bmatrix},$$

where $|\epsilon_1|, |\epsilon_2| \leq u$, $|\theta_{pi}|, |\theta'_{pi}| \leq \gamma_{p-1}$, $|\theta_{qi}|, |\theta'_{qi}| \leq \gamma_{q-1}$ and

$$\Delta E = \begin{bmatrix} \delta e_{pp} & \delta e_{pq} \\ \delta e_{qp} & \delta e_{qq} \end{bmatrix},$$

is the roundoff error matrix from solving the 2×2 linear system by GEPP.

Notice that the growth factor of this matrix is $\rho_2 \leq 2$. Therefore using (4.6) and norm properties $\|\Delta E\|_M \leq \|\Delta E\|_\infty \leq 16\gamma_6\|E\|_M$,

$$(4.8) \quad |\delta e_{pp}| \leq \|\Delta E\|_M \leq 16\gamma_6\|E\|_M = 16\gamma_6 \max\{|x_{pp}|, |x_{pq}|\}.$$

Since E is symmetric centrosymmetric, the same bound holds for other entries of ΔE . Using Lemma 2, we obtain:

$$(4.9) \quad \begin{bmatrix} (x_{pp} + \delta e_{pp})(1 + \eta_{pp}) & (x_{pq} + \delta e_{pq})(1 + \eta_{pq}) \\ (x_{pq} + \delta e_{qp})(1 + \eta_{qp}) & (x_{pp} + \delta e_{qq})(1 + \eta_{qq}) \end{bmatrix} \begin{bmatrix} w_p \\ w_q \end{bmatrix} \\ + \begin{bmatrix} \sum_{i=1}^{p-1} x_{pi}w_i(1 + \eta_{pi}) - \sum_{i=q+1}^n x_{pi}w_i(1 + \eta'_{pi}) \\ \sum_{i=1}^{p-1} x_{qi}w_i(1 + \eta_{qi}) - \sum_{i=q+1}^n x_{qi}w_i(1 + \eta'_{qi}) \end{bmatrix} = \begin{bmatrix} b_p \\ b_q \end{bmatrix},$$

where

$$(4.10) \quad |\eta_{pp}|, |\eta_{pq}|, |\eta_{qp}|, |\eta_{qq}| \leq \gamma_1,$$

and

$$(4.11) \quad |\eta_{pi}|, |\eta'_{pi}| \leq \gamma_p, \quad |\eta_{qi}|, |\eta'_{qi}| \leq \gamma_q, \quad i = 1, \dots, p-1, q+1, \dots, n.$$

We write

$$(x_{pp} + \delta e_{pp})(1 + \eta_{pp}) = x_{pp} + \delta e_{pp} + \eta_{pp}x_{pp} + \delta e_{pp}\eta_{pp} =: x_{pp} + \delta x_{pp},$$

and then, a bound for $|\delta x_{pp}|$ is

$$(4.12) \quad |\delta x_{pp}| \leq |\delta e_{pp}| + |\eta_{pp}||x_{pp}| + |\delta e_{pp}||\eta_{pp}|.$$

Let us consider the case where X is a 4×4 centrosymmetric h-double-cone matrix. Assume the computed solution is the exact solution of $(X + \Delta X)w = b$. Incorporating Equations (4.8)–(4.11) into (4.12), we have

$$\begin{aligned}
 |\Delta X| \leq & \begin{bmatrix} \gamma_1|x_{11}| & 0 & 0 & \gamma_1|x_{14}| \\ \gamma_2|x_{21}| & \gamma_1|x_{22}| & \gamma_1|x_{23}| & \gamma_2|x_{24}| \\ \gamma_2|x_{24}| & \gamma_1|x_{23}| & \gamma_1|x_{22}| & \gamma_2|x_{21}| \\ \gamma_1|x_{41}| & 0 & 0 & \gamma_1|x_{11}| \end{bmatrix} \\
 & + 16(\gamma_6 + \gamma_6\gamma_1) \begin{bmatrix} \max\{|x_{11}|, |x_{14}|\} & 0 & 0 & \max\{|x_{11}|, |x_{14}|\} \\ 0 & \max\{|x_{22}|, |x_{23}|\} & \max\{|x_{22}|, |x_{23}|\} & 0 \\ 0 & \max\{|x_{22}|, |x_{23}|\} & \max\{|x_{22}|, |x_{23}|\} & 0 \\ \max\{|x_{11}|, |x_{14}|\} & 0 & 0 & \max\{|x_{11}|, |x_{14}|\} \end{bmatrix}.
 \end{aligned}$$

Using Lemma 3,

$$\begin{aligned}
 \|\Delta X\|_M & \leq (\gamma_2 + 16\gamma_6 + 16\gamma_6\gamma_1)\|X\|_M \\
 & \leq (\gamma_2 + 16(\gamma_1 + \gamma_6 + \gamma_6\gamma_1))\|X\|_M \\
 & \leq (\gamma_2 + 16\gamma_7)\|X\|_M \\
 & \leq (\gamma_2 + \gamma_{112})\|X\|_M.
 \end{aligned}$$

This can be generalized to $n \times n$ centrosymmetric h-double-cone. If $k = \lfloor n/2 \rfloor \geq 112$, as long as assumption (4.1) holds, we have

$$\|\Delta X\|_M \leq \gamma_{2k}\|X\|_M,$$

which completes the proof. \square

The constant γ_{2k} in this theorem is analogous to that of a triangular solver, which is γ_n [12]. A similar result holds for Algorithm 2.

THEOREM 6. *Let Y be an $n \times n$ centrosymmetric v-double-cone matrix and $w \in \mathbb{R}^n$. Assume that (4.1), (4.4), and (4.6) hold, and $k = \lfloor n/2 \rfloor \geq 112$. If the solution \hat{z} to the linear system $Yz = w$ is computed using the modified backward substitution of Algorithm 2, then*

$$(4.13) \quad (Y + \Delta Y)\hat{z} = w, \quad \|\Delta Y\|_M \leq \gamma_{2k}\|Y\|_M,$$

Proof. The proof is similar to that of Theorem 5. \square

4.3. Roundoff error analysis for the double-cone factorization. Next, we analyze the double cone factorization of a centrosymmetric matrix. The goal is to calculate the accumulative roundoff error and see if the constants in (4.3), (4.5) and (4.7) propagate stably into the final error bound.

THEOREM 7. *Under the assumptions (4.3) and (4.5), the X and Y factors of a centrosymmetric matrix $\mathcal{A} \in \mathcal{C}_n$ computed by (2.5) and (2.6) satisfy*

$$\hat{X}\hat{Y} = Q\mathcal{A} + \Delta\mathcal{A},$$

where

$$(4.14) \quad \|\Delta\mathcal{A}\|_M \leq c(k)u (\|\mathcal{A}\|_M + \|\hat{X}\|_M \|\hat{Y}\|_M) + O(u^2),$$

with $c(k)$ a constant depending polynomially on $k = \lfloor n/2 \rfloor$.

Proof. Given a centrosymmetric matrix $\mathcal{A} \in C_n$, with n even, of block form (2.1), consider Q, X and Y given by (2.5) and (2.6) such that

$$(4.15) \quad \mathcal{A} = \begin{bmatrix} A & JCJ \\ C & JAJ \end{bmatrix} = Q^T XY.$$

Let $B_1 := A + JC$ and $B_2 := A - JC$ be the diagonal blocks in the similarity transformation of \mathcal{A} in (2.2). Then the computed B_1 and B_2 satisfy

$$(4.16) \quad \hat{B}_i = B_i + \bar{E}_i, \quad \|\bar{E}_i\|_M \leq u(\|A\|_M + \|C\|_M) + O(u^2), \quad i = 1, 2.$$

Next, assume the LU factorization is done in such a way that the computed LU factors of B_1 and B_2 satisfy

$$(4.17) \quad \hat{L}_i \hat{U}_i = \hat{B}_i + \bar{F}_i, \quad \|\bar{F}_i\|_M \leq c_2(k)u \|\hat{L}_i\|_M \|\hat{U}_i\|_M + O(u^2), \quad i = 1, 2.$$

Considering (4.16)–(4.17) we have

$$(4.18) \quad \hat{L}_i \hat{U}_i = B_i + \Delta B_i, \quad i = 1, 2.$$

where

$$(4.19) \quad \|\Delta B_i\|_M \leq u(\|A\|_M + \|C\|_M + c_2(k)\|\hat{L}_i\|_M \|\hat{U}_i\|_M) + O(u^2), \quad i = 1, 2.$$

We assume

$$(4.20) \quad \hat{A} = A + \Delta A,$$

where $\Delta A = F + G$ and F is the accumulated roundoff error from the previous steps and G is the roundoff associated with multiplication and addition/subtraction. Then

$$\begin{aligned} \mathcal{A} = Q^T XY &= \frac{1}{8} \begin{bmatrix} P_1^T + P_2^T & (P_1^T - P_2^T)J \\ J(P_1^T - P_2^T) & J(P_1^T + P_2^T)J \end{bmatrix} \\ &\quad \cdot \begin{bmatrix} L_1 + L_2 & (L_1 - L_2)J \\ J(L_1 - L_2) & J(L_1 + L_2)J \end{bmatrix} \begin{bmatrix} U_1 + U_2 & (U_1 - U_2)J \\ J(U_1 - U_2) & J(U_1 + U_2)J \end{bmatrix} \\ &= \frac{1}{4} \begin{bmatrix} P_1^T + P_2^T & (P_1^T - P_2^T)J \\ J(P_1^T - P_2^T) & J(P_1^T + P_2^T)J \end{bmatrix} \begin{bmatrix} L_1 U_1 + L_2 U_2 & L_1 U_1 J - L_2 U_2 J \\ JL_1 U_1 - JL_2 U_2 & JL_1 U_1 J + JL_2 U_2 J \end{bmatrix}. \end{aligned}$$

Then $A = \frac{1}{2}(P_1^T L_1 U_1 + P_2^T L_2 U_2)$. By adding (4.18) for $i = 1, 2$ and including the error of adding and multiplying by scalar, we obtain a bound for F , namely

$$(4.21) \quad \begin{aligned} \|F\|_M &\leq \frac{1}{2}u \left(2\|\hat{L}_1 \hat{U}_1\|_M + 2\|\hat{L}_2 \hat{U}_2\|_M \right. \\ &\quad \left. + 2\|A\|_M + 2\|C\|_M + c_2(k)\|\hat{L}_1\|_M \|\hat{U}_1\|_M + c_2(k)\|\hat{L}_2\|_M \|\hat{U}_2\|_M \right) + O(u^2) \\ &\leq u(\|\mathcal{A}\|_M + (c_2(k) + 2k)\|\hat{X}\|_M \|\hat{Y}\|_M) + O(u^2). \end{aligned}$$

Notice that $\|\hat{L}_1 \hat{U}_1\|_M \leq k\|\hat{L}_1\|_M \|\hat{U}_1\|_M$ and by the definition of \hat{X} and \hat{Y} we have $\|\hat{L}_1\|_M, \|\hat{L}_2\|_M \leq \|\hat{X}\|_M$ and $\|\hat{U}_1\|_M, \|\hat{U}_2\|_M \leq \|\hat{Y}\|_M$, also $\|A\|_M, \|C\|_M \leq \|\mathcal{A}\|_M$.

Next, we calculate the accumulated roundoff error resulting from multiplication and addition/subtraction in (4.20). Consider the first term. Let

$$H_1 := (\hat{L}_1 + \hat{L}_2)(\hat{U}_1 + \hat{U}_2) \quad \text{and} \quad H_2 := (\hat{L}_1 - \hat{L}_2)J^2(\hat{U}_1 - \hat{U}_2).$$

Then $\hat{H}_1 = H_1 + E_1$ and $\hat{H}_2 = H_2 + E_2$, where E_1 and E_2 are the errors of the product. Applying (4.3),

$$\begin{aligned} \|E_1\|_M &\leq c_1(k)u\|\hat{L}_1 + \hat{L}_2\|_M\|\hat{U}_1 + \hat{U}_2\|_M + O(u^2) \leq c_1(k)u\|\hat{X}\|_M\|\hat{Y}\|_M + O(u^2), \\ \|E_2\|_M &\leq c_1(k)u\|\hat{L}_1 - \hat{L}_2\|_M\|\hat{U}_1 - \hat{U}_2\|_M + O(u^2) \leq c_1(k)u\|\hat{X}\|_M\|\hat{Y}\|_M + O(u^2). \end{aligned}$$

Then,

$$\begin{aligned} (4.22) \quad H &:= (P_1^T + P_2^T)((\hat{L}_1 + \hat{L}_2)(\hat{U}_1 + \hat{U}_2) + (\hat{L}_1 - \hat{L}_2)J^2(\hat{U}_1 - \hat{U}_2)) \\ &= (P_1^T + P_2^T)(\hat{H}_1 + \hat{H}_2) \\ &= (P_1^T + P_2^T)(H_1 + E_1 + H_2 + E_2 + E_3) + E_4, \end{aligned}$$

where E_3 is the error resulting from adding H_1 and H_2 , satisfying

$$\begin{aligned} \|E_3\|_M &\leq c_1(k)u\|P_1^T + P_2^T\|_M(\|H_1\|_M + \|H_2\|_M) + O(u^2) \\ &\leq 2c_1(k)ku\|\hat{X}\|_M\|\hat{Y}\|_M + O(u^2), \end{aligned}$$

where we used the fact $\|P_1^T + P_2^T\|_M \leq 2$. Also E_4 is the error resulting from the product in second line of (4.22),

$$\begin{aligned} \|E_4\|_M &\leq c_1(k)ku\|P_1^T + P_2^T\|_M(\|H_1\|_M + \|H_2\|_M) + O(u^2) \\ &\leq 2c_1(k)k^2u\|\hat{X}\|_M\|\hat{Y}\|_M + O(u^2). \end{aligned}$$

Notice that errors propagated by the product of $(P_1^T + P_2^T)$ and E_1, E_2 and E_3 are of order $O(u^2)$. Then the total roundoff error from calculating H is G_1 , where

$$\|G_1\|_M \leq 2c_1(k)k^2u\|\hat{X}\|_M\|\hat{Y}\|_M + O(u^2).$$

Doing the same for the second term of Equation (4.20), with the roundoff error G_2 , the total error from multiplication and addition/subtraction is $G = G_1 + G_2 + G_3$, and

$$\begin{aligned} \|G_3\|_M &\leq \frac{1}{8}(2u) \left(\|(P_1^T + P_2^T)((\hat{L}_1 + \hat{L}_2)(\hat{U}_1 + \hat{U}_2) + (\hat{L}_1 - \hat{L}_2)J^2(\hat{U}_1 - \hat{U}_2))\|_M \right. \\ &\quad \left. + \|(P_1^T - P_2^T)((\hat{L}_1 - \hat{L}_2)(\hat{U}_1 + \hat{U}_2)J + (\hat{L}_1 + \hat{L}_2)J^2(\hat{U}_1 - \hat{U}_2)J)\|_M \right) \\ &\leq 2uk^2\|\hat{X}\|_M\|\hat{Y}\|_M + O(u^2). \end{aligned}$$

Therefore,

$$(4.23) \quad \|G\|_M \leq (4c_1(k)k + 2k^2)u\|\hat{X}\|_M\|\hat{Y}\|_M + O(u^2).$$

Considering the bounds for F and G given by (4.21) and (4.23), we have

$$\|\Delta A\|_M \leq u(\|\mathcal{A}\|_M + c_3(k)(\|\hat{X}\|_M\|\hat{Y}\|_M)) + O(u^2),$$

where $c_3(k) = 4c_1(k)k + c_2(k) + 2k^2 + 2k$. In the same way we can find bounds for the three other blocks of \mathcal{A} . Overall we obtain the bound in (4.14). \square

4.4. Roundoff error analysis for Algorithm 3. We now perform a roundoff error analysis for Algorithm 3 and derive a bound for the relative error.

THEOREM 8. *Let $\mathcal{A} \in \mathbb{C}_n$ and suppose Algorithm 3 produces computed \hat{X} , \hat{Y} , and a computed solution \hat{z} to $\mathcal{A}z = b$. Then*

$$(\mathcal{A} + \Delta\mathcal{A})\hat{z} = b, \quad \|\Delta\mathcal{A}\|_M \leq d(k)u (\|\mathcal{A}\|_M + \|\hat{X}\|_M \|\hat{Y}\|_M) + O(u^2),$$

where $d(k)$ is a constant depending polynomially on $k = \lfloor n/2 \rfloor$ while k is assumed to be large enough while still satisfying the assumption (4.1).

Proof. The solution of the linear system $\mathcal{A}z = b$ can be computed by solving $Xw = \tilde{b}$, where $\tilde{b} = Qb$, followed by $Yz = w$. We calculate $\tilde{b} = Qb$ by using (2.9). Let $\hat{b} = \text{fl}(\tilde{b})$, then

$$\hat{b} = \frac{1}{2} \begin{bmatrix} (P_1b_1 + P_2b_2 + P_1Jb_1 - P_2Jb_2)(1 + \delta_1) \\ (JP_1b_1 - JP_2b_2 + JP_1Jb_1 + JP_2Jb_2)(1 + \delta_2) \end{bmatrix},$$

where $|\delta_1|, |\delta_2| \leq \gamma_4$. Thus $\hat{b} = (Q + \Delta Q)b$, where $|\Delta Q| \leq \gamma_4|Q|$, which results in $\|\Delta Q\|_2 \leq \gamma_4$. By Theorem 4, the modified substitution produces \hat{w} and \hat{z} satisfying

$$\begin{aligned} (\hat{X} + \Delta\hat{X})\hat{w} &= \hat{b}, & \|\Delta\hat{X}\|_M &\leq \gamma_{2k}\|\hat{X}\|_M, \\ (\hat{Y} + \Delta\hat{Y})\hat{z} &= \hat{w}, & \|\Delta\hat{Y}\|_M &\leq \gamma_{2k}\|\hat{Y}\|_M, \end{aligned}$$

where k is assumed to be large enough while still satisfying (4.1). Therefore,

$$\hat{b} = (\hat{X} + \Delta\hat{X})(\hat{Y} + \Delta\hat{Y})\hat{z} = (\hat{X}\hat{Y} + \hat{X}\Delta\hat{Y} + \Delta\hat{X}\hat{Y} + \Delta\hat{X}\Delta\hat{Y})\hat{z}.$$

Multiplying both sides by $(Q + \Delta Q)^{-1}$,

$$(Q + \Delta Q)^{-1}\hat{b} = (Q + \Delta Q)^{-1}(\hat{X}\hat{Y} + \hat{X}\Delta\hat{Y} + \Delta\hat{X}\hat{Y} + \Delta\hat{X}\Delta\hat{Y})\hat{z},$$

we thus have

$$\begin{aligned} (4.24) \quad b &= (Q + \Delta Q)^{-1}(\hat{X}\hat{Y} + \hat{X}\Delta\hat{Y} + \Delta\hat{X}\hat{Y} + \Delta\hat{X}\Delta\hat{Y})\hat{z} \\ &= ((Q + \Delta Q)^{-1}\hat{X}\hat{Y} + F)\hat{z}, \end{aligned}$$

where $F := (Q + \Delta Q)^{-1}(\hat{X}\Delta\hat{Y} + \Delta\hat{X}\hat{Y} + \Delta\hat{X}\Delta\hat{Y})$ and

$$\|F\|_M \leq n\|(Q + \Delta Q)^{-1}\|_M (\|\hat{X}\Delta\hat{Y}\|_M + \|\Delta\hat{X}\hat{Y}\|_M + \|\Delta\hat{X}\Delta\hat{Y}\|_M).$$

We have the following bound for the norm of $(Q + \Delta Q)^{-1}$,

$$\begin{aligned} \|(Q + \Delta Q)^{-1}\|_M &\leq \|(Q + \Delta Q)^{-1}\|_2 = \|(I + Q^T\Delta Q)^{-1}Q^T\|_2 \\ &\leq \frac{1}{1 - \|Q^T\Delta Q\|_2} \leq \frac{1}{1 - \gamma_4} = 1 + \gamma_4 + O(u^2). \end{aligned}$$

Therefore,

$$\begin{aligned} (4.25) \quad \|F\|_M &\leq (1 + \gamma_4 + O(u^2))n^2(2\gamma_{2k}\|\hat{X}\|_M\|\hat{Y}\|_M + \gamma_{2k}^2\|\hat{X}\|_M\|\hat{Y}\|_M) \\ &\leq (2n^2\gamma_{2k} + O(u^2))\|\hat{X}\|_M\|\hat{Y}\|_M. \end{aligned}$$

Also notice that $(Q + \Delta Q)\mathcal{A} = Q\mathcal{A} + \Delta Q\mathcal{A} = \hat{X}\hat{Y} + E + \Delta Q\mathcal{A}$, so

$$(4.26) \quad \mathcal{A} = (Q + \Delta Q)^{-1}\hat{X}\hat{Y} + (Q + \Delta Q)^{-1}(E + \Delta Q\mathcal{A}) =: (Q + \Delta Q)^{-1}\hat{X}\hat{Y} + G,$$

and by using Theorem 5 and under the assumption (4.1),

$$\begin{aligned}
 (4.27) \quad \|G\|_M &\leq n\|(Q + \Delta Q)^{-1}\|_M(\|E\|_M + n\|\Delta Q\|_M \|\mathcal{A}\|_M) \\
 &\leq n(1 + \gamma_4 + O(u^2)) \left(c(k)u(\|\mathcal{A}\|_M + \|\hat{X}\|_M \|\hat{Y}\|_M) + n\gamma_4\|\mathcal{A}\|_M \right) \\
 &\leq c'(k)u(\|\mathcal{A}\|_M + \|\hat{X}\|_M \|\hat{Y}\|_M) + O(u^2).
 \end{aligned}$$

From (4.26) we have $(Q + \Delta Q)^{-1}\hat{X}\hat{Y} = \mathcal{A} - G$, and we substitute this in (4.24):

$$b = (\mathcal{A} - G + F)\hat{z} =: (\mathcal{A} + \Delta\mathcal{A})\hat{z}.$$

By (4.25) and (4.27),

$$\|\Delta\mathcal{A}\|_M \leq d(k)u(\|\mathcal{A}\|_M + \|\hat{X}\|_M \|\hat{Y}\|_M) + O(u^2).$$

where $d(k)$ is a constant depending polynomially on $k = \lfloor n/2 \rfloor$. \square

In Theorem 8, a roundoff error analysis was presented for Algorithm 3. Next, we derive a bound for the relative error of the direct solver. Let $\|\cdot\|_p$ be a consistent matrix norm, $\|AB\|_p \leq \|A\|_p\|B\|_p$. Recall that the M norm does not satisfy this inequality. Suppose that $\mathcal{A}z = b$ and $(\mathcal{A} + \Delta\mathcal{A})\hat{z} = b + \Delta b$, then we apply the following well-known result in perturbation theory to find a bound for the relative error [9, Section 2.2],

$$\frac{\|z - \hat{z}\|_p}{\|z\|_p} \leq \frac{\kappa_p(\mathcal{A})}{1 - \|\mathcal{A}^{-1}\|_p \|\Delta\mathcal{A}\|_p} \left(\frac{\|\Delta\mathcal{A}\|_p}{\|\mathcal{A}\|_p} + \frac{\|\Delta b\|_p}{\|b\|_p} \right),$$

assuming $\|\mathcal{A}^{-1}\|_p \|\Delta\mathcal{A}\|_p < 1$. Here $\kappa_p(\mathcal{A}) = \|\mathcal{A}\|_p \|\mathcal{A}^{-1}\|_p$ is the condition number of the matrix \mathcal{A} . We set $\Delta b = 0$ in this inequality, and by employing Theorem 8 we establish the following bound on the relative error for the approximate solution.

COROLLARY 4.1. *Given the assumptions of Theorem 8, and suppose also that $\|\mathcal{A}^{-1}\|_\infty \|\Delta\mathcal{A}\|_\infty \leq c < 1$. Then*

$$\frac{\|z - \hat{z}\|_\infty}{\|z\|_\infty} \leq d'(k)u\kappa_\infty(\mathcal{A}) \left(1 + \frac{\|\hat{X}\|_\infty \|\hat{Y}\|_\infty}{\|\mathcal{A}\|_\infty} \right) + O(u^2),$$

where $d'(k)$ is a constant depending polynomially on $k = \lfloor n/2 \rfloor$.

5. Numerical results. In this section we examine the performance of the proposed algorithms for solving linear systems that arise in the spectral discretization of a set of linear PDEs. Spectral methods [5, 23] are known to solve PDEs with spectral accuracy if the solution is smooth. Compared to other methods such as finite difference and finite element methods, these methods require fewer degrees of freedom to achieve the same accuracy.

We have listed the PDEs used in this study in Table 5.1; additional details on the spectral discretization can be found in [14]. To discretize the PDEs, we use the spectral collocation method with either Chebyshev or Legendre Gauss-Lobatto nodes. This results in linear systems that have a centrosymmetric structure. In some special cases, such as for Poisson and Helmholtz equations with Robin boundary conditions, centrosymmetry is lost but the matrix is *nearly centrosymmetric*; specifically, it is centrosymmetric except along the diagonal.

Table (5.1) List of PDEs used for the numerical experiments

name	PDE	structure
1DP	1D Poisson equation $-u'' = f, u = 0$ on $\partial\Omega$	dense, centrosym
2DP	2D Poisson equation $-\Delta u = f, u = 0$ on $\partial\Omega$	sparse, centrosym
3DP	3D Poisson equation $-\Delta u = f, u = 0$ on $\partial\Omega$	sparse, centrosym
2DPV	2D diffusion equation $-\nabla \cdot (a(x, y)\nabla u) = f, \text{ on } \partial\Omega$	sparse, centrosym depending on $a(x, y)$
2DPN	2D Neumann problem $-\Delta u + u = f, \frac{\partial u}{\partial \nu} = 0$ on $\partial\Omega$	sparse, SPD, centrosym
2DPR	2D Poisson equation with Robin BC $-\Delta u = f, \frac{\partial u}{\partial \nu} + a(x, y)u = 0$ on $\partial\Omega$	sparse, nearly centrosym depending on $a(x, y)$
1DS	1D Singular perturbation problem $-\epsilon u_{xx} + u_x = f$ on $\partial\Omega$	dense, nearly skew-centrosym depending on ϵ
1DB	1D biharmonic equation $u'''' = f, u = u' = 0$ on $\partial\Omega$	dense, centrosym
2DB	2D biharmonic equation $\Delta^2 u = f, u = \frac{\partial u}{\partial \nu} = 0$ on $\partial\Omega$	sparse+dense, centrosym
2DBV	2D biharmonic with variable coefficient $\Delta(a(x, y)\Delta u) = f, u = \frac{\partial u}{\partial \nu} = 0$ on $\partial\Omega$	sparse+dense, centrosym
2DH	2D Helmholtz equation $-(\Delta + k^2)u = f, u = 0$ on $\partial\Omega$	sparse, centrosym
3DH	3D Helmholtz equation $-(\Delta + k^2)u = f, \frac{\partial u}{\partial \nu} + a(x, y)u = 0$ on $\partial\Omega$	sparse, nearly centrosym depending on $a(x, y)$

5.1. Direct solvers with double-cone factorization. We now apply the direct solver to the centrosymmetric spectral differentiation matrices of Table 5.1. To fully exploit the structure, as previously described, we have written our own MATLAB code for the LU factorization, the Cholesky factorization, backward substitution, and forward substitution. Before running the solver, we apply equilibration (Algorithm 4) to improve the condition number of the matrix.

In Figure 5.1 we illustrate the effect of the equilibration. We (experimentally) observe that the condition number of the equilibrated 1D Poisson discrete operator seems to decrease from $O(N^4)$ to $O(N^2)$, where N is the number of collocation nodes. For the 1D biharmonic operator the condition number of the equilibrated matrix appears to be $O(N^4)$, whereas that of the original matrix is $O(N^8)$. Thus, there seems to be a significant improvement in the conditioning.

In Table 5.2 we present the relative error from using Algorithm 3 for centrosymmetric linear systems arising from the Chebyshev collocation method for 2D and 3D Poisson, 2D diffusion, and 3D Helmholtz equations. For 2DP (in the notation used in the table) we choose f such that $u = \sin(w\pi x)\sin(w\pi y)$ is the exact solution of the PDE with homogeneous boundary condition and $w = 10$. For 2DPV, we set the same exact solution with $a(x, y) = 1 + w^2x^2y^2$. In 3DP and 3DH we set f such that $u = \sin(w\pi x)\sin(w\pi y)\sin(w\pi z)$ is the exact solution, where $w = 3$ is the wave number and $k = w^2$ for 3DH. In Table 5.2 δ_{XY} denotes the relative errors from solving the linear system $\mathcal{A}z = b$ using the XY factorization:

$$\delta_{XY} = \frac{\|z - \hat{z}\|_{\infty}}{\|z\|_{\infty}},$$

where \hat{z} is the computed solution.

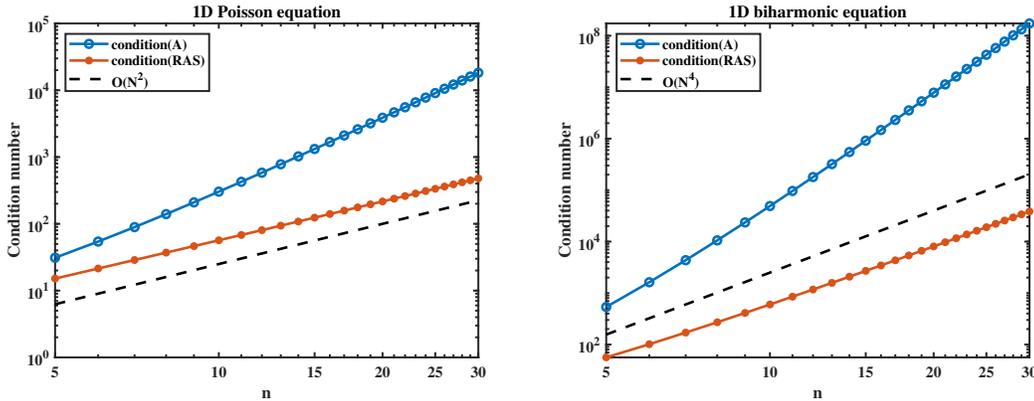


Fig. (5.1) Condition number of equilibrated 1D second- (left) and fourth-order (right) Chebyshev spectral differentiation matrices compared with the original matrices.

Table (5.2) The relative error using the XY solver for centrosymmetric linear systems arising from Chebyshev collocation method for 2D and 3D Poisson (2DP, 3DP), 2D diffusion (2DPV) and 3D Helmholtz (3DH) equations.

PDE	N	size(\mathcal{A})	$\kappa(\mathcal{A})$	δ_{XY}
2DP	101	10000	4.13×10^6	5.15×10^{-14}
2DPV	101	10000	1.21×10^8	6.63×10^{-14}
3DP	26	15625	2.30×10^4	2.06×10^{-11}
3DH	26	15625	1.65×10^5	2.07×10^{-11}

In Table 5.3 we show the results for SPD centrosymmetric linear systems arising from symmetric Legendre collocation methods for 2D Poisson equation (2DPSym) and 2D Neumann problem. Here we use the XX^T factorization. For 2DPSym we choose f such that $u = \sin(w\pi x) \sin(w\pi y)$ is the exact solution of the PDE with homogeneous boundary condition and $w = 10$. To maintain the symmetry and positive definiteness, we use a quadrature formula with Legendre collocation nodes for the weak form of the PDE. Similarly to 2DPSym, we discretize the problem with Neumann boundary conditions. We set f such that $u = (1 - x^2)^2 \cos(w\pi y)$ is the exact solution of the PDE with Neumann boundary condition where $w = 10$.

Table (5.3) The relative error using the XX^T solver for SPD centrosymmetric linear systems arising from symmetric Legendre collocation methods for 2D Poisson equation (2DPSym) and 2D Neumann problem (2DPN).

PDE	N	size(\mathcal{A})	$\kappa(\mathcal{A})$	δ_{XX^T}
2DPSym	121	14400	4.44×10^6	1.68×10^{-13}
2DPN	119	14400	1.01×10^6	4.65×10^{-10}

5.2. Skew-centrosymmetric linear systems. Consider the singular perturbation problem

$$\begin{aligned}
 (5.1) \quad & -\epsilon u_{xx} + u_x = f(x), & \text{in } \Omega = (-1, 1), \\
 & u = 0 & \text{on } \partial\Omega,
 \end{aligned}$$

where ϵ is a small positive constant. This is an advection-dominated PDE and it is particularly difficult to numerically solve if ϵ is very small. A boundary layer appears along the right side of the domain, whose thickness is tied to the magnitude of ϵ .

Let D and D^2 represent the first and second-order spectral differentiation matrices. To enforce homogeneous boundary conditions, we remove their first and last rows and columns and denote the resulting matrices by $[[D]]$ and $[[D^2]]$, respectively. Then, spectral discretization leads to a linear system involving the matrix $\mathcal{A}_{1DS} = -\epsilon[[D^2]] + [[D]]$. A relatively large value of N , number of collocation nodes, is necessary in order to resolve the boundary layer.

To assess the performance of a direct solver for a skew-centrosymmetric system, a spectral discretization of (5.1) is considered, where ϵ is very small, so the associated linear system is nearly skew-centrosymmetric. Then the linear system is solved by the direct solver described in Section 2. The results of applying XY direct solver are given in Table 5.4 for different values of ϵ . In our numerical experiments, we set f such that $u = (1+x)(1 - e^{(x-1)/\sqrt{\epsilon}})$ is the exact solution of the PDE with Dirichlet boundary conditions.

Table (5.4) The relative error using the XY solver for linear systems arising from using Chebyshev collocation methods for the 1D singular perturbation problem (1DS).

ϵ	size(\mathcal{A})	$\kappa(\mathcal{A})$	δ_{XY}
10^{-6}	1200	3.27×10^5	4.59×10^{-12}
10^{-7}	1500	3.59×10^5	6.12×10^{-12}

5.3. Iterative refinement with mixed precision. We have implemented Algorithm 5 for spectral collocation methods applied to the 1D and 2D biharmonic and 2D variable-coefficient biharmonic equation with homogeneous boundary conditions. The corresponding matrices are dense and ill-conditioned.

In Table 5.5 we show the result of applying Algorithm (5) when $u_f =$ single precision and $u = u_r =$ double precision. We have computed the XY factorization in single precision and the residual in double precision. We set the convergence criterion for the refinement process in step 8 of Algorithm (5) as $\hat{\delta} \leq \text{tol}$ where

$$\hat{\delta} = \frac{\|b - \mathcal{A}z_{i+1}\|_{\infty}}{\|\mathcal{A}\|_{\infty}\|z_{i+1}\|_{\infty} + \|b\|_{\infty}},$$

is the normwise backward error, see [15, Theorem 7.1]. We set $\text{tol} = nu$ in our experiments, where n is the dimension of the matrix. The tolerance in GMRES in each iteration of the refinement is set to be $\text{tol}_{\text{gmres}} = 10^{-2}$ and $\text{tol}_{\text{gmres}} = 10^{-4}$ respectively for the single and double precision X and Y factors.

For 1DB and 2DB we set f such that $u = 1 + \cos(\pi x)$ and $u = (1 + \cos(\pi x))(1 - 2y^2 + y^4)$ are the exact solutions of the PDEs, respectively. For 2DBV we set f such that $u = \sin^2(\pi x) \sin^2(\pi y)$ is the exact solution of the PDE, where $a(x, y) = 1 + kx^2y^2$, and $k = 1000$. This results in an ill-conditioned centrosymmetric linear system.

For 2DBV problems, the relative errors are larger than for the other examples, at approximately 10^{-9} . The stopping criterion for the iterative refinement is $nu \approx 10^{-14}$ in this case and the relative error is expected to be bounded by the product of the condition number of the equilibrated preconditioned matrix and the relative residual.

We note that we have also applied Algorithm 5 with the double-precision XY factorization and observed that in this case, the refinement algorithm converges in one iteration for all cases. Comparing this with the results in Table 5.5, we observe that for condition numbers up to

Table (5.5) The result of applying Algorithm 5 to different centrosymmetric matrices. The parameters n_{ird} and n_{irs} denote the number of iterations in the refinement using equilibration where the X and Y factors are calculated in double and single precision, respectively. δ_{ird} and δ_{irs} denote the relative errors in Algorithm 5 where the X and Y factors are calculated in double and single precision, respectively. δ_s denotes the relative error in the direct solver where the X and Y factors are calculated in single precision.

\mathcal{A}	size(\mathcal{A})	$\kappa(\mathcal{A})$	$\kappa(RAS)$	n_{ird}	δ_{ird}	δ_s	n_{irs}	δ_{irs}
1DB	20	7.82×10^6	8.05×10^3	1	1.99×10^{-14}	2.00×10^{-5}	2	7.09×10^{-14}
1DB	22	1.61×10^7	1.16×10^4	1	7.32×10^{-14}	1.59×10^{-4}	2	4.27×10^{-14}
2DB	324	3.06×10^6	8.21×10^3	1	3.74×10^{-14}	3.52×10^{-5}	2	3.46×10^{-14}
2DB	400	6.74×10^6	1.22×10^4	1	6.64×10^{-14}	9.24×10^{-5}	2	1.22×10^{-13}
2DBV	729	6.87×10^9	2.28×10^4	1	1.02×10^{-13}	2.32×10^{-5}	1	1.68×10^{-9}
2DBV	900	1.48×10^{10}	3.25×10^4	1	2.25×10^{-13}	7.43×10^{-5}	1	6.04×10^{-9}

approximately 10^{10} , Algorithm 5 with a single-precision factorization achieves the desired convergence in almost the same number of iterations as using a double-precision factorization.

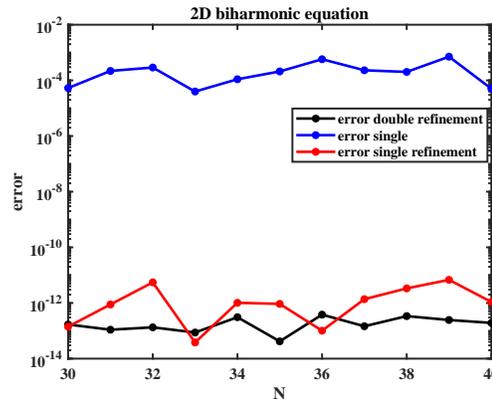


Fig. (5.2) Accuracy of Algorithm 5 compared with direct solvers using X and Y factors calculated in single and double precision for 2D biharmonic.

In Figure 5.2, the relative errors of Algorithm 5, where the X and Y factors are calculated in double and single precision are compared. With an extra refinement step the single precision calculation yields a relative error similar to the double precision calculation.

6. Concluding remarks. The double-cone factorization [14] ensures that symmetry with respect to the center is retained and yields an efficient and robust algorithm. As such, it preserves structure while maintaining the computational cost and memory required to solve the linear system. The bound we have derived on the relative error in solving centrosymmetric linear systems shows the numerical stability of the factorization and the modified substitution.

Equilibration and mixed precision further improve performance of the numerical solution procedure in cases where the matrix is ill-conditioned.

Several interesting questions remain open for further exploration. One of them is how to fully exploit the savings that arise from using the double-cone structure; each of the double-cone factors requires only half of the storage of a dense matrix.

Exploiting the double-cone structure has potential important implications also on backward error, as it may allow for structured perturbations, the advantages of which have been discussed in [21, 22], for example.

The numerical code to solve the problems described in this paper is available at <http://tinyurl.com/2e2paedr>

Acknowledgment. The authors thank the referees for their careful reading and valuable suggestions. This work was partially funded by the Natural Sciences and Engineering Research Council of Canada (NSERC).

REFERENCES

- [1] I. T. ABU-JEIB, *Algorithms for centrosymmetric and skew-centrosymmetric matrices*, Missouri Journal of Mathematical Sciences, 18 (2006).
- [2] A. L. ANDREW, *Eigenvectors of certain matrices*, Linear Algebra Appl., 7 (1973), pp. 151–162.
- [3] A. L. ANDREW, *Solution of equations involving centrosymmetric matrices*, Technometrics, 15 (1973), pp. 405–407.
- [4] K. BURNIK, *A structure-preserving QR factorization for centrosymmetric real matrices*, Linear Algebra Appl., 484 (2015), pp. 356–378.
- [5] C. CANUTO, M. Y. HUSSAINI, A. QUARTERONI, AND T. A. ZANG, *Spectral methods, fundamentals in single domains*, Scientific computation, Springer-Verlag, Berlin, 2006.
- [6] E. CARSON AND N. J. HIGHAM, *A new analysis of iterative refinement and its application to accurate solution of ill-conditioned sparse linear systems*, SIAM J. Sci. Comput., 39 (2017), pp. A2834–A2856.
- [7] R. H.-F. CHAN AND X.-Q. JIN, *An introduction to iterative Toeplitz solvers*, vol. 5 of Fundamentals of Algorithms, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2007.
- [8] L. DATTA AND S. D. MORGERA, *On the reducibility of centrosymmetric matrices, applications in engineering problems*, Circuits Systems Signal Process., 8 (1989), pp. 71–96.
- [9] J. W. DEMMEL, *Applied numerical linear algebra*, Society for Industrial and Applied Mathematics, Philadelphia, 1997.
- [10] M. EL-MIKKAWY AND F. ATLAN, *On solving centrosymmetric linear systems*, Applied mathematics (Irvine, Calif.), 4 (2013), pp. 21–32.
- [11] P. GAUDREAU AND H. SAFOUHI, *Centrosymmetric matrices in the sinc collocation method for Sturm-Liouville problems*, EPJ Web of Conferences, Mathematical Modeling and Computational Physics (MMCP 2015), 108 (2016), p. 01004.
- [12] G. H. GOLUB AND C. F. VAN LOAN, *Matrix computations*, Johns Hopkins Studies in the Mathematical Sciences, Johns Hopkins University Press, Baltimore, MD, fourth ed., 2013.
- [13] I. J. GOOD, *The inverse of a centrosymmetric matrix*, Technometrics, 12 (1970), pp. 925–928.
- [14] C. GREIF, S. NATAJ, AND M. TRUMMER, *Incomplete double-cone factorizations of centrosymmetric matrices arising in spectral methods*, Numerical algorithms, 95 (2024), pp. 1359–1386.
- [15] N. J. HIGHAM, *Accuracy and stability of numerical algorithms*, Society for Industrial and Applied Mathematics, Philadelphia, second ed., 2002.
- [16] M. KIMURA, *Some problems of stochastic processes in genetics*, Ann. Math. Statist., 28 (1957), pp. 882–901.
- [17] P. A. KNIGHT, D. RUIZ, AND B. UÇAR, *A symmetry preserving algorithm for matrix scaling*, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 931–955.
- [18] P. LV AND B. ZHENG, *Perturbation analysis for the QX factorization for centrosymmetric matrices*, Linear and multilinear algebra, (2020), pp. 1–24.
- [19] D. S. MACKEY, N. MACKEY, AND D. M. DUNLAVY, *Structure preserving algorithms for perplectic eigenproblems*, Electron. J. Linear Algebra, 13 (2005), pp. 10–39.
- [20] D. S. MACKEY, N. MACKEY, AND F. TISSEUR, *Structured tools for structured matrices*, Electron. J. Linear Algebra, 10 (2003), pp. 106–145.
- [21] S. M. RUMP, *Structured perturbations part I: Normwise distances*, SIAM Journal on Matrix Analysis and Applications, 25 (2003), pp. 1–30.
- [22] ———, *Structured perturbations part II: Componentwise distances*, SIAM Journal on Matrix Analysis and Applications, 25 (2003), pp. 31–56.
- [23] J. SHEN, T. TANG, AND L.-L. WANG, *Spectral methods: algorithms, analysis and applications*, vol. 41 of Springer series in computational mathematics, Springer-Verlag, Berlin, Heidelberg, 1. Aufl. ed., 2011.
- [24] A. STEELE, J. YALIM, AND B. WELFERT, *QX factorization of centrosymmetric matrices*, Appl. Numer. Math., 134 (2018), pp. 11–16.
- [25] D. TAO AND M. YASUDA, *A spectral characterization of generalized real symmetric centrosymmetric and generalized real symmetric skew-centrosymmetric matrices*, SIAM journal on matrix analysis and applications, 23 (2002), pp. 885–895.

- [26] W. F. TRENCH, *Characterization and properties of matrices with generalized symmetry or skew symmetry*, Linear Algebra Appl., 377 (2004), pp. 207–218.
- [27] J. R. WEAVER, *Centrosymmetric (cross-symmetric) matrices, their basic properties, eigenvalues, and eigenvectors*, Amer. Math. Monthly, 92 (1985), pp. 711–717.
- [28] J. H. WILKINSON, *Rounding errors in algebraic processes*, Prentice-Hall, Inc., Englewood Cliffs, N.J., 1963.