

# Nonsingular systems of generalized Sylvester equations: an algorithmic approach<sup>\*</sup>

Fernando De Terán<sup>1</sup>, Bruno Iannazzo<sup>2</sup>, Federico Poloni<sup>3</sup>, and Leonardo Robol<sup>4,5</sup>

<sup>1</sup>Departamento de Matemáticas, Universidad Carlos III de Madrid, Avda. Universidad 30, 28911 Leganés, Spain. [fteran@math.uc3m.es](mailto:fteran@math.uc3m.es).

<sup>2</sup>Dipartimento di Matematica e Informatica, Università di Perugia, Via Vanvitelli 1, 06123 Perugia, Italy. [bruno.iannazzo@dmf.unipg.it](mailto:bruno.iannazzo@dmf.unipg.it).

<sup>3</sup>Dipartimento di Informatica, Università di Pisa, Largo B. Pontecorvo 3, 56127 Pisa, Italy. [federico.poloni@unipi.it](mailto:federico.poloni@unipi.it).

<sup>4</sup>Dipartimento di Matematica, Università di Pisa, Largo B. Pontecorvo 5, 56127 Pisa, Italy. [leonardo.robol@unipi.it](mailto:leonardo.robol@unipi.it).

<sup>5</sup>Institute of Information Science and Technologies “A. Faedo”, ISTI-CNR, Via G. Moruzzi, 1, 56124 Pisa, Italy.

## Abstract

We consider the uniqueness of solution (i.e., nonsingularity) of systems of  $r$  generalized Sylvester and  $\star$ -Sylvester equations with  $n \times n$  coefficients. After several reductions, we show that it is sufficient to analyze periodic systems having, at most, one generalized  $\star$ -Sylvester equation. We provide characterizations for the nonsingularity in terms of spectral properties of either matrix pencils or formal matrix products, both constructed from the coefficients of the system. The proposed approach uses the periodic Schur decomposition, and leads to a backward stable  $O(n^3 r)$  algorithm for computing the (unique) solution.

**Keywords:** Sylvester and  $\star$ -Sylvester equations, systems of linear matrix equations, matrix pencils, periodic Schur decomposition, periodic QR/QZ algorithm, formal matrix product

## 1 Introduction

The *generalized Sylvester* equation

$$AXB - CXD = E, \tag{1}$$

goes back to, at least, the early 20th century [35]. Here the unknown  $X$ , the coefficients  $A, B, C, D$ , and the right-hand side  $E$  are complex matrices of appropriate size. This equation

---

<sup>\*</sup>2010 *Mathematics Subject Classification*. Primary 15A22, 15A24, 65F15. This work was partially supported by the Ministerio de Economía y Competitividad of Spain through grants MTM2015-68805-REDT, and MTM2015-65798-P (F. De Terán), by an INdAM/GNCS Research Project 2016 (B. Iannazzo, F. Poloni, and L. Robol), and by the Research project of the Università di Perugia “Soluzione numerica di problemi di algebra lineare strutturata” (B. Iannazzo). Part of this work was done during a visit of the first author to the Università di Perugia as a Visiting Researcher.

has attracted much attention since the 1970s, mainly due to its appearance in applied problems (see, for instance, [11, 25, 27, 30, 32]).

Another related equation, whose interest is growing recently (see, for instance, [10, 13–16, 19]), arises when introducing the  $\star$  operator in the second appearance of the unknown. This equation is the *generalized  $\star$ -Sylvester equation*

$$AXB - CX^{\star}D = E, \quad (2)$$

where the unknown  $X$ , the coefficients  $A, B, C, D$ , and the right-hand side  $E$  are again complex matrices of appropriate size, and  $\star$  can be either the transpose (T) or the conjugate transpose (H) operator. When  $\star = \text{T}$ , the equation can be seen as a linear system in the entries of the unknown  $X$ , while if  $\star = \text{H}$ , the equation is no more linear in the entries of  $X$  because of conjugation. Nevertheless, with the usual isomorphism  $\mathbb{C} \cong \mathbb{R}^2$ , obtained by splitting the real and imaginary parts, it turns out to be a linear system with respect to the real entries of  $\text{re}(X)$  and  $\text{im}(X)$ .

One could argue that, in some sense, solving generalized Sylvester and  $\star$ -Sylvester equations is an elementary problem both from the theoretical and the computational point of view, since they are equivalent to linear systems. Nevertheless, there has been great interest in giving conditions on the existence and uniqueness of solutions based just on properties of certain small-sized matrix pencils constructed from the coefficients. For instance, when all coefficients are square, it is known that (1) has a unique solution if and only if the two pencils  $A - \lambda C$  and  $D - \lambda B$  have disjoint spectra [11, Th. 1], whereas the uniqueness of solutions of (2) depends on spectral properties of the matrix pencil  $\begin{bmatrix} \lambda D^{\star} & -B^{\star} \\ A & -\lambda C \end{bmatrix}$  (see [15, Th. 15]).

On the other hand, if all coefficients are square and of size  $n$ , then the resulting linear system has size  $n^2$  or  $2n^2$ . From the computational point of view, solving a linear system of size  $n^2$  with standard (non-structured) algorithms may be prohibitive, since they result in a method which approximates the solution in  $O(n^6)$  (floating point) arithmetic operations (flops). However, dealing with the coefficients it is possible to get algorithms requiring only  $O(n^3)$  flops, such as the one given in [11].

Recently, systems of coupled generalized Sylvester and  $\star$ -Sylvester equations have been considered, and useful conditions on the existence of solutions have been derived in [18]. Here, we consider the same kind of systems and provide further characterizations for the uniqueness of their solution, for any right-hand side, based on certain spectral conditions on their coefficients. It is worth to emphasize that, while in [18] non-square coefficients are allowed, as long as the matrix products are well-defined, here we assume that all coefficients, as well as the unknowns, are square of size  $n \times n$ . This choice has been made because the problem of nonsingularity, even for just one equation, presents certain additional subtleties when the coefficients are not square or they are square with different sizes (see [16]). In the assumption that all coefficients are square and of size  $n \times n$ , such a system of matrix equations is equivalent to a square linear system, which has a unique solution, for any right-hand side, if and only if the coefficient matrix is nonsingular. For this reason, we will use the term *nonsingular system* as a synonym of a system having a unique solution (for any right-hand side).

The *systems of generalized Sylvester and  $\star$ -Sylvester equations* that we consider are of the form

$$A_k X_{\alpha_k}^{s_k} B_k - C_k X_{\beta_k}^{t_k} D_k = E_k, \quad k = 1, \dots, r, \quad (3)$$

where all matrices involved are complex and of size  $n \times n$ , the indices  $\alpha_i, \beta_i$  of the unknowns are positive integers and can be equal or different to each other, and  $s_i, t_i \in \{1, \star\}$ .

Our approach starts by reducing the problem on the nonsingularity of (3) to the special case

of *periodic* systems of the form

$$\begin{cases} A_k X_k B_k - C_k X_{k+1} D_k &= E_k, & k = 1, \dots, r-1, \\ A_r X_r B_r - C_r X_1^s D_r &= E_r, \end{cases} \quad (4)$$

where  $s \in \{1, \star\}$ . We provide an explicit characterization of nonsingularity only for periodic systems like (4). However, our reduction allows one to get a characterization for any system like (3) after undoing all changes that take the system (3) into (4). Since these systems can be seen as linear systems with a square matrix coefficient, the criteria for nonsingularity do not depend on the right-hand sides  $E_k$ , but only on the coefficients  $A_k, B_k, C_k, D_k$ , for  $k = 1, \dots, r$ .

Periodic systems of Sylvester equations naturally arise in the context of discrete-time periodic systems, and they have been analyzed by several authors (see, for instance, [1, 20, 21, 33]). Prior to our work, Byers and Rhee provided in the unpublished work [9] a characterization for the nonsingularity of (4) with  $s = 1$ , together with an  $O(n^3 r)$  algorithm to compute the solution.

The first contribution of the present work is the reduction of a nonsingular system of Sylvester and  $\star$ -Sylvester equations (3) to several disjoint systems of periodic type (4), where all equations are generalized Sylvester, with the exception of the last one that may be either a generalized Sylvester or a generalized  $\star$ -Sylvester equation. We note that neither the coefficients, nor the number of equations in the original and the reduced system necessarily coincide.

As a second contribution, we provide a characterization for the nonsingularity of (4) for  $s = \mathbf{H}, \mathbf{T}$  (i. e.,  $s = \star$ , according to our notation). This characterization appears in two different formulations. The first one is given in terms of the spectrum of *formal products* constructed from the coefficients of the system (we include the case  $s = 1$ , treated in Theorem 5, and the case  $s = \star$ , treated in Theorem 6). The second formulation, valid for  $s = \star$ , is given in terms of spectral properties of a block-partitioned  $(2rn) \times (2rn)$  matrix pencil constructed in an elementary way from the coefficients (Theorem 7). This characterization extends the one in [15] for the single equation (2), and it is in the same spirit as the one in [9] for periodic systems with  $s = 1$ .

The third contribution of the paper is to provide an  $O(n^3 r)$  algorithm to compute the unique solution of a nonsingular system. Our algorithm is a Bartels-Stewart like algorithm, based on the periodic Schur form [7]. It extends the one in [9] for systems of Sylvester equations only, the one in [13] for the  $\star$ -Sylvester equation  $AX + X^\star D = E$ , and the one outlined in [10, §4.2] for (2).

We note that extending the results of [9] to include  $\star$ -Sylvester equations is not a trivial endeavour: the presence of transpositions creates additional dependencies between the data, hence we need a different strategy to reduce the coefficients to a triangular form, and the resulting criteria have a significantly different form.

Throughout the manuscript,  $\mathbf{i}$  denotes the imaginary unit, that is,  $\mathbf{i}^2 = -1$ . By  $M^{-\star}$  we denote the inverse of the invertible matrix  $M^\star$ , with  $\star = \mathbf{H}, \mathbf{T}$ . A pencil  $\mathcal{Q}(\lambda)$  is *regular* if it is square and  $\det \mathcal{Q}(\lambda)$  is not identically zero. We use the symbol  $\Lambda(\mathcal{Q})$  to denote the *spectrum* of a regular matrix pencil  $\mathcal{Q}(\lambda)$ , that is the set of values  $\lambda$  such that  $\mathcal{Q}(\lambda)$  is singular (including  $\infty$  if the degree of  $\det \mathcal{Q}(\lambda)$  is smaller than the size of the pencil). For simplicity, we use the term *system of Sylvester-like equations* for a system of generalized Sylvester and  $\star$ -Sylvester equations.

The paper is organized as follows. In Section 2 we present some applications of systems of Sylvester and  $\star$ -Sylvester equations; in Section 3 the periodic Schur decomposition and the concept of formal matrix product are recalled. Section 4 hosts the main theoretical results of the paper, whose proofs are deferred to Section 7, after Sections 5 and 6, that are devoted to some successive simplifications of the problem which are useful for the proofs. Section 8 is devoted to describe and analyze an efficient algorithm for the solution of systems of Sylvester-like equations. Finally, in Section 9 we draw some conclusions.

## 2 Applications

Sylvester-like equations appear in various fields of applied mathematics. In some cases, the applications have natural “periodic extensions”, where systems of these equations come into play.

As an example, consider a  $2 \times 2$  block upper triangular matrix  $M = \begin{bmatrix} A & C \\ 0 & B \end{bmatrix}$ , and assume that we want to block diagonalize it, setting  $C$  to zero with a similarity transformation. This problem arises, for instance, when  $M$  is the block Schur form of a given matrix and we want to decouple the action of the parts of the spectrum contained in  $A$  and  $B$ . Then, we can look for a matrix  $V$  such that

$$V^{-1}MV = \begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix}, \quad V = \begin{bmatrix} I & X \\ 0 & I \end{bmatrix}. \quad (5)$$

This problem can be solved by finding a solution to the Sylvester equation  $AX - XB + C = 0$ , and admits a natural extension in periodic form, when we want to block diagonalize the product of  $2 \times 2$  block upper triangular matrices, as the one arising in a periodic Schur form. We start from

$$M = M_1 \cdots M_r, \quad M_i := \begin{bmatrix} A_i & C_i \\ 0 & B_i \end{bmatrix}, \quad (6)$$

where the blocks have the same size for each  $i$ , and we want to block diagonalize  $M$ . For stability reasons, rather than working directly on the product  $M$ , it is often preferable to look for matrices  $V_i$  such that  $V_i^{-1}M_iV_{i+1}$  are all block diagonal, with  $V_{r+1} = V_1$  (see, e.g. [34]). If we impose on  $V_i = \begin{bmatrix} I & X_i \\ 0 & I \end{bmatrix}$  the same block upper triangular structure we had for  $V$  in (5), then we obtain the periodic system of Sylvester equations  $A_iX_{i+1} - X_iB_i + C_i = 0$ , for  $i = 1, \dots, r$ , with  $X_{r+1} = X_1$ .

Similarly, decoupling saddle-point matrices (as quadratic forms) given in product form

$$N = N_1N_2N_3 = \begin{bmatrix} A_1 & 0 \\ C_1 & B_1 \end{bmatrix} \begin{bmatrix} 0 & A_2 \\ B_2 & C_2 \end{bmatrix} \begin{bmatrix} B_3 & C_3 \\ 0 & A_3 \end{bmatrix} = \begin{bmatrix} 0 & A_1A_2A_3 \\ B_1B_2B_3 & C_1A_2A_3 + B_1C_2A_3 + B_1B_2C_3 \end{bmatrix}$$

(see e.g. [31] for similar factorizations) naturally leads to systems of  $\star$ -Sylvester equations: one can choose the following change of bases to eliminate the blocks  $C_i$

$$U_1^*NU_1 = U_1^* \begin{bmatrix} A_1 & 0 \\ C_1 & B_1 \end{bmatrix} V_2^{-1}V_2 \begin{bmatrix} 0 & A_2 \\ B_2 & C_2 \end{bmatrix} V_3V_3^{-1} \begin{bmatrix} B_3 & C_3 \\ 0 & A_3 \end{bmatrix} U_1,$$

$$U_1 = \begin{bmatrix} I & X_1 \\ 0 & I \end{bmatrix}, \quad V_2 = \begin{bmatrix} I & 0 \\ X_2 & I \end{bmatrix}, \quad V_3 = \begin{bmatrix} I & X_3 \\ 0 & I \end{bmatrix};$$

then the factors become

$$\begin{aligned} U_1^* \begin{bmatrix} A_1 & 0 \\ C_1 & B_1 \end{bmatrix} V_2^{-1} &= \begin{bmatrix} A_1 & 0 \\ X_1^*A_1 - B_1X_2 + C_1 & B_1 \end{bmatrix}, \\ V_2 \begin{bmatrix} 0 & A_2 \\ B_2 & C_2 \end{bmatrix} V_3 &= \begin{bmatrix} 0 & A_2 \\ B_2 & X_2A_2 + B_2X_3 + C_2 \end{bmatrix}, \\ V_3^{-1} \begin{bmatrix} B_3 & C_3 \\ 0 & A_3 \end{bmatrix} U_1 &= \begin{bmatrix} B_3 & -X_3A_3 + B_3X_1 + C_3 \\ 0 & A_3 \end{bmatrix}. \end{aligned}$$

Hence the blocks in the position of the  $C_i$  vanish if the  $X_i$  solve the periodic system of  $\star$ -Sylvester equations

$$\begin{cases} X_1^*A_1 - B_1X_2 + C_1 = 0, \\ X_2A_2 + B_2X_3 + C_2 = 0, \\ -X_3A_3 + B_3X_1 + C_3 = 0. \end{cases}$$

Another relevant application is the reordering of periodic Schur forms. In order to swap the diagonal blocks of  $M$  in (6) it may be convenient to swap the blocks of the factors  $M_i$ , for  $i = 1, \dots, r$ . While the problem of swapping the blocks of  $M$  can be reduced to a Sylvester equation [4], the problem of swapping the blocks of the factors can be reduced to a periodic system of Sylvester equations. Indeed, swapping diagonal entries of matrices given in products form, without forming the product, is an essential step in the eigenvector recovery procedures of some fast methods for matrix polynomial eigenvalue problems (see [2, 3]).

### 3 Periodic Schur decomposition of formal matrix products

In order to state and prove the nonsingularity results for a system of Sylvester-like equations and to design an efficient algorithm to compute the solution, we need to introduce several results and definitions that extend the ideas of matrix pencils and generalized eigenvalues to products of matrices of an arbitrary number of factors. These are standard tools in the literature (see, for instance, [20, 21]).

**Theorem 1** (Periodic Schur decomposition [7]). *Let  $M_k, N_k$ , for  $k = 1, \dots, r$ , be two sequences of  $n \times n$  complex matrices. Then there exist unitary matrices  $Q_k, Z_k$ , for  $k = 1, \dots, r$ , such that*

$$Q_k^H M_k Z_k = T_k, \quad Q_k^H N_k Z_{k+1} = R_k, \quad k = 1, \dots, r \quad (7)$$

where  $T_k, R_k$  are upper triangular and  $Z_{r+1} = Z_1$ .

If the matrices  $N_k$  are invertible, Theorem 1 means that we can apply suitable unitary changes of bases to the product

$$\Pi = N_r^{-1} M_r N_{r-1}^{-1} M_{r-1} \cdots N_1^{-1} M_1 \quad (8)$$

to make all its factors upper triangular simultaneously. More precisely,

$$Z_1^{-1} \Pi Z_1 = R_r^{-1} T_r R_{r-1}^{-1} T_{r-1} \cdots R_1^{-1} T_1.$$

In this case, the eigenvalues of  $\Pi$  are

$$\lambda_i = \frac{(T_1)_{ii}(T_2)_{ii} \cdots (T_r)_{ii}}{(R_1)_{ii}(R_2)_{ii} \cdots (R_r)_{ii}}, \quad i = 1, 2, \dots, n. \quad (9)$$

Even when some of the  $N_k$  matrices are not invertible, we call the expression (8) a *formal matrix product*, and (7) a *formal periodic Schur form* of the product. If  $(T_1)_{ii}(T_2)_{ii} \cdots (T_r)_{ii} = (R_1)_{ii}(R_2)_{ii} \cdots (R_r)_{ii} = 0$ , for some  $i \in \{1, 2, \dots, n\}$ , we call the formal product *singular*; otherwise, we call it *regular*. If  $\Pi$  is regular, it makes sense to consider the ratios  $\lambda_i$  defined in (9), with the convention that  $\frac{a}{0} = \infty$  for  $a \neq 0$ . We call these ratios the *eigenvalues* of the regular formal matrix product  $\Pi$ . The set of eigenvalues of  $\Pi$  is called, as usual, the *spectrum* of  $\Pi$ , and we denote it by  $\Lambda(\Pi)$ .

We also define the eigenvalues of a formal matrix product of the form

$$\tilde{\Pi} = M_r N_{r-1}^{-1} M_{r-1} \cdots N_1^{-1} M_1 N_r^{-1}$$

(i. e., one in which the exponent  $-1$  appears in the factors in *even* positions) by the same formula (9).

*Remark 2.* For the notion of eigenvalues of formal products to be well defined, one should prove that it does not depend on the choice of the (non-unique) decomposition (7). If all  $N_i$  matrices

are nonsingular, then this is evident because they coincide with the eigenvalues obtained by performing the inversions and computing the actual product  $\Pi$ . If some of the  $N_i$  are singular, then we can use a continuity argument to show that the  $\lambda_i$  are the limits, as  $\varepsilon \rightarrow 0$ , of the eigenvalues of

$$(N_r + \varepsilon P_r)^{-1} M_r (N_{r-1} + \varepsilon P_{r-1})^{-1} M_{r-1} \cdots (N_1 + \varepsilon P_1)^{-1} M_1$$

for each choice of the nonsingular matrices  $P_1, P_2, \dots, P_r$  that make the factors  $N_k + \varepsilon P_k$  invertible, for all  $k = 1, \dots, r$  and sufficiently small  $\varepsilon > 0$ .

**Lemma 3.** *Let  $\Pi = M_1^{-1} N_1 \cdots M_r^{-1} N_r$  be a formal matrix product. Then, the matrix pencil*

$$\mathcal{Q}(\lambda) := \begin{bmatrix} \lambda M_1 & -N_1 & & & \\ & \lambda M_2 & \ddots & & \\ & & \ddots & -N_{r-1} & \\ -N_r & & & & \lambda M_r \end{bmatrix}$$

is regular if and only if  $\Pi$  is regular. In this case, the eigenvalues of  $\mathcal{Q}(\lambda)$  are the  $r$ -th roots of the eigenvalues of  $\Pi$ , with the convention that  $\sqrt[r]{\infty} = \infty$ .

*Proof.* Let us start by considering the case when  $\Pi$  is regular with distinct (simple) eigenvalues, and all matrices  $M_i, N_i$  are invertible. Let  $\mu \in \mathbb{C}$  be an eigenvalue of  $\Pi$ , with  $v$  a corresponding right eigenvector, and let  $\lambda \in \mathbb{C}$  be such that  $\lambda^r = \mu$ . We set  $v_1 := v$ , and define

$$v_j := \lambda N_{j-1}^{-1} M_{j-1} v_{j-1}, \quad j = 2, \dots, r.$$

Then, the relation  $\Pi v = \lambda^r v$  implies  $\mathcal{Q}(\lambda) \hat{v} = 0$ , where  $\hat{v} := [v_1^\top \ v_2^\top \ \dots \ v_r^\top]^\top$ , which can be verified by a direct computation. In particular, all the  $r$ -th distinct roots of  $\mu$  are eigenvalues of  $\mathcal{Q}(\lambda)$ .

This implies that  $q(\lambda) := \det(\mathcal{Q}(\lambda)) = \det(\Pi - \lambda^r I) \cdot \det(M_1 \cdots M_r)$ , since  $\det(M_1 \cdots M_r)$  is the leading coefficient of the degree  $nr$  polynomial  $\det \mathcal{Q}(\lambda)$ . Let  $Q_k^H M_k Z_k = T_k$  and  $Q_k^H N_k Z_{k+1} = R_k$  be a periodic Schur decomposition of  $\Pi$ . Then, we may write

$$q(\lambda) = \det(T_1^{-1} R_1 \cdots T_r^{-1} R_r - \lambda^r I) \cdot \det(T_1 \cdots T_r) = \det(R_1 \cdots R_r - \lambda^r T_1 \cdots T_r),$$

where we have swapped the factors inside the determinant using the fact that all the matrices are upper triangular. Using a continuity argument like the one in Remark 2, we see that the identity  $\det(\mathcal{Q}(\lambda)) = \det(R_1 \cdots R_r - \lambda^r T_1 \cdots T_r) =: p(\lambda)$  also holds when some of the  $T_i, R_i$  are singular, and even when  $\Pi$  has multiple eigenvalues. This proves the second claim in the statement. In addition,  $p(\lambda) \equiv 0$  if and only if  $T_i, R_j$  have a common diagonal zero for some  $i, j$ . Since  $\mathcal{Q}(\lambda)$  is singular if and only if  $p(\lambda) \equiv 0$ , this concludes the proof.  $\square$

## 4 Main results

Here we state the characterizations for the nonsingularity of a periodic system of type (4) for each of the three possible cases  $s \in \{1, \mathbb{T}, \mathbb{H}\}$  (the proofs will be given in Section 7). Later, in Section 5, we will show that these characterizations are enough to get a characterization of nonsingularity of the general system (3).

We recall the following definition.



- if  $\star = \mathbf{H}$ , then  $\Lambda(\mathcal{Q})$  is  $\mathbf{H}$ -reciprocal-free, and
- if  $\star = \mathbf{T}$ , then  $\Lambda(\mathcal{Q}) \setminus \mathfrak{A}_{2r}$  is reciprocal free and the multiplicity of  $\xi$ , for any  $\xi \in \mathfrak{A}_{2r}$ , is at most 1.

The proof of Theorem 7 can be readily obtained by means of the following result combined with Theorem 6.

**Lemma 8.** *Let  $\mathcal{S}$  be a subset of  $\mathbb{C} \cup \{\infty\}$ , let  $p \in \mathbb{N}$ , and define the sets:*

$$-\mathcal{S} := \{-z \mid z \in \mathcal{S}\}, \quad \mathcal{S}^{-1} := \{z^{-1} \mid z \in \mathcal{S}\}, \quad \sqrt[p]{\mathcal{S}} := \{z \in \mathbb{C} \cup \{\infty\} \mid z^p \in \mathcal{S}\}$$

(we set  $\infty^p = \infty, -\infty = \infty$ , and  $\infty^{-1} = 0, 0^{-1} = \infty$ ). Then the following statements are equivalent:

- $\mathcal{S}$  is  $\star$ -reciprocal free.
- $-\mathcal{S}$  is  $\star$ -reciprocal free.
- $\mathcal{S}^{-1}$  is  $\star$ -reciprocal free.
- $\sqrt[p]{\mathcal{S}}$  is  $\star$ -reciprocal free.

The equivalence between claims (a) and (d) in Lemma 8 can be found in [15, Lemma 3] for  $p = 2$ . The extension to arbitrary  $p$ , as well as the other equivalences, are straightforward.

*Proof of Theorem 7.* By Lemma 3,  $\Lambda(\mathcal{Q}) = \sqrt[p]{-\Lambda(\Pi^{-1})} = \sqrt[p]{-\Lambda(\Pi)^{-1}}$ , with  $\Pi$  as in Theorem 6 (the second identity is immediate). From this, we also get  $\sqrt[p]{-\Lambda(\Pi) \setminus \{-1\}}^{-1} = \sqrt[p]{-\Lambda(\Pi)^{-1} \setminus \{1\}} = \Lambda(\mathcal{Q}) \setminus \mathfrak{A}_{2r}$ .

Now, Theorem 7 is an immediate consequence of Theorem 6 and Lemma 8.  $\square$

Theorem 7 is an extension of [15, Th. 15], where the case of a single generalized  $\star$ -Sylvester equation is treated. It also resembles the characterization obtained in [9, Th. 3] for systems of generalized Sylvester equations (i.e., without  $\star$ ). We reproduce this last result here, for completeness.

**Theorem 9** (Byers and Rhee, [9]). *The system (4), with  $s = 1$ , is nonsingular if and only if the matrix pencils*

$$\begin{bmatrix} \lambda A_1 & C_1 & & & \\ & \lambda A_2 & \ddots & & \\ & & \ddots & C_{r-1} & \\ C_r & & & \lambda A_r & \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} \lambda D_1 & B_1 & & & \\ & \lambda D_2 & \ddots & & \\ & & \ddots & B_{r-1} & \\ B_r & & & \lambda D_r & \end{bmatrix}$$

*are regular and have disjoint spectra.*

Our strategy to prove Theorems 5 and 6 for periodic systems (4) relies on several steps. First, we use the fact that the system is equivalent to a system with triangular coefficients, as shown in Section 6.1. Second, in Section 6.2, when  $s = 1$  or  $s = \mathbf{T}$ , we transform the system of matrix equations with triangular coefficients to an equivalent linear system that is block upper triangular in a suitable basis (given by an appropriate order of the unknowns). The remaining case  $s = \mathbf{H}$  is reduced to the case  $s = 1$  in Section 6.3. Third, we prove in Section 7 that the diagonal blocks of the matrix coefficient of the resulting block triangular system are invertible if and only if the conditions in the statement of Theorems 5 and 6 hold.



## 5 Reducing the problem to periodic systems

In this section, we are going to show how to reduce the problem of nonsingularity of a general system (3) to the question on nonsingularity of periodic systems (4) with at most one  $\star$  in the last equation.

### 5.1 Reduction to an irreducible system

We say that the system (3) of  $r$  equations in  $s$  unknowns is *reducible* if there are  $0 < k < s$  unknowns appearing only in  $0 < h < r$  equations and the remaining  $s - k$  unknowns appear only in the remaining  $r - h$  equations. In other words, a reducible system can be partitioned into two systems with no common unknowns. A system is said to be *irreducible* if it is not reducible.

Let  $\mathbb{S}$  be a system of  $r$  ordered equations like (3). Let  $\{1, \dots, r\} = \mathcal{I}_1 \cup \dots \cup \mathcal{I}_\ell$  be a partition of the set of indices. Then we denote by  $\mathbb{S}(\mathcal{I}_j)$ , for  $j = 1, \dots, \ell$ , the system of equations comprising the equations with indices in  $\mathcal{I}_j$ .

**Proposition 10.** *Let  $\mathbb{S}$  be a system (3) with  $r$  equations. There exists a partition  $\mathcal{I}_1 \cup \dots \cup \mathcal{I}_\ell$  of  $\{1, \dots, r\}$  such that, for each  $j = 1, \dots, \ell$ , the system  $\mathbb{S}(\mathcal{I}_j)$  is irreducible.*

*Proof.* We proceed by strong induction on  $r$ . If  $r = 1$  the system has only one equation and thus it is irreducible. Let  $r > 1$  and consider a system with  $r$  equations. If it is irreducible, then we can choose  $\ell = 1$  and  $\mathcal{I}_1 = \{1, \dots, r\}$ . Otherwise, it can be split (by definition) into two systems with indices in two disjoint nonempty index sets  $\mathcal{I}$  and  $\mathcal{J}$ , respectively, such that  $\mathcal{I} \cup \mathcal{J} = \{1, \dots, r\}$ . The systems  $\mathbb{S}(\mathcal{I})$  and  $\mathbb{S}(\mathcal{J})$  have strictly less than  $r$  equations, and therefore, relying on the inductive hypothesis, they can be split further into irreducible subsystem using the partitions

$$\mathcal{I} = \mathcal{I}_1 \cup \dots \cup \mathcal{I}_{\ell_1}, \quad \mathcal{J} = \mathcal{J}_1 \cup \dots \cup \mathcal{J}_{\ell_2}.$$

Then,  $\{1, \dots, r\} = \mathcal{I}_1 \cup \dots \cup \mathcal{I}_{\ell_1} \cup \mathcal{J}_1 \cup \dots \cup \mathcal{J}_{\ell_2}$  yields a decomposition into irreducible systems with  $\ell := \ell_1 + \ell_2$  components, and this concludes the proof.  $\square$

Proposition 10 shows that every system can be split into irreducible systems. To determine if a system is nonsingular, it is sufficient to answer the same question for its irreducible components, as stated in the following result.

**Proposition 11.** *Let  $\mathbb{S}$  be the system (3) with  $s$  matrix unknowns, and let  $\mathcal{I}_1 \cup \dots \cup \mathcal{I}_\ell$  be a partition of  $\{1, \dots, r\}$  such that each system  $\mathbb{S}(\mathcal{I}_j)$  is irreducible, for  $j = 1, \dots, \ell$ . The system  $\mathbb{S}$  is nonsingular if and only if the system  $\mathbb{S}(\mathcal{I}_j)$  is nonsingular, for each  $j = 1, \dots, \ell$ .*

*Proof.* We shall show directly that  $\mathbb{S}$  has a unique solution if and only if  $\mathbb{S}(\mathcal{I}_j)$  has a unique solution for each  $j = 1, \dots, \ell$ . Any solution of  $\mathbb{S}$  yields a solution of  $\mathbb{S}(\mathcal{I}_j)$ , for each  $j = 1, \dots, \ell$ , and viceversa. Let us assume that  $\mathbb{S}$  has two different solutions  $(X_1, \dots, X_s)$  and  $(Y_1, \dots, Y_s)$ . Then there exists some  $1 \leq p \leq s$  such that  $X_p \neq Y_p$ . If  $p \in \mathcal{I}_q$ , for some  $1 \leq q \leq \ell$ , then  $\mathbb{S}(\mathcal{I}_q)$  has two different solutions, the first one containing  $X_p$  and the second one containing  $Y_p$ . Conversely, if not every system  $\mathbb{S}(\mathcal{I}_j)$  is nonsingular, then there is some  $1 \leq q \leq \ell$  such that either  $\mathbb{S}(\mathcal{I}_q)$  is not consistent or it has two different solutions. In the first case, the whole system  $\mathbb{S}$  would not be consistent either. If  $\mathbb{S}(\mathcal{I}_q)$  has two different solutions,  $(X_1, \dots, X_{s_q})$  and  $(Y_1, \dots, Y_{s_q})$ , and  $\mathbb{S}(\mathcal{I}_j)$  is consistent, for any  $j \neq q$ , then we can construct two different solutions of  $\mathbb{S}$  by completing with  $(X_1, \dots, X_{s_q})$  and  $(Y_1, \dots, Y_{s_q})$ , respectively, a solution of the remaining  $\mathbb{S}(\mathcal{I}_j)$  for  $j \neq q$ .  $\square$

Finally, we show that for nonsingular systems, the number of equations and unknowns in each irreducible subsystem is the same.

**Proposition 12.** *Let  $\mathbb{S}$  be the system (3) with  $r$  matrix unknowns with size  $n \times n$  and let  $\mathcal{I}_1 \cup \dots \cup \mathcal{I}_\ell$  be a partition of  $\{1, \dots, r\}$  such that each system  $\mathbb{S}(\mathcal{I}_j)$  is irreducible, for  $j = 1, \dots, \ell$ . Let  $r_j$  and  $s_j$  be the number of matrix equations and unknowns, respectively, of  $\mathbb{S}(\mathcal{I}_j)$ . If the system  $\mathbb{S}$  has a unique solution then  $r_j = s_j$ , for  $j = 1, \dots, \ell$ .*

*Proof.* If an irreducible system with  $\widehat{r}$  equations and  $\widehat{s}$  unknowns has a unique solution, then  $\widehat{s} \leq \widehat{r}$ , since otherwise this system, considered as a linear system on the entries of the matrix unknowns, would have more unknowns than equations.

Now, by contradiction, assume that  $r_j \neq s_j$ , for some  $1 \leq j \leq \ell$ . Then, since  $\sum_{j=1}^{\ell} r_j = \sum_{j=1}^{\ell} s_j = r$ , there exists some  $1 \leq p \leq \ell$  such that  $r_p < s_p$ . Thus the system  $\mathbb{S}(\mathcal{I}_p)$  cannot have a unique solution, and this contradicts Proposition 11.  $\square$

The previous results show that, in order to analyze the nonsingularity of a system of  $r$  matrix equations in  $r$  matrix unknowns, we may assume that the system is irreducible.

Moreover, Proposition 11 shows that a first step to compute the unique solution of a system of type (3) consists in splitting the system into irreducible systems and solving them separately.

## 5.2 Reduction to a system where every unknown appears twice

We consider a nonsingular irreducible system of Sylvester-like equations and we want to prove that the system can be reduced to another one in which each unknown appears exactly twice (and in different equations, when the system has at least two equations). For this purpose, we need the following result.

**Theorem 13.** *Let  $\mathbb{S}$  be an irreducible system of equations in the form (3) with  $r > 1$  equations and unknowns. If the unknown  $X_{\alpha_k}$  appears in just one equation, say  $A_k X_{\alpha_k}^{s_k} B_k - C_k X_{\beta_k}^{t_k} D_k = E_k$ , then  $\mathbb{S}$  is nonsingular if and only if  $A_k$  and  $B_k$  are invertible and the system  $\widetilde{\mathbb{S}}$  formed by the remaining  $r - 1$  equations is nonsingular. Moreover  $\widetilde{\mathbb{S}}$  is irreducible.*

*Proof.* Note, first, that  $\beta_k \neq \alpha_k$ , and that the variable  $X_{\beta_k}$  appears again in  $\widetilde{\mathbb{S}}$ , otherwise  $\mathbb{S}$  would be reducible. Suppose first that  $\widetilde{\mathbb{S}}$  is nonsingular and  $A_k, B_k$  are invertible. Then, the unique solution of  $\mathbb{S}$  is obtained by first solving  $\widetilde{\mathbb{S}}$  to get the value of all the variables except  $X_{\alpha_k}$ , and then computing  $X_{\alpha_k}$  from

$$X_{\alpha_k}^{s_k} = A_k^{-1} (C_k X_{\beta_k}^{t_k} D_k + E_k) B_k^{-1}. \quad (14)$$

If  $\widetilde{\mathbb{S}}$  has more than one solution, for  $A_k$  and  $B_k$  invertible, then (14) produces multiple solutions to  $\mathbb{S}$ . If  $\widetilde{\mathbb{S}}$  has no solution, then clearly  $\mathbb{S}$  has no solution either. If  $A_k$  is singular, let  $v$  be a nonzero vector such that  $A_k v = 0$ ; then, given any solution to (3) we can replace  $X_{\alpha_k}^{s_k}$  with  $X_{\alpha_k}^{s_k} + v u^T$ , for any  $u \in \mathbb{C}^n$ , obtaining a new solution of (3), so  $\mathbb{S}$  does not have a unique solution. A similar argument can be used if  $B_k$  is singular.

Moreover,  $\widetilde{\mathbb{S}}$  is irreducible. Otherwise, it could be split in two systems with different unknowns, and just one of them would contain  $X_{\beta_k}$ ; adding the  $k$ th equation to this last system would give a partition of the original system  $\mathbb{S}$  in two systems with different unknowns.  $\square$

The proof of Theorem 13 shows that, if an irreducible nonsingular system  $\mathbb{S}$  having  $r > 1$  unknowns contains an unknown appearing just once in  $\mathbb{S}$ , then we can remove this unknown, together with its corresponding equation, to get a new irreducible system with  $r - 1$  equations

and  $r - 1$  unknowns. Notice that the new system may have unknowns appearing just once, that can be removed if  $r > 2$ , using Theorem 13 again.

This elimination procedure can be repeated as long as the number of equations is greater than one and there is an unknown appearing just once. After a finite number of reductions (using Theorem 13 repeatedly), we arrive at an irreducible system  $\tilde{\mathbb{S}}$ , which has the same number  $\tilde{r}$  of equations and unknowns and either  $\tilde{r} = 1$  or no unknown appears in just one equation. In both cases, all unknowns in  $\tilde{\mathbb{S}}$  appear just twice. Moreover,  $\tilde{\mathbb{S}}$  is nonsingular. Therefore, we can focus, from now on, on irreducible systems with the same number of equations and unknowns, and where each unknown appears exactly twice.

### 5.3 Reduction to a periodic system with at most one $\star$

In Section 5.2 we have proved that, without loss of generality, and regarding nonsingularity, we can consider irreducible systems of  $r$  Sylvester-like equations with  $r$  matrix unknowns, any of which appearing just twice. Now, we want to show that from any system of the latter form, we can get an equivalent periodic system of the form (4).

We first note that, by renaming the unknowns if necessary, under these assumptions the system (3) can be written as

$$\begin{cases} A_k X_k^{s_k} B_k - C_k X_{k+1}^{t_k} D_k & = E_k, \quad k = 1, \dots, r-1, \\ A_r X_r^{s_r} B_r - C_r X_1^{t_r} D_r & = E_r, \end{cases} \quad (15)$$

where  $s_k, t_k \in \{1, \star\}$ . A way to show this is as follows. Let us start with  $X_1$  and choose one of the two equations containing this unknown (there are at least two as long as the system contains at least two equations). Let this equation, with appropriate relabeling of the coefficients if needed, be  $A_1 X_1^{s_1} B_1 - C_1 X_{\alpha_1}^{t_1} D_1 = E_1$ . Now we look for the other equation containing  $X_{\alpha_1}$ . With a relabeling of the coefficients if needed, this equation is  $A_2 X_{\alpha_1}^{s_2} B_2 - C_2 X_{\alpha_2}^{t_2} D_2 = E_2$ , and we proceed in this way with  $X_{\alpha_2}$  and so on with the remaining unknowns. Note that, during this process, it cannot happen that  $\alpha_i = \alpha_j$  for  $i \neq j$ , since otherwise  $X_{\alpha_i}$  would appear more than twice in the system. Therefore, at some point we end up with  $\alpha_t = 1$ . If there were some  $1 \leq j \leq r$  such that  $j \neq \alpha_i$ , for all  $i = 1, \dots, t$ , then the system would be reducible. Hence, it must be  $t = r$  and, by relabeling the unknowns as  $\alpha_k = k + 1$ , for  $k = 1, \dots, r - 1$ , and  $\alpha_r = 1$ , we get the system in the form (15).

We now show that each periodic irreducible system of the form (15) can be reduced to the simpler form (4), with at most one  $\star$ . This can be obtained by applying a sequence of  $\star$  operations and renaming of variables, without further linear algebraic manipulations. This is stated in the following result.

**Lemma 14.** *Given the system of generalized  $\star$ -Sylvester equations (15), there exists a system of the type*

$$\begin{cases} \tilde{A}_k Y_k \tilde{B}_k - \tilde{C}_k Y_{k+1} \tilde{D}_k & = \tilde{E}_k, \quad k = 1, \dots, r-1, \\ \tilde{A}_r Y_r \tilde{B}_r - \tilde{C}_r Y_1 \tilde{D}_r & = \tilde{E}_r, \end{cases} \quad (16)$$

with  $s \in \{1, \star\}$ , and  $u_k \in \{1, \star\}$ , for  $k = 1, \dots, r$ , such that  $Y_1, \dots, Y_r$  is a solution of (16) if and only if  $X_1, \dots, X_r$ , with  $X_k = Y_k^{u_k}$ , is a solution of (15).

Moreover,  $s = 1$  if the number of  $\star$  symbols appearing among  $s_i, t_i$  in the original system (15) is even, and  $s = \star$  if it is odd.

*Proof.* The proof of this result is constructive, i.e., it is presented in an algorithmic way that produces the system (16) from (15) by a sequence of transpositions and substitutions of the type  $Y_k = X_k^{u_k}$ , from which the statement follows.

The procedure has  $r$  steps. At the first step we consider the first equation. If  $s_1 = \star$  then we apply the  $\star$  operator to both sides of the equation, obtaining a new equivalent equation with no star on the first unknown:

$$A_1 X_1^\star B_1 - C_1 X_2^{t_1} D_1 = E_2 \iff B_1^\star X_1 A_1^\star - D_1^\star (X_2^{t_1})^\star C_1^\star = E_1^\star.$$

We set  $Y_1 = X_1$  and  $(\tilde{A}_1, \tilde{B}_1, \tilde{C}_1, \tilde{D}_1, \tilde{E}_1) = (B_1^\star, A_1^\star, D_1^\star, C_1^\star, E_1^\star)$ . If  $s_1 = 1$ , then we set  $Y_1 = X_1$  as well and  $(\tilde{A}_1, \tilde{B}_1, \tilde{C}_1, \tilde{D}_1, \tilde{E}_1) = (A_1, B_1, C_1, D_1, E_1)$ . In both cases,  $u_1 = 1$  and the first equation has been replaced by  $\tilde{A}_1 Y_1 \tilde{B}_1 - \tilde{C}_1 (X_2^{t_1})^{s_1} \tilde{D}_1 = \tilde{E}_1$ . Notice that, for  $r = 1$ , we get an equivalent periodic system of the type (16) and then we are done.

If  $r > 1$ , then we continue the first step of the procedure and check the second unknown of the first equation, namely  $(X_2^{t_1})^{s_1}$ , that can be  $X_2$  or  $X_2^\star$ . If the second unknown is  $X_2$ , then we set  $Y_2 = X_2$  and  $u_2 = 1$ , otherwise we set  $Y_2 = X_2^\star$  and  $u_2 = \star$ . In both cases we get an equation of the type  $\tilde{A}_1 Y_1 \tilde{B}_1 - \tilde{C}_1 Y_2 \tilde{D}_1 = \tilde{E}_1$ , with no  $\star$  in the unknowns. Replacing  $X_2$  by  $Y_2^{u_2}$  also in the second equation we get a system equivalent to (15) but with no  $\star$  in the first equation.

The procedure can be repeated for the remaining equations. The second step works on the second equation, that now is of the form  $A_2 (Y_2^{u_2})^{s_2} B_2 - C_2 X_3^{t_2} D_2 = E_2$ . If  $(Y_2^{u_2})^{s_2} = X_2$ , then we can take  $(\tilde{A}_2, \tilde{B}_2, \tilde{C}_2, \tilde{D}_2, \tilde{E}_2) = (A_2, B_2, C_2, D_2, E_2)$ ; otherwise,  $(Y_2^{u_2})^{s_2} = X_2^\star$ , so we apply the operator  $\star$  to the second equation, obtaining an equivalent one, and hence set  $(\tilde{A}_2, \tilde{B}_2, \tilde{C}_2, \tilde{D}_2, \tilde{E}_2) = (B_2^\star, A_2^\star, D_2^\star, C_2^\star, E_2^\star)$ . Then we check if the other unknown appearing in the resulting equation is  $X_3$  or  $X_3^\star$ , and proceed analogously. After  $r - 1$  steps we arrive at the last equation, which is of the form  $A_r X_r^{s_r} B_r - C_r X_1^{t_r} D_r = E_r$ , with  $X_1 = Y_1$  and either  $X_r = Y_r$  or  $X_r = Y_r^\star$ . Therefore, there are four possible cases

$$A_r Y_r B_r - C_r Y_1 D_r = E_r, \quad (17)$$

$$A_r Y_r B_r - C_r Y_1^\star D_r = E_r, \quad (18)$$

$$A_r Y_r^\star B_r - C_r Y_1 D_r = E_r, \quad (19)$$

$$A_r Y_r^\star B_r - C_r Y_1^\star D_r = E_r. \quad (20)$$

Cases (17) and (18) are already in the form required in (16). For case (19) we apply the  $\star$  operator to this equation and arrive at

$$\tilde{A}_r Y_r \tilde{B}_r - \tilde{C}_r Y_1^\star \tilde{D}_r = \tilde{E}_r,$$

with  $(\tilde{A}_r, \tilde{B}_r, \tilde{C}_r, \tilde{D}_r, \tilde{E}_r) = (B_r^\star, A_r^\star, D_r^\star, C_r^\star, E_r^\star)$ , and in case (20) we apply again the  $\star$  operator to this equation and we get

$$\tilde{A}_r Y_r \tilde{B}_r - \tilde{C}_r Y_1 \tilde{D}_r = \tilde{E}_r,$$

with  $(\tilde{A}_r, \tilde{B}_r, \tilde{C}_r, \tilde{D}_r, \tilde{E}_r) = (B_r^\star, A_r^\star, D_r^\star, C_r^\star, E_r^\star)$ , as above. Therefore, in all cases we arrive at a system (16).

Each of the transformations performed by the algorithm preserves the parity of the number of  $\star$  symbols appearing within the equations, since each change of variables may swap the exponent, from  $\star$  to 1 or vice versa, in the two appearances of each unknown. Therefore, the second part of the statement follows.  $\square$

The above results show that we can reduce the problem on the nonsingularity of (3) either to the problem of the nonsingularity of a periodic system of  $r$  generalized Sylvester equations or to the problem of the nonsingularity of a periodic system of  $r - 1$  generalized Sylvester and one generalized  $\star$ -Sylvester equation.

---

**Algorithm 1** Transformation of a periodic system into a system with just one  $\star$ . Vectors  $s$  and  $t$  contain the transpositions in the original system. The procedure returns the new coefficients, the vector  $u$  so that  $Y_k = X_k^{u_k}$ , and the symbol  $t_r$  on  $X_{r+1} = X_1$  in the last equation (which is the only entry in both  $s$  and  $t$  that could be a  $\star$  after the procedure).

---

```

1: procedure GENERATESYSTEM( $A_k, B_k, C_k, D_k, E_k, s, t$ )
2:    $u_1 \leftarrow 1$  ▷  $u_1$  is always 1, since  $Y_1 = X_1$ 
3:   for  $k = 1, \dots, r$  do
4:     if  $s_k = 1$  then
5:        $(\tilde{A}_k, \tilde{B}_k, \tilde{C}_k, \tilde{D}_k, \tilde{E}_k) \leftarrow (A_k, B_k, C_k, D_k, E_k)$ 
6:     else
7:        $(\tilde{A}_k, \tilde{B}_k, \tilde{C}_k, \tilde{D}_k, \tilde{E}_k) \leftarrow (B_k^*, A_k^*, D_k^*, C_k^*, E_k^*)$ 
8:       Swap  $t_k$  ▷ Swap the value of  $t_k$  between 1 and  $\star$ 
9:     end if
10:    if  $k < r$  then
11:       $u_{k+1} \leftarrow t_k$ 
12:      if  $t_k = \star$  then
13:        Swap  $s_{k+1}$  ▷ Swap the value of  $s_{k+1}$  between 1 and  $\star$ 
14:      end if
15:    end if
16:  end for
17:  return  $\tilde{A}_k, \tilde{B}_k, \tilde{C}_k, \tilde{D}_k, \tilde{E}_k, u, t_r$ 
18: end procedure

```

---

## 6 Reduction to a block triangular linear system

In Section 5 we have seen how a nonsingular system of general type (3) can be reduced to one or more independent periodic systems of the type (4), where all equations are generalized Sylvester equations except the last one, that is either a generalized Sylvester or a generalized  $\star$ -Sylvester equation.

Here we focus on a periodic system of type (4). First, we show in Section 6.1 that it can be transformed into an equivalent periodic system with triangular coefficients. Then, in Section 6.2 we show that, in the cases  $s = 1$  and  $s = \mathbf{T}$ , the latter system is a linear system whose coefficient matrix is block triangular with diagonal blocks of order  $r$  or  $2r$ . Finally, in Section 6.3 we show that the case  $s = \mathbf{H}$  can be reduced to the case  $s = 1$ .

The reduction to a special linear system allows one to deduce useful conditions for the nonsingularity of a system of generalized Sylvester equations and, moreover, to design an efficient numerical algorithm for its solution.

### 6.1 Reduction to a system with triangular coefficients

We can multiply by suitable unitary matrices and perform a change of variables on the system (4) which simultaneously make the matrices  $A_k, B_k, C_k, D_k$  upper or lower (quasi-)triangular.

**Lemma 15.** *There exists a change of variables of the form  $\hat{X}_k = Z_k^H X_k \hat{Z}_k$ , with  $Z_k, \hat{Z}_k \in \mathbb{C}^{n \times n}$  unitary, for  $k = 1, 2, \dots, r$ , which simultaneously makes the coefficients  $A_k, C_k$  of (4) upper triangular, and the coefficients  $B_k, D_k$  lower triangular, after pre-multiplying and post-multiplying the  $k$ th equation by appropriate unitary matrices  $Q_k$  and  $\hat{Q}_k$ , respectively.*

*Proof.* We distinguish the cases  $s = 1$  and  $s \in \{\mathbf{T}, \mathbf{H}\}$ . For both cases, we provide an appropriate

change of variables to take the system in upper/lower triangular form, based on the periodic Schur form of certain formal matrix products (see Section 3).

**Case**  $s = 1$  is already treated in [9]; we report it here for completeness. Let

$$Q_k^H A_k Z_k = \widehat{A}_k, \quad Q_k^H C_k Z_{k+1} = \widehat{C}_k, \quad Z_{r+1} = Z_1, \quad k = 1, 2, \dots, r,$$

with  $\widehat{A}_k, \widehat{C}_k$  upper triangular, be a periodic Schur form of  $C_r^{-1} A_r C_{r-1}^{-1} A_{r-1} \cdots C_1^{-1} A_1$ , and

$$\widehat{Q}_k^H B_k^H \widehat{Z}_k = \widehat{B}_k^H, \quad \widehat{Q}_k^H D_k^H \widehat{Z}_{k+1} = \widehat{D}_k^H, \quad Z_{r+1} = Z_r, \quad k = 1, 2, \dots, r,$$

with  $\widehat{B}_k^H, \widehat{D}_k^H$  upper triangular, be a periodic Schur form of  $D_r^{-H} B_r^H D_{r-1}^{-H} B_{r-1}^H \cdots D_1^{-H} B_1^H$ . Setting  $\widehat{X}_k = Z_k^H X_k \widehat{Z}_k$  and multiplying the equations in (4) by  $Q_k^H$  from the left, and by  $\widehat{Q}_k$  from the right yields a transformed system of equations with unknowns  $\widehat{X}_k$  and upper/lower triangular coefficients, as claimed.

**Case**  $s \in \{\mathbb{H}, \mathbb{T}\}$  can be handled by considering the periodic Schur form

$$\begin{aligned} Q_k^H A_k Z_k &= \widehat{A}_k, & Q_k^H C_k Z_{k+1} &= \widehat{C}_k, & Z_{2r+1} &= Z_1, \\ Q_{r+k}^H B_k^s Z_{r+k} &= \widehat{B}_k^s, & Q_{r+k}^H D_k^s Z_{r+k+1} &= \widehat{D}_k^s, & k &= 1, 2, \dots, r, \end{aligned}$$

of  $D_r^{-s} B_r^s D_{r-1}^{-s} B_{r-1}^s \cdots D_1^{-s} B_1^s C_r^{-1} A_r C_{r-1}^{-1} A_{r-1} \cdots C_1^{-1} A_1$ .

Performing the change of variables  $\widehat{X}_k = Z_k^H X_k (Z_{r+k}^s)^H$  and multiplying the equations in (4) by  $Q_k$  on the left and by  $(Q_{r+k}^s)^H$  on the right yields a system with upper/lower triangular coefficients in the unknowns  $\widehat{X}_k$ . Note that, for any matrix  $M$ ,  $(M^s)^H$  is equal to  $M$  if  $s = \mathbb{H}$  and  $\overline{M}$  (the complex conjugate) if  $s = \mathbb{T}$ .

□

## 6.2 Reduction to a block upper triangular linear system for $s = 1, \mathbb{T}$

A system like (4) can be seen as a system of  $n^2 r$  equations in  $n^2 r$  unknowns in terms of the entries of the unknown matrices. This is a linear system for  $s = 1$  or  $s = \mathbb{T}$ , while in the case  $s = \mathbb{H}$  it is not linear over  $\mathbb{C}$  due to the conjugation. Nevertheless, it can be either transformed into a linear system over  $\mathbb{R}$ , by splitting the real and imaginary parts of both the coefficients and the unknowns (see Section 8.2), or into a linear system over  $\mathbb{C}$  by doubling the size (see Section 6.3).

A standard approach to get explicitly the matrix coefficient of the (linear) system associated with a system of Sylvester-like equations is to exploit the relation  $\text{vec}(AXB) = (B^T \otimes A) \text{vec} X$  [24, Lemma 4.3.1] where the  $\text{vec}(\cdot)$  operator maps a matrix into the vector obtained by stacking its columns one on top of the other, and  $A \otimes B$  is the Kronecker product of  $A$  and  $B$ , namely the block matrix with blocks of the type  $[a_{ij}B]$  (see [24, Ch. 4]).

Relying on the reduction scheme that we have presented in Section 6.1, we may assume that the coefficients  $A_k, C_k$ , and  $B_k, D_k$ , in (4) are upper and lower triangular matrices, respectively. In this case the matrix of the linear system obtained after applying the  $\text{vec}(\cdot)$  operator has a nice structure; indeed, performing appropriate row and column permutations to the matrix (in other words, choosing an appropriate ordering of the unknowns), in Section 6.2.1, we get a block upper triangular coefficient matrix, with diagonal blocks of dimensions  $r$  or  $2r$ .

In the case where  $s = 1$ , a characterization for nonsingularity was obtained in [9] (see Theorem 9). The approach followed in that reference is similar to the one we follow here.

We first deal with the cases  $s \in \{1, \top\}$ , which are both linear, and for which we can directly give conditions based on the matrix representing the linear system in the entries of the unknowns. This is the aim of Section 6.2.1. The case  $s = \mathbf{H}$  can be reduced to the case  $s = 1$  by using specific developments which are contained in Section 6.3.

### 6.2.1 Making the matrix coefficient block triangular

We assume that  $A_k, C_k$  are upper triangular and  $B_k, D_k$  are lower triangular, for  $k = 1, \dots, r$ .

Using the relation  $\text{vec}(AXB) = (B^\top \otimes A) \text{vec} X$  we can rewrite the system (4), for the case  $s = 1$ , with  $r > 1$ , as the linear system

$$\begin{bmatrix} B_1^\top \otimes A_1 & -D_1^\top \otimes C_1 & & & \\ & \ddots & \ddots & & \\ & & & B_{r-1}^\top \otimes A_{k-1} & -D_{r-1}^\top \otimes C_{r-1} \\ & & & & B_r^\top \otimes A_r \\ -D_r^\top \otimes C_r & & & & \end{bmatrix} \mathcal{X} = \mathcal{E}, \quad (21)$$

where the empty block entries should be understood as zero blocks, and

$$\mathcal{X} := \begin{bmatrix} \text{vec} X_1 \\ \vdots \\ \text{vec} X_r \end{bmatrix}, \quad \mathcal{E} := \begin{bmatrix} \text{vec} E_1 \\ \vdots \\ \text{vec} E_r \end{bmatrix}.$$

In the case  $s = \top$ , with  $r > 1$ , we have, instead

$$\begin{bmatrix} B_1^\top \otimes A_1 & -D_1^\top \otimes C_1 & & & \\ & \ddots & \ddots & & \\ & & & B_{r-1}^\top \otimes A_{k-1} & -D_{r-1}^\top \otimes C_{r-1} \\ & & & & B_r^\top \otimes A_r \\ -(D_r^\top \otimes C_r)P_{n,n} & & & & \end{bmatrix} \mathcal{X} = \mathcal{E}, \quad (22)$$

where  $P_{a,b}$  denotes the *commutation matrix*, i.e., the permutation matrix such that  $P_{a,b} \text{vec} X = \text{vec}(X^\top)$  for each  $X \in \mathbb{R}^{a \times b}$  [24, Th. 4.3.8].

In the case  $r = 1$ , the system is  $(B_1^\top \otimes A_1 - D_1^\top \otimes C_1)\mathcal{X} = \mathcal{E}$  for  $s = 1$  and  $(B_1^\top \otimes A_1 - (D_1^\top \otimes C_1)P_{n,n})\mathcal{X} = \mathcal{E}$  for  $s = \top$ .

In the following, we index the components of  $\mathcal{X}$  by means of the triple  $(i, j, k)$ , that denotes the  $(i, j)$  entry of  $X_k$ . This is just a shorthand for the component  $(k-1)n^2 + (j-1)n + i$  of  $\mathcal{X}$ . Notice that each coordinate of any of the systems (21) and (22) can be obtained by multiplying one of the  $r$  equations of (4) by  $e_i^\top$  on the left and by  $e_j$  on the right, for appropriate  $1 \leq i, j \leq n$ .

We are interested in performing a permutation on systems (21) and (22) that takes them to block upper triangular form (independently on the presence of the permutation matrix  $P_{n,n}$ ). The next Lemma shows that this is always possible.

**Lemma 16.** *Let  $A_k, C_k$  be  $n \times n$  upper triangular matrices and  $B_k, D_k$  be  $n \times n$  lower triangular matrices, for  $k = 1, \dots, r$ . Let  $\mathbb{S}$  be the system of  $n^2 r$  equations*

$$\begin{cases} e_i^\top (A_k X_k B_k - C_k X_{k+1} D_k) e_j &= (E_k)_{ij}, & i, j = 1, \dots, n, & k = 1, \dots, r-1, \\ e_i^\top (A_r X_r B_r - C_r X_1^s D_r) e_j &= (E_r)_{ij}, & i, j = 1, \dots, n, \end{cases} \quad (23)$$

in the  $n^2 r$  unknowns  $x_{ijk}$ , for  $i, j = 1, \dots, n$  and  $k = 1, \dots, r$ , where  $x_{ijk}$  is the  $(i, j)$  entry of  $X_k$ . With a suitable ordering of the equations and unknowns, the coefficient matrix  $M \in \mathbb{C}^{n^2 r \times n^2 r}$  of the system is block upper triangular, with diagonal blocks of size either  $r \times r$  or  $2r \times 2r$

*Proof.* We define an ordering of the triples  $(i, j, k)$  as follows. Define the following ordered sublists

$$\begin{aligned}\mathcal{L}_{ii} &= (i, i, 1), (i, i, 2), \dots, (i, i, r), & 1 \leq i \leq n, \\ \mathcal{L}_{ij} &= (i, j, 1), (i, j, 2), \dots, (i, j, r), (j, i, 1), (j, i, 2), \dots, (j, i, r), & 1 \leq j < i \leq n;\end{aligned}$$

then, we concatenate these sublists in lexicographic order of their index,

$$\mathcal{L}_{11}, \mathcal{L}_{21}, \mathcal{L}_{22}, \mathcal{L}_{31}, \mathcal{L}_{32}, \mathcal{L}_{33}, \mathcal{L}_{41}, \mathcal{L}_{42}, \mathcal{L}_{43}, \mathcal{L}_{44}, \dots, \mathcal{L}_{n1}, \mathcal{L}_{n2}, \dots, \mathcal{L}_{nn}. \quad (24)$$

In the matrix  $M$ , we sort the equations (23) (corresponding to rows) and the unknowns  $x_{ijk}$  (corresponding to columns) according to this order (24) of the triples  $(i, j, k)$ . Grouping together the triples that belong to the same sublist  $\mathcal{L}_{ij}$ , we obtain a block partition of  $M$  with  $\frac{n(n+1)}{2}$  block rows and columns, each of size  $r$  or  $2r$ , depending on whether  $i \neq j$  or  $i = j$ .

In order to simplify the notation, we set  $x_{i,j,r+1} = x_{ij1}$  if  $s = 1$  and  $x_{i,j,r+1} = x_{ji1}$  if  $s = \star$ . With this choice,  $x_{ijk}$  and  $x_{i,j,k'}$  belong to the same sublist ( $\mathcal{L}_{ij}$  or  $\mathcal{L}_{ji}$ ) for any  $k, k'$ , and whenever  $i \leq \ell$  and  $j \leq t$  the unknown  $x_{ijk}$  belongs to a sublist that comes before  $x_{\ell tk}$ .

Since  $A_k$  is upper triangular and  $B_k$  is lower triangular, for a given  $(i, j, k)$  we have

$$(A_k X_k B_k)_{ij} = \sum_{\ell=1}^n (A_k)_{i\ell} \sum_{t=1}^n (X_k)_{\ell t} (B_k)_{tj} = \sum_{\ell=i}^n \sum_{t=j}^n (A_k)_{i\ell} (X_k)_{\ell t} (B_k)_{tj},$$

and similarly for  $(C_k X_{k+1} D_k)_{ij}$ . Thus the  $(i, j, k)$  equation of the system is

$$\sum_{\substack{i \leq \ell \\ j \leq t}} ((A_k)_{i\ell} x_{\ell tk} (B_k)_{tj} - (C_k)_{i\ell} x_{\ell, t, k+1} (D_k)_{tj}) = (E_k)_{ij}.$$

Hence an equation with index in  $\mathcal{L}_{ij}$  contains only unknowns belonging to the sublist  $\mathcal{L}_{ij}$  and to sublists that follow it in the order of (24). This proves that  $M$  is block upper triangular.  $\square$

### 6.2.2 Characterizing the diagonal blocks

Both from the computational and from the theoretical point of view we are interested in characterizing the structure of the diagonal blocks of the coefficient matrix  $M$  associated with the linear system obtained by applying the permutation of Lemma 16.

Theoretically, this is interesting because the system (4) is nonsingular if and only if the determinants of all diagonal blocks of  $M$  are nonzero. This will allow us to prove Theorems 5 and 6.

Computationally, this is relevant because these are the matrices that allow one to carry out the block back substitution process to compute the solution of (4), when it is unique.

As already pointed out in Section 6.2.1 the diagonal blocks can be obtained by choosing a pair  $(i, j)$  and selecting the equations given by

$$\begin{cases} e_i^\top (A_k X_k B_k - C_k X_{k+1} D_k) e_j &= (E_k)_{ij}, \quad k = 1, \dots, r-1, \\ e_i^\top (A_r X_r B_r - C_r X_1^s D_r) e_j &= (E_r)_{ij}, \end{cases}$$

and the ones obtained by the pair  $(j, i)$ , and removing all the variables with indices different from  $(i, j)$  and  $(j, i)$ . As mentioned in the proof of Lemma 16, these other variables have indices  $(i', j', k')$  belonging to a subset  $\mathcal{L}_{i', j'}$  that follows  $\mathcal{L}_{ij}$  in the given order, and hence their value has already been computed in the back substitution process. When  $i = j$  this gives us an  $r \times r$  linear system, otherwise we obtain a  $2r \times 2r$  linear system. We denote them with  $\mathbb{S}_{ij}$ , for  $i \geq j$ .

Notice that this procedure can be carried out both in the case  $s \in \{1, \mathsf{T}\}$  and in the  $s = \mathsf{H}$  case, even if in the latter these systems are nonlinear.



**Lemma 17.** *Let  $M$  be the following matrix:*

$$M = \begin{bmatrix} \alpha_1 & \beta_1 & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & \beta_{p-1} \\ \beta_p & & & & \alpha_p \end{bmatrix}.$$

Then,  $\det M = \prod_{k=1}^p \alpha_k - (-1)^p \prod_{k=1}^p \beta_k$ .

*Proof.* Use Laplace's determinant expansion on the first column.  $\square$

In the cases  $s \in \{1, \top\}$ ,  $\mathbb{S}_{ii}$  is an  $r \times r$  linear system in the variables  $(X_1)_{ii}, \dots, (X_r)_{ii}$  with coefficient matrix:

$$M_{ii} := \begin{bmatrix} (A_1)_{ii}(B_1)_{ii} & -(C_1)_{ii}(D_1)_{ii} & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & (A_{r-1})_{ii}(B_{r-1})_{ii} & -(C_{r-1})_{ii}(D_{r-1})_{ii} \\ -(C_r)_{ii}(D_r)_{ii} & & & & (A_r)_{ii}(B_r)_{ii} \end{bmatrix}, \quad (25)$$

for  $r > 1$  and  $M_{ii} = (A_1)_{ii}(B_1)_{ii} - (C_1)_{ii}(D_1)_{ii}$  for  $r = 1$ .

According to Lemma 17 we have:

$$\det M_{ii} = \prod_{k=1}^r (A_k)_{ii}(B_k)_{ii} - \prod_{k=1}^r (C_k)_{ii}(D_k)_{ii}. \quad (26)$$

A similar relation holds also when  $i > j$  in the  $s = 1$  case, since  $\mathbb{S}_{ij}$  can be decoupled into two  $r \times r$  systems. More precisely, in the case  $s = 1$ , the coefficient matrix of  $\mathbb{S}_{ij}$  is block diagonal with two diagonal blocks, the top left block is

$$M_{ij} := \begin{bmatrix} (A_1)_{ii}(B_1)_{jj} & -(C_1)_{ii}(D_1)_{jj} & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & (A_{r-1})_{ii}(B_{r-1})_{jj} & -(C_{r-1})_{ii}(D_{r-1})_{jj} \\ -(C_r)_{ii}(D_r)_{jj} & & & & (A_r)_{ii}(B_r)_{jj} \end{bmatrix}, \quad (27)$$

for  $r > 1$  and  $M_{ij} = (A_1)_{ii}(B_1)_{jj} - (C_1)_{ii}(D_1)_{jj}$  for  $r = 1$ , while the lower bottom block,  $M_{ji}$ , is obtained exchanging the roles of  $i$  and  $j$ . From Lemma 17 we get:

$$\det M_{ij} = \prod_{k=1}^r (A_k)_{ii}(B_k)_{jj} - \prod_{k=1}^r (C_k)_{ii}(D_k)_{jj}. \quad (28)$$

In the case  $s = \top$ , instead, the systems  $\mathbb{S}_{ij}$  form a  $2r \times 2r$  linear system in the variables  $(X_k)_{ij}, (X_k)_{ji}$ , for  $k = 1, \dots, r$ , with coefficient matrix

$$M_{ij} := \begin{bmatrix} \mathcal{B}_{ij} & -(C_r)_{ii}(D_r)_{jj}e_re_1^\top \\ -(C_1)_{jj}(D_1)_{ii}e_re_1^\top & \mathcal{B}_{ji} \end{bmatrix}, \quad (29)$$

where

$$\mathcal{B}_{ij} = \begin{bmatrix} (A_1)_{ii}(B_1)_{jj} & -(C_1)_{ii}(D_1)_{jj} & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & -(C_{r-1})_{ii}(D_{r-1})_{jj} \\ & & & & (A_r)_{ii}(B_r)_{jj} \end{bmatrix}.$$

Thanks, again, to Lemma 17, this matrix has determinant equal to

$$\det M_{ij} = \prod_{k=1}^r (A_k)_{ii}(B_k)_{ii}(A_k)_{jj}(B_k)_{jj} - \prod_{k=1}^r (C_k)_{ii}(D_k)_{ii}(C_k)_{jj}(D_k)_{jj}. \quad (30)$$

### 6.3 Linearizing the case $s = \mathbb{H}$

We have already mentioned that, when  $s = \mathbb{H}$ , the system (4) is not linear over the complex field, since it involves not only the entries of the matrix  $X_1$  but also their conjugates. A method to transform it into a linear system over  $\mathbb{C}$  is as follows: in addition to the equations of the system, we consider the equations obtained by taking their conjugate transpose, namely

$$\begin{aligned} B_k^{\mathbb{H}} X_k^{\mathbb{H}} A_k^{\mathbb{H}} - D_k^{\mathbb{H}} X_{k+1}^{\mathbb{H}} C_k^{\mathbb{H}} &= E_k^{\mathbb{H}}, \quad k = 1, \dots, r-1, \\ B_r X_r^{\mathbb{H}} A_r - D_r^{\mathbb{H}} X_1 C_r^{\mathbb{H}} &= E_r^{\mathbb{H}}. \end{aligned}$$

If we consider  $X_k$  and  $X_k^{\mathbb{H}}$  as two separate variables, then this is a system of  $2r$  generalized Sylvester equations in  $2r$  matrix unknowns. We prove more formally that this process produces an equivalent system.

**Lemma 18.** *The system (4) is nonsingular if and only if the system*

$$\begin{cases} A_k X_k B_k - C_k X_{k+1} D_k &= E_k, \quad k = 1, \dots, r-1, \\ A_r X_r B_r - C_r X_{r+1} D_r &= E_r, \\ B_k^{\mathbb{H}} X_{r+k} A_k^{\mathbb{H}} - D_k^{\mathbb{H}} X_{r+k+1} C_k^{\mathbb{H}} &= E_k^{\mathbb{H}}, \quad k = 1, \dots, r-1, \\ B_r^{\mathbb{H}} X_{2r} A_r^{\mathbb{H}} - D_r^{\mathbb{H}} X_1 C_r^{\mathbb{H}} &= E_r^{\mathbb{H}} \end{cases} \quad (31)$$

is nonsingular.

*Proof.* We may consider only the case in which  $E_k = 0$ : checking nonsingularity corresponds to checking that there are no solutions to this homogenous system apart from the trivial one  $X_k = 0$ , for  $k = 1, \dots, r$ .

Let us first assume that (4) has a nonzero solution  $(X_1, \dots, X_r)$ . Then  $(X_1, \dots, X_r, X_1^{\mathbb{H}}, \dots, X_r^{\mathbb{H}})$  is a nonzero solution of (31).

Conversely, if  $(X_1, \dots, X_r, X_{r+1}, \dots, X_{2r})$  is a nonzero solution of (31), then  $(X_1 + X_{r+1}^{\mathbb{H}}, \dots, X_r + X_{2r}^{\mathbb{H}})$  is a solution of (4). If  $(X_1 + X_{r+1}^{\mathbb{H}}, \dots, X_r + X_{2r}^{\mathbb{H}}) = 0$ , then  $X_{r+i} = -X_i^{\mathbb{H}}$ , for  $i = 1, \dots, r$ , and then  $i(X_1, \dots, X_r)$  is a nonzero solution of (4).  $\square$

*Remark 19.* The proof of Lemma 18 does not work if one replaces  $\mathbb{H}$  with  $\mathbb{T}$  everywhere: it breaks in the final part, because  $i(X_1, \dots, X_r)$  is not necessarily a solution of (4) with  $\star = \mathbb{T}$ . Indeed, Lemma 18 is false with  $\mathbb{T}$  instead of  $\mathbb{H}$ . Let us consider, for instance, the case  $n = r = 1$  and the equation  $x_1 + x_1^{\mathbb{T}} = 2x_1 = 0$ . This equation has only the trivial solution, but the linearized system

$$\begin{cases} z_1 + z_2 = 0 \\ z_1 + z_2 = 0 \end{cases}$$

has infinitely many solutions.

Another relevant difference between the  $\star = \mathbb{T}$  and the  $\star = \mathbb{H}$  cases is the following. System (4) is nonsingular if and only if the system obtained after replacing the minus sign in the last equation by a plus sign

$$\begin{cases} A_k X_k B_k - C_k X_{k+1} D_k &= E_k, \quad k = 1, \dots, r-1, \\ A_r X_r B_r + C_r X_1^{\mathbb{H}} D_r &= E_r \end{cases} \quad (32)$$

is nonsingular. To see this, reduce again to the case  $E_k = 0$  for all  $k = 1, \dots, r$  and note that if  $(X_1, \dots, X_r)$  is a nonzero solution of (4) then  $\mathfrak{i}(X_1, \dots, X_r)$  is a nonzero solution of (32), and viceversa. This property no longer holds true with  $s = \mathbb{T}$ .

## 7 Proofs of the main results

Here we prove Theorems 5–6, with the aid of all previous developments. We start with Theorem 5.

*Proof of Theorem 5.* We can consider only the case in which  $E_i = 0$ ,  $i = 1, 2, \dots, r$ . Using the periodic Schur form of the formal products (10) we may consider the equivalent system (see the proof of Lemma 15)

$$\begin{cases} \widehat{A}_k X_k \widehat{B}_k - \widehat{C}_k X_{k+1} \widehat{D}_k &= 0, \quad k = 1, \dots, r-1, \\ \widehat{A}_r X_r \widehat{B}_r - \widehat{C}_r X_1 \widehat{D}_r &= 0, \end{cases}$$

where, for each  $k$ , the matrices  $\widehat{A}_k$  and  $\widehat{C}_k$  are upper triangular and  $\widehat{B}_k$  and  $\widehat{D}_k$  are lower triangular. If the formal products (10) are regular, then their eigenvalues are the ratios  $\lambda_i := \prod_{k=1}^r \frac{(\widehat{A}_k)_{ii}}{(\widehat{C}_k)_{ii}}$ ,  $\mu_i := \prod_{k=1}^r \frac{(\widehat{D}_k)_{ii}}{(\widehat{B}_k)_{ii}}$ , respectively, for  $i = 1, \dots, n$  (they are allowed to be  $\infty$ ).

With this triangularity assumption, in Lemma 16 we have shown that the system of Sylvester equations is equivalent to a block upper triangular system whose matrix coefficient has determinant  $\delta := \prod_{i,j=1}^n \det(M_{ij})$ , where  $M_{ij}$  is defined in (25) and (27).

In summary, the system of Sylvester equations is nonsingular if and only if  $\delta \neq 0$ , which, using (26) and (28), is equivalent to requiring

$$\prod_{k=1}^r (\widehat{A}_k)_{ii} (\widehat{B}_k)_{jj} \neq \prod_{k=1}^r (\widehat{C}_k)_{ii} (\widehat{D}_k)_{jj}, \quad i, j = 1, \dots, n. \quad (33)$$

If  $\delta \neq 0$ , then it cannot happen that  $\prod_k (\widehat{A}_k)_{ii}$  and  $\prod_k (\widehat{C}_k)_{ii}$  are both zero or that  $\prod_k (\widehat{B}_k)_{ii}$  and  $\prod_k (\widehat{D}_k)_{ii}$  are both zero and thus the formal products are regular. Moreover, condition (33) implies that  $\lambda_i \neq \mu_j$  for any  $i, j = 1, \dots, n$  and thus the two products have disjoint spectra.

On the contrary, if  $\delta = 0$  then the equality holds in (33) for some  $i$  and  $j$ . One can check that this condition implies that either one of the two formal products is singular or  $\lambda_i = \mu_j$  and they cannot have disjoint spectra.  $\square$

We now give the proof of Theorem 6 separating the cases  $\star = \mathbb{T}$  and  $\star = \mathbb{H}$  since the techniques we use are different.

*Proof of Theorem 6 for  $\star = \mathbb{T}$ .* Proceeding as in the proof of Theorem 5, we use the periodic Schur form of the formal product (11) to get the equivalent system (see the proof of Lemma 15)

$$\begin{cases} \widehat{A}_k X_k \widehat{B}_k - \widehat{C}_k X_{k+1} \widehat{D}_k &= 0, \quad k = 1, \dots, r-1, \\ \widehat{A}_r X_r \widehat{B}_r - \widehat{C}_r X_1^{\mathbb{T}} \widehat{D}_r &= 0, \end{cases}$$

where, for each  $k$ , the matrices  $\widehat{A}_k$  and  $\widehat{C}_k$  are upper triangular and  $\widehat{B}_k$  and  $\widehat{D}_k$  are lower triangular. If the formal product (10) is regular, then its eigenvalues are the ratios  $\lambda_i := \prod_{k=1}^r \frac{(\widehat{A}_k)_{ii}(\widehat{B}_k)_{ii}}{(\widehat{C}_k)_{ii}(\widehat{D}_k)_{ii}}$ , for  $i = 1, \dots, n$ .

With this triangularity assumption, in Lemma 16 we have shown that the previous system is equivalent to a block upper triangular system whose coefficient matrix has determinant  $\delta := \prod_{i=1}^n \det(M_{ii}) \prod_{\substack{i,j=1 \\ i < j}}^n \det(M_{ij})$ , with  $M_{ii}$  as in (25) and  $M_{ij}$ , for  $i \neq j$ , as in (29).

In summary, the system of Sylvester-like equations is nonsingular if and only if  $\delta \neq 0$ , that, using (26) and (30), is equivalent to requiring

$$\begin{aligned} \prod_{k=1}^r (\widehat{A}_k)_{ii}(\widehat{B}_k)_{ii} &\neq \prod_{k=1}^r (\widehat{C}_k)_{ii}(\widehat{D}_k)_{ii}, & i = 1, \dots, n, \\ \prod_{k=1}^r (\widehat{A}_k)_{ii}(\widehat{B}_k)_{ii}(\widehat{A}_k)_{jj}(\widehat{B}_k)_{jj} &\neq \prod_{k=1}^r (\widehat{C}_k)_{ii}(\widehat{D}_k)_{ii}(\widehat{C}_k)_{jj}(\widehat{D}_k)_{jj}, & i \neq j. \end{aligned} \quad (34)$$

If  $\delta \neq 0$ , then it cannot happen that  $\prod_k (\widehat{A}_k)_{ii}(\widehat{B}_k)_{ii}$  and  $\prod_k (\widehat{C}_k)_{ii}(\widehat{D}_k)_{ii}$  are both zero, for some  $i$ , thus the formal product (10) is regular. Moreover, conditions (34) imply that

$$\begin{cases} \lambda_i \neq 1, & i = 1, \dots, n \\ \lambda_i \neq \lambda_j^{-1}, & i \neq j, \end{cases}$$

and this implies in turn that the spectrum  $\Lambda(\Pi) \setminus \{-1\}$  is reciprocal free and the multiplicity of  $\{-1\}$  is at most one.

On the contrary, if  $\delta = 0$  then the equality holds in (34) above for some  $i$  or below for some pair  $(i, j)$ , with  $i \neq j$ . One can check that this condition implies that one of the following cases holds: (a) the formal product is singular; (b)  $\lambda_i = 1$ , for some  $i$ , and thus  $\Lambda(\Pi) \setminus \{-1\}$  is not reciprocal free; (c)  $\lambda_i = 1/\mu_j \neq -1$ , for some  $i \neq j$ , and thus  $\Lambda(\Pi) \setminus \{-1\}$  is not reciprocal free; (d)  $\lambda_i = 1/\mu_j = -1$  and the multiplicity of  $-1$  is greater than 1.  $\square$

Using Lemma 18, the following argument allows us to obtain Theorem 6 with  $\star = \mathbb{H}$  directly as a consequence of Theorem 5.

*Proof of Theorem 6 for  $\star = \mathbb{H}$ .* Let us start from a system of the form (4) with  $s = \mathbb{H}$ . Lemma 18 shows that it is nonsingular if and only if the larger linear system (31) is nonsingular. System (31) is a system of  $2r$  generalized Sylvester equations with  $s = 1$ . Hence we can apply Theorem 5 to this system, obtaining that (31) is nonsingular if and only if the two formal products

$$\Pi_1 := \Pi = D_r^{-\mathbb{H}} B_r^{\mathbb{H}} D_{r-1}^{-\mathbb{H}} B_{r-1}^{\mathbb{H}} \cdots D_1^{-\mathbb{H}} B_1^{\mathbb{H}} C_r^{-1} A_r C_{r-1}^{-1} A_{r-1} \cdots C_1^{-1} A_1$$

and

$$\Pi_2 := C_r^{\mathbb{H}} A_r^{-\mathbb{H}} C_{r-1}^{\mathbb{H}} A_{r-1}^{-\mathbb{H}} \cdots C_1^{\mathbb{H}} A_1^{-\mathbb{H}} D_r B_r^{-1} D_{r-1} B_{r-1}^{-1} \cdots D_1 B_1^{-1}$$

are regular and have no common eigenvalues. If  $\lambda_1, \lambda_2, \dots, \lambda_n$  denote the eigenvalues of  $\Pi_1$ , then the eigenvalues of the formal product

$$\Pi_2^{-\mathbb{H}} := C_r^{-1} A_r C_{r-1}^{-1} A_{r-1} \cdots C_1^{-1} A_1 D_r^{-\mathbb{H}} B_r^{\mathbb{H}} D_{r-1}^{-\mathbb{H}} B_{r-1}^{\mathbb{H}} \cdots D_1^{-\mathbb{H}} B_1^{\mathbb{H}}$$

are again  $\lambda_1, \lambda_2, \dots, \lambda_n$ , because  $\Pi_2^{-\mathbb{H}}$  differs from  $\Pi_1$  only by a cyclic permutation of the factors. This proves that the eigenvalues of  $\Pi_2$  are  $(\bar{\lambda}_1)^{-1}, (\bar{\lambda}_2)^{-1}, \dots, (\bar{\lambda}_n)^{-1}$ , so they are distinct from those of  $\Pi_1$  if and only if  $\Lambda(\Pi_1)$  is a  $\mathbb{H}$ -reciprocal free set.  $\square$

This proof shows clearly the connection between the condition on a single formal product in Theorem 5 and the condition on two products in Theorem 6. Unfortunately, we were unable to find a simple modification of this argument that works for the case  $\star = \mathbb{T}$ , mostly due to the issue presented in Remark 19.

## 8 An $O(n^3 r)$ algorithm for computing the solution

Here we describe an efficient algorithm for the solution of a nonsingular system of  $r$  Sylvester-like equations (3) of size  $n \times n$ . We follow the big-oh notation  $O(\cdot)$ , as in [23], for both large and small quantities, and we use the number of floating point operations (flops) as a complexity measure.

The tools needed to develop the algorithm are the same used, in the previous sections, for the nonsingularity results. In the description of the algorithm we focus on the complex case and so we consider triangular coefficients. However, a solution with quasitriangular forms in case of real data can be done following a similar procedure.

We proceed through the following steps:

1. (Step 1) We perform a suitable number of substitutions, changes and elimination of variables, in order to transform the system into irreducible systems of periodic form (4), as described in Section 5.
2. (Step 2) For each (irreducible) periodic system, we compute a periodic Schur decomposition to reduce the coefficients, say  $A_k, B_k, C_k, D_k$ , to upper and lower triangular forms, as described in Section 6.1.
3. (Step 3) Since the resulting systems can be seen as essentially block triangular linear systems (as described in Section 6.2.1), we solve them by back substitution.
4. (Step 4) We compute the value of the variables that have been eliminated in Step 1 (using Theorem 13).

This section describes how to handle these steps algorithmically. Moreover, we perform an analysis of the computational costs, showing that the solution can be computed in  $O(n^3 r)$  flops, and we prove a backward stability result for the computed solution.

We discuss Step 1 in Section 8.1. Step 2 amounts to computing a periodic Schur factorization, which can be carried out in  $O(n^3 r)$  flops; we refer to [7] for details concerning it.

Step 3 is the one that requires more discussion; we devote Sections 8.2–8.4 to it. Moreover, we perform a backward error analysis for the resulting algorithm in Section 8.6. We focus on the case  $s = \star$ , since the case  $s = 1$  can be found in [9]. The cases  $\star = \mathbb{T}$  and  $\star = \mathbb{H}$  are handled in a similar way, but the former is easier to describe since the associated system is linear, without the need of separating the real and imaginary parts. We describe accurately the procedure for  $\star = \mathbb{T}$ , and briefly explain the modifications needed for  $\star = \mathbb{H}$ . The procedure for  $r = 1$  is the same as the one proposed in [13], and thus our algorithm can be seen as a generalization of the one presented in [13].

Finally, Step 4 amounts to applying formula (14) several times.

### 8.1 An algorithm for the reduction step

We describe how Step 1 can be implemented in  $O(r)$  operations. This requires concepts and tools from graph theory, that can be found in [12]. Technically, there are no floating-point operations, so one could argue that this step has cost 0 in our model, but nevertheless it is useful to have an efficient way to perform it on a real-world computer.

Consider the undirected multigraph with self loops in which the nodes are the unknowns  $X_1, \dots, X_r$ , and there is an edge  $(X_i, X_j)$  for each equation in which  $X_i$  and  $X_j$  appear. A self loop arises when an equation contains just one variable, and multiple edges arise when the same two unknowns appear in several equations.

Reducing the system into irreducible subsystems corresponds to identifying the connected components of this graph, which can be done with  $O(r)$  operations, since it has  $r$  edges. We now consider each connected component  $\mathbb{S}(\mathcal{I}_k)$  separately; if the system is irreducible, the corresponding subgraph  $(V_k, \mathcal{E}_k)$  has  $r_k$  nodes and  $r_k$  edges (see Theorem 12). Removing from  $\mathcal{E}_k$  the self loops and the repeated edges (leaving just one of them for each occurrence), we get a connected subgraph  $(V_k, \tilde{\mathcal{E}}_k)$ . If  $(V_k, \mathcal{E}_k)$  had two self loops or one self-loop and a multiple edge or two multiple edges or a multiple edge with more than two edges, then  $(V_k, \tilde{\mathcal{E}}_k)$  would be a connected graph with less than  $r_k - 1$  edges and  $r_k$  nodes and this cannot happen. Thus, there are three possible cases:

- Case 1.  $(V_k, \mathcal{E}_k)$  has no self loops and no multiple edges;
- Case 2.  $(V_k, \mathcal{E}_k)$  has one self loop and no multiple edges;
- Case 3.  $(V_k, \mathcal{E}_k)$  has no self loops and one double edge.

After removing the self loop or the double edge (if any), choose an arbitrary node of the resulting graph  $(V_k, \tilde{\mathcal{E}}_k)$  as root, and perform a graph visit using breadth-first search (BFS, [12]). Since  $(V_k, \tilde{\mathcal{E}}_k)$  is connected, this visit will find all its vertices and form a predecessor subgraph  $\mathcal{T}$  that contains  $r_k - 1$  edges of  $(V_k, \tilde{\mathcal{E}}_k)$  [12]. In any of the three cases above,  $\mathcal{T}$  is a tree obtained from  $(V_k, \tilde{\mathcal{E}}_k)$  removing one edge; let  $(i, j)$  be this missing edge.

The two nodes  $i, j$  are connected by a path in  $\mathcal{T}$  via their least common ancestor. In Case 1 this path can be determined from the predecessor subgraph structure: for instance, build the paths from  $i$  and  $j$  to the root of  $\mathcal{T}$  and remove their common final part; in Case 2, we have  $i = j$  and the path is empty; in Case 3, the path is the edge in  $\mathcal{T}$  connecting  $i$  and  $j$ . Together with the removed edge  $(i, j)$ , this path forms a cycle  $(\mathcal{C}, \mathcal{E}_{\mathcal{C}})$  in  $(V_k, \mathcal{E}_k)$ . The graph  $(V_k, \mathcal{E}_k \setminus \mathcal{E}_{\mathcal{C}})$  contains no cycles, because it is a subgraph of the predecessor subgraph. Moreover, in  $(V_k, \mathcal{E}_k \setminus \mathcal{E}_{\mathcal{C}})$  each node is connected to exactly one node of the cycle  $\mathcal{C}$  (because if it were connected to more than one, this would form a cycle in  $\mathcal{T}$ ). Hence,  $(V_k, \mathcal{E}_k \setminus \mathcal{E}_{\mathcal{C}})$  is a collection of trees, each containing exactly one node of  $\mathcal{C}$ . We perform a visit of each of these trees, starting from its unique node  $c \in \mathcal{C}$ . The variables corresponding to the nodes other than  $c$  in this tree can be eliminated one by one, starting from the leaves (in the reverse of the order in which they are discovered by the BFS), with the elimination step described in Section 5.2, which removes a degree-1 tree from the graph. This elimination procedure reduces the system of equations associated to  $V_k$  to the one associated to  $\mathcal{C}$ , which is a periodic system.

All the steps described above can easily be implemented with  $O(r)$  operations— $O(r_k)$  for each connected component—just by operations on the indices. Once we have identified which cycles are formed, the coefficients can be swapped, transposed and conjugated as needed in  $O(n^2 r)$  operations (in-place, if one wishes to minimize the space overhead).

## 8.2 Solving the triangular system

We consider the block-triangular system (23) with  $s = \text{T}$ , ordered according to (24), as described in Lemma 16. This system is block upper triangular with  $\frac{n(n+1)}{2}$  diagonal blocks of order  $r$  and  $2r$ . We refer to the linear systems corresponding to these diagonal blocks as the *small systems*  $\mathbb{S}_{ij}$ .

We provide in this section a high-level overview of the solution of this system by block back substitution, and in Sections 8.3 and 8.4 we describe how to perform it within the required computational cost.

At each of the  $\frac{n(n+1)}{2}$  steps of the back substitution process, we need to solve a square linear system of the form:

$$M_{ij}\mathcal{X}_{ij} = \mathcal{E}_{ij} - \mathcal{F}_{ij}, \quad (35)$$

with  $M_{ij}$  as in (25) (when  $i = j$ ) or (29) (when  $i \neq j$ ); the vector  $\mathcal{X}_{ij}$  has  $r$  (if  $i = j$ ) or  $2r$  (if  $i \neq j$ ) components, obtained by stacking vertically all the entries  $(X_1)_{ii}, \dots, (X_r)_{ii}$  (when  $i = j$ ) or  $(X_1)_{ij}, \dots, (X_r)_{ij}$  followed by  $(X_1)_{ji}, \dots, (X_r)_{ji}$  (when  $i \neq j$ ); the vector  $\mathcal{F}_{ij}$  is defined as

$$\mathcal{F}_{ij} := \begin{cases} w_{ii} & \text{if } i = j, \\ \begin{bmatrix} w_{ij} \\ w_{ji} \end{bmatrix} & \text{otherwise,} \end{cases}$$

where  $w_{ij}$  is given by

$$w_{ij} := \begin{bmatrix} v_{ij1} \\ \vdots \\ v_{ijr} \end{bmatrix}, \quad v_{ijk} := \sum_{\substack{s \geq i, t \geq j \\ (s,t) \neq (i,j)}} ((A_k)_{is}(X_k)_{st}(B_k)_{tj} - (C_k)_{is}(X_{k+1})_{st}(D_k)_{tj}); \quad (36)$$

and  $\mathcal{E}_{ij}$  contains all the entries in position  $(i, j)$  (when  $i = j$ ) or  $(i, j)$  and  $(j, i)$  (when  $i \neq j$ ) of  $E_1, \dots, E_r$  stacked vertically, according to the order in  $\mathcal{F}_{ij}$ . We identify  $X_{r+1}$  with  $X_1^*$  for simplicity.

Note that the values of the unknowns appearing in  $\mathcal{F}_{ij}$  have been already computed if the linear systems are solved by block back substitution in the reverse of the order in (24).

The case  $s = \mathbb{H}$  can be handled in a similar way, even if the associated system  $\mathbb{S}$  is nonlinear. In Section 6.3, we have seen how the system can be linearized over  $\mathbb{C}$  by doubling the number of equations. In order to use real arithmetic, here we follow a different approach: we consider it as a larger linear system over  $\mathbb{R}$  of double the dimension in the variables  $\text{re}(\mathcal{X}_{ij})$  and  $\text{im}(\mathcal{X}_{ij})$ . More precisely, the system  $\mathbb{S}_{ii}$ , when  $s = \mathbb{H}$ , is equivalent to the linear system over  $\mathbb{R}$  defined, for  $r > 1$ , by

$$\begin{bmatrix} \alpha_1 & \beta_1 & & & \\ & \ddots & \ddots & & \\ & & \ddots & \beta_{r-1} & \\ \beta_r & & & & \alpha_r \end{bmatrix} \begin{bmatrix} Z_1 \\ \vdots \\ Z_{r-1} \\ Z_r \end{bmatrix} = \begin{bmatrix} U_1 \\ \vdots \\ U_{r-1} \\ U_r \end{bmatrix}, \quad \begin{cases} Z_k & = \begin{bmatrix} \text{re}(X_k)_{ii} \\ \text{im}(X_k)_{ii} \end{bmatrix}, \\ U_k & = \begin{bmatrix} \text{re}((E_k)_{ii} - (v_{ii})_k) \\ \text{im}((E_k)_{ii} - (v_{ii})_k) \end{bmatrix}, \end{cases}$$

where  $\alpha_k, \beta_k$  are  $2 \times 2$  matrices defined, respectively, by

$$\begin{bmatrix} \text{re}((A_k)_{ii}(B_k)_{ii}) & -\text{im}((A_k)_{ii}(B_k)_{ii}) \\ \text{im}((A_k)_{ii}(B_k)_{ii}) & \text{re}((A_k)_{ii}(B_k)_{ii}) \end{bmatrix}, \quad - \begin{bmatrix} \text{re}((C_k)_{ii}(D_k)_{ii}) & -\text{im}((C_k)_{ii}(D_k)_{ii}) \\ \text{im}((C_k)_{ii}(D_k)_{ii}) & \text{re}((C_k)_{ii}(D_k)_{ii}) \end{bmatrix},$$

when  $k < r$ , and by

$$\begin{bmatrix} \text{re}((A_r)_{ii}(B_r)_{ii}) & -\text{im}((A_r)_{ii}(B_r)_{ii}) \\ \text{im}((A_r)_{ii}(B_r)_{ii}) & \text{re}((A_r)_{ii}(B_r)_{ii}) \end{bmatrix}, \quad - \begin{bmatrix} \text{re}((C_r)_{ii}(D_r)_{ii}) & \text{im}((C_r)_{ii}(D_r)_{ii}) \\ \text{im}((C_r)_{ii}(D_r)_{ii}) & -\text{re}((C_r)_{ii}(D_r)_{ii}) \end{bmatrix},$$

when  $k = r$ . Notice that the only differences between the two cases are the signs in the matrix on the right; this is due to the conjugation appearing in the last equation.

For  $r = 1$ , the matrix coefficient is

$$\begin{bmatrix} \operatorname{re}((A_1)_{ii}(B_1)_{ii}) & -\operatorname{im}((A_1)_{ii}(B_1)_{ii}) \\ \operatorname{im}((A_1)_{ii}(B_1)_{ii}) & \operatorname{re}((A_1)_{ii}(B_1)_{ii}) \end{bmatrix} - \begin{bmatrix} \operatorname{re}((C_1)_{ii}(D_1)_{ii}) & \operatorname{im}((C_1)_{ii}(D_1)_{ii}) \\ \operatorname{im}((C_1)_{ii}(D_1)_{ii}) & -\operatorname{re}((C_1)_{ii}(D_1)_{ii}) \end{bmatrix}.$$

The systems obtained for  $\mathbb{S}_{ij}$  are defined similarly.

We will show, in Section 8.3, that the components  $v_{ijk}$  can be computed recursively so that, for each  $(i, j)$ , the computation of  $\mathcal{F}_{ij}$  requires only  $O(nr)$  flops.

Moreover, we will show, in Section 8.4, that the system  $M_{ij}\mathcal{X}_{ij} = \mathcal{E}_{ij} - \mathcal{F}_{ij}$ , once the right-hand side term has been computed, can be solved in linear time, that is in  $O(r)$  flops, thanks to the special structure of the matrix  $M_{ij}$ .

With all the above tools we can formulate Algorithm 2 to compute the solution of a periodic system of  $r$  generalized Sylvester equations whose coefficients are in upper and lower triangular form as in Section 6.1. Besides the computation of the solution  $X_k$ , the routine also computes the matrices  $X_k B_k$  and  $X_{k+1} D_k$ , here denoted  $X_k^B$  and  $X_k^D$ , respectively, which are needed for an efficient computation of the right-hand side  $\mathcal{E}_{ij} - \mathcal{F}_{ij}$  of the linear system.

---

**Algorithm 2** Solution of a periodic system of generalized  $\star$ -Sylvester equations

---

```

1: procedure GENERALIZEDSTARSYLVESTERSYSTEM( $A_k, B_k, C_k, D_k, E_k$ )
2:   for  $k = 1, \dots, r$  do
3:      $X_k \leftarrow 0_{n \times n}$  ▷ we store the solution here
4:      $X_k^B \leftarrow 0_{n \times n}$  ▷ storage for  $X_k^B$ 
5:      $X_k^D \leftarrow 0_{n \times n}$  ▷ storage for  $X_k^D$ 
6:   end for
7:   for  $(i, j) \in \{1, 2, \dots, n\}^2$  with  $i \geq j$ , in the reverse of the ordering (24) do
8:      $\mathcal{F}_{ij} \leftarrow \text{COMPUTE}\mathcal{F}(X_k, X_k^B, X_k^D, A_k, B_k, C_k, D_k, i, j)$ 
9:      $x \leftarrow \text{SOLVE}\text{INTERMEDIATE}\text{SYSTEM}(M_{ij}, \mathcal{E}_{ij} - \mathcal{F}_{ij})$ 
10:    for  $k = 1, \dots, r$  do
11:       $[X_k]_{ij} \leftarrow x_k$ 
12:       $[X_k^B]_{ij} \leftarrow (e_i^\top X_k)(B_k e_j)$ 
13:       $[X_k^D]_{ij} \leftarrow (e_i^\top X_{k+1})(D_k e_j)$  ▷ with the convention  $X_{r+1} = X_1^*$ 
14:      if  $i \neq j$  then
15:         $[X_k]_{ji} \leftarrow x_{r+k}$ 
16:         $[X_k^B]_{ji} \leftarrow (e_j^\top X_k)(B_k e_i)$ 
17:         $[X_k^D]_{ji} \leftarrow (e_j^\top X_{k+1})(D_k e_i)$  ▷ with the convention  $X_{r+1} = X_1^*$ 
18:      end if
19:    end for
20:  end for
21:  return  $X_k$ 
22: end procedure

```

---

Section 8.3 is devoted to describe the routine `COMPUTE $\mathcal{F}$` , that computes  $\mathcal{F}_{ij}$  in the right-hand side of the systems  $\mathbb{S}_{ij}$ , while Section 8.4 describes the solution of the system, that is the routine `SOLVEINTERMEDIATESYSTEM`. An algorithmic description of the former is given in Algorithm 3, while the latter procedure is outlined in algorithmic form in the proof of Lemma 22. A FORTRAN implementation of the algorithm is available at <https://github.com/numpy/starsylv/>.



### 8.3 Computing the term $\mathcal{F}_{ij}$

The computation of the term  $\mathcal{F}_{ij}$ , if evaluated directly using Equation (36), requires  $O(n^2r)$  multiplications and additions. However, by reusing some intermediate quantities computed in the previous steps, the computation can be carried out in  $O(nr)$  flops.

We rearrange the first term in the definition of  $v_{ijk}$  (and similarly for  $v_{jik}$ ) as follows:

$$\sum_{\substack{s \geq i, t \geq j \\ (s,t) \neq (i,j)}} (A_k)_{is} (X_k)_{st} (B_k)_{tj} = \sum_{t > j} (A_k)_{ii} (X_k)_{it} (B_k)_{tj} + \sum_{s > i, t \geq j} (A_k)_{is} (X_k)_{st} (B_k)_{tj}.$$

The first summand in the right-hand side of the above equation can be computed in  $O(n)$  flops for a given  $k$ , so we only need to deal with the efficient evaluation of the latter summand. We can re-arrange it as follows:

$$\sum_{s > i, t \geq j} (A_k)_{is} (X_k)_{st} (B_k)_{tj} = \sum_{s > i} (A_k)_{is} \underbrace{\sum_{t \geq j} (X_k)_{st} (B_k)_{tj}}_{:= (X_k^B)_{sj}},$$

and this can be computed in  $O(n)$  flops if  $(X_k^B)_{sj}$ , for  $s > i$ , is known. After solving the block with indices  $\mathcal{L}_{ij}$ , we compute and store  $(X_k^B)_{ij}$  and  $(X_k^B)_{ji}$  (if different), and use them in the subsequent steps. Notice that the computation of  $(X_k^B)_{ij}$  requires only  $O(n)$  operations since  $(X_k^B)_{ij}$  is the element in position  $(i, j)$  of the product  $X_k B_k$ , and it depends only on entries of  $X_k$  that have already been computed, thanks to the triangular structure of  $B_k$ .

In Algorithm 2,  $(X_k^B)_{sj}$  has been precomputed in the previous steps, after the computation of  $(X_k)_{sj}$ . Thus, we can evaluate the first addend of  $v_{ijk}$  by computing a summation of  $O(n)$  elements, so by means of  $O(n)$  flops.

Setting  $X_{r+1} := X_1^*$ , a similar formula holds for the second term, which can be written as

$$\sum_{\substack{s \geq i, t \geq j \\ (s,t) \neq (i,j)}} (C_k)_{is} (X_{k+1})_{st} (D_k)_{tj} = \sum_{t > j} (C_k)_{ii} (X_{k+1})_{it} (D_k)_{tj} + \sum_{s > i} (C_k)_{is} \underbrace{\sum_{t > j} (X_{k+1})_{st} (D_k)_{tj}}_{:= (X_k^D)_{sj}},$$

and can be computed in  $O(n)$  by storing the computed  $(X_k^D)_{sj}$  at every step, as with  $(X_k^B)_{sj}$ .

An algorithmic description of the above process, which can be plugged directly into Algorithm 2, is given in Algorithm 3, and clearly requires  $O(nr)$  arithmetic operations. Notice that in Algorithm 3 all scalar products are computed on the complete rows and columns of the matrices  $X_1, \dots, X_r$ . This is for notational convenience, but the formulation of Algorithm 3 is equivalent to (36), thanks to the initialization to zero of  $X_k, X_k^B$ , and  $X_k^D$ , for  $k = 1, \dots, r$ . Nevertheless, in the implementation it is convenient to skip all the entries that are known to be zero.

*Remark 20.* In Algorithm 2, we have shown that it is possible to compute  $(X_k^B)_{ij}$  and  $(X_k^D)_{ij}$  after the solution of the linear system. In fact, a careful look at the algorithm shows that the scalar products

$$[X_k^B]_{ij} \leftarrow (e_i^T X_k)(B_k e_j), \quad [X_k^D]_{ij} \leftarrow (e_i^T X_{k+1})(D_k e_j)$$

can be avoided. All non-zero elements in the above summations, except the ones corresponding to the diagonal entries of  $X_k$  and  $B_k$  or  $D_k$ , are already computed and summed up in COMPUTEF. Thus, the entries in position  $(i, j)$  of  $X_k^B$  and  $X_k^D$  can be computed with an  $O(1)$  update of these

---

**Algorithm 3** Subroutines used to compute the entries of  $\mathcal{F}_{ij}$ , which is part of the right-hand side of the linear system.

---

```

1: procedure COMPUTEF( $X_k, X_k^B, X_k^D, A_k, B_k, C_k, D_k, i, j$ )
2:   if  $i = j$  then
3:      $F \leftarrow$  COMPUTEW( $X_k, X_k^B, X_k^D, A_k, B_k, C_k, D_k, i, j$ )
4:   else
5:      $F(1 : r) \leftarrow$  COMPUTEW( $X_k, X_k^B, X_k^D, A_k, B_k, C_k, D_k, i, j$ )
6:      $F(r + 1 : 2r) \leftarrow$  COMPUTEW( $X_k, X_k^B, X_k^D, A_k, B_k, C_k, D_k, j, i$ )
7:   end if
8:   return  $F$ 
9: end procedure
10: procedure COMPUTEW( $X_k, X_k^B, X_k^D, A_k, B_k, C_k, D_k, i, j$ )
11:    $F \leftarrow 0_r$ 
12:   for  $k = 1, \dots, r$  do
13:      $f_1 \leftarrow (A_k)_{ii}(e_i^\top X_k)(B_k e_j) + (e_i^\top A_k)(X_k^B e_j)$ 
14:      $f_2 \leftarrow (C_k)_{ii}(e_i^\top X_{k+1})(D_k e_j) + (e_i^\top C_k)(X_k^D e_j)$  ▷ With  $X_{r+1} = X_1^\star$ 
15:      $F_k \leftarrow f_1 + f_2$ 
16:   end for
17:   return  $F$ 
18: end procedure

```

---

partial sums. This does not change the asymptotic cost, but slightly improves the timing and it has been exploited in the implementation. However, we decided to avoid describing it in detail in the pseudocode for the sake of simplicity.

*Remark 21.* For simplicity, both here in the pseudocode and in the implementation used in the experiments, we have allocated  $2rn^2$  additional memory entries to store the matrices  $X_k^B$  and  $X_k^D$ . However, it is possible to implement the algorithm allocating with only  $O(r+n)$  additional memory if one can overwrite the input matrices  $A_k, B_k, C_k, D_k, E_k$ . Indeed, while computing the periodic Schur form as described in Lemma 15, one can use the upper triangular part of  $A_k, B_k, C_k, D_k$  to store  $\widehat{A}_k, \widehat{B}_k, \widehat{C}_k, \widehat{D}_k$  and their lower triangular parts to store in compressed format the orthogonal matrices  $Q_k, \widehat{Q}_k, Z_k, \widehat{Z}_k$ . Then, one overwrites  $E_k$  with  $\widehat{E}_k$ . Afterwards, the matrices  $Q_k, \widehat{Q}_k$  are not needed anymore, and with some index juggling one can overwrite the  $rn(n-1)$  entries used to store them with the entries of  $X_k^B$  and  $X_k^D$ , discarding those that are not needed anymore. The entries of the solution  $X_k$  can overwrite those of  $\widehat{E}_k$ .

## 8.4 Solving the small linear systems

We describe how to efficiently solve the linear system (35) involving the matrix  $M_{ij}$ . The cases  $i = j$  and  $i \neq j$  are different in the dimension of the matrix, but share the same structure, so we can handle them at the same time. More precisely, we have the following result for  $\star = \text{T}$ .

**Lemma 22.** *Let  $M$  be an  $\ell \times \ell$  matrix such that the elements in position  $(i, j)$  are allowed to be nonzero only if  $0 \leq j - i \leq 1$  or if  $(i, j) = (\ell, 1)$ . Then  $M$  admits a QR factorization  $M = QR$  where  $R$  is upper bidiagonal except in the last column, and  $Q$  is a product of  $\ell - 1$  plane rotations.*

*Proof.* The proof is constructive and by induction. The case  $\ell = 1$  is trivial, so let us assume that we have an  $(\ell + 1) \times (\ell + 1)$  matrix  $M$ , so that we can compute a rotation  $G$  acting on the

first and last row that annihilates the elements in position  $(\ell + 1, 1)$ . More precisely

$$GM = G \begin{bmatrix} \times & \times & & & \\ & \ddots & \ddots & & \\ & & \times & \times & \\ \times & & & & \times \end{bmatrix} = \left[ \begin{array}{c|cc} a_1 & b_1 & x_1 \\ \hline & \widetilde{M} & \end{array} \right],$$

where  $\widetilde{M}$  has the same shape as  $M$ , but is of size  $\ell \times \ell$ . Therefore, we can factorize  $\widetilde{M} = \widetilde{Q}\widetilde{R}$ , with  $\widetilde{Q}$  being the product of  $\ell - 1$  rotations. Setting  $Q := G^* \begin{bmatrix} 1 & 0 \\ 0 & \widetilde{Q} \end{bmatrix}$  and

$$R = \left[ \begin{array}{c|cc} a_1 & b_1 & x_1 \\ \hline & \widetilde{R} & \end{array} \right]$$

concludes the proof.  $\square$

The above proof shows that the matrices  $Q$  and  $R$  can be computed in  $O(\ell)$ , and then the linear system  $Mx = QRx = y$  can be solved in  $O(\ell)$  by the application of  $O(\ell)$  rotations to  $y$  (each of these operations can be done in  $O(1)$ ) and by a back substitution, that, thanks to the sparsity of  $R$ , can be computed in  $O(\ell)$  as well.

In our case the matrix of the linear system has  $\ell \in \{r, 2r\}$ , so we can solve each intermediate linear system in  $O(r)$ .

The case  $\star = \text{H}$  is not much different, since the matrices  $M_{i,j}$  of the linear system are block bidiagonal (except for the block at the end of the first column), with  $2 \times 2$  blocks. In fact, the matrices  $M_{i,j}$  can be brought into upper triangular form using about  $5r$  rotations, and the upper triangular form enjoys a block bidiagonal form that allows us to solve the linear system in  $O(r)$ .

Lemma 22 can be easily converted into a routine and provides a possible implementation for SOLVEINTERMEDIATESYSTEM in Algorithm 2. An implementation for this routine can be found in the code used for the tests, available at <https://github.com/numpy/starsylv/>.

## 8.5 Computational cost and storage

We evaluate the total computational cost of the algorithm (in terms of floating-point operations) by taking into account the cost of all single steps.

Step 1 requires only some bookkeeping and possibly swapping and transposing matrices in memory, but no floating point operations. This step produces several periodic systems; let  $r_1, r_2, \dots, r_m$  be their sizes, with  $r_1 + \dots + r_m \leq r$ . We prove that each of these systems is solved using  $O(n^3 r_i)$  flops.

Step 2 (for the  $i$ th periodic system of size  $r_i$ ) requires computing a periodic Schur form, which costs  $O(n^3 r_i)$  with the algorithm of [7]. Once the periodic Schur form has been computed, the changes of variables amount to  $O(r_i)$  products between  $n \times n$  matrices.

In Step 3, the method described in Section 8.3 allows one to compute each of the  $\frac{n(n+1)}{2}$  terms  $\mathcal{F}_{i,j}$  in  $O(nr_i)$  flops, and Section 8.4 shows how to solve in  $O(r_i)$  flops the linear systems required in each of the  $\frac{n(n+1)}{2}$  back substitution steps. The total amount of flops required by this step is, thus,  $O(n^3 r_i)$ .

Step 4 requires applying formula (14) (which costs  $O(n^3)$  to compute) once for each remaining variable, that is, at most  $r - 1$  times.

Combining all the above steps we obtain an algorithm with a total cost of  $O(n^3 r)$  flops. Moreover, the only storage required is the one of  $O(r)$  matrices of size  $n \times n$ , so the storage required is  $O(n^2 r)$ , which is optimal (given that the same amount of storage is required to store the solutions).

*Remark 23.* Step 1 requires some discrete computations on the indices to identify the periodic systems and eliminate variables and equations; we have ignored them here since they involve no floating-point operations, but we have shown in Section 8.1 that they can be performed in  $O(r)$  operations with the help of a graph algorithm.

## 8.6 Backward error analysis

Here we provide a backward error analysis of the algorithm described in the previous sections. We use the standard floating point number model with unit roundoff  $u$  and, for an expression  $\ell$ , we denote by  $\text{fl}(\ell)$  the computed value of  $\ell$  using floating point operations. We use the notation

$$\gamma_k := \frac{cku}{1 - cku},$$

where  $c$  denotes a small constant, whose exact value is not relevant (see [23, p. 68]).

We assume that all linear systems  $Ax = b$  that are encountered are solved using a backward stable method. More precisely, we say that an algorithm to solve a linear system  $Ax = b$ , with  $A \in \mathbb{C}^{m \times m}$ , has *backward error*  $\varepsilon_A$  if the computed solution  $\tilde{x} = \text{fl}(A^{-1}b)$  is the exact solution of a perturbed system  $(A + \delta A)\tilde{x} = b$ , with  $\|\delta A\|_2 / \|A\|_2 \leq \varepsilon_A$ . Note that only the coefficient matrix is perturbed (see [23, Th. 19.5] and the following discussion for an explanation). In the case of solving the system with the QR factorization using  $s$  Givens rotations, as we do in Section 8.4 with  $s = O(r)$ , this quantity can be taken as  $\varepsilon_A = m \cdot \gamma_s$  (see p. 368 and Theorem 19.10 in [23]). The factor  $m$  comes from the fact that the bound in [23] is only given column-wise and

$$\|\text{Col}_j A\|_2 \leq \|A\|_2 \leq \sqrt{m} \|A\|_1 = \sqrt{m} \max_{j=1, \dots, m} \|\text{Col}_j A\|_1 \leq m \max_{j=1, \dots, m} \|\text{Col}_j A\|_2, \quad (37)$$

for all  $j = 1, \dots, m$ , where  $\text{Col}_j A$  is the  $j$ th column of  $A$  (see, for instance, [23, Tables 6.1 and 6.2] for the last two inequalities).

We obtain a backward error result formulating the problem as a vectorized linear system. For simplicity, we will focus on periodic systems with upper and lower triangular coefficients in Theorem 24. The general case will be commented right after the proof.

**Theorem 24.** *Consider a system of equations of the form (4), with  $A_k, C_k, B_k^\top, D_k^\top$  being upper triangular, and let  $M\mathcal{X} = \mathcal{E}$  be its vectorized form, where  $M \in \mathbb{C}^{rn^2 \times rn^2}$  if  $\star = \text{T}$ , or  $M \in \mathbb{R}^{2rn^2 \times 2rn^2}$  if  $\star = \text{H}$ .*

*When implemented in standard floating-point arithmetic, the algorithm described in Sections 8.2–8.4 produces a result  $\tilde{\mathcal{X}}$  satisfying*

$$(M + \delta M)\tilde{\mathcal{X}} = \mathcal{E} + \delta \mathcal{E}, \quad (38)$$

*with  $\|\delta M\|_2 / \|M\|_2 \leq r \gamma_r + \gamma_n^2 (1 + r \gamma_r)$ ,  $\|\delta \mathcal{E}\|_2 / \|\mathcal{E}\|_2 \leq \gamma_n^2$ .*

*Remark 25.* The reader may wonder if a stronger form of structured backward stability holds: the algorithm should produce matrices that satisfy

$$(A_k + \delta A_k) \tilde{X}_{\alpha_k}^{s_k} (B_k + \delta B_k) - (C_k + \delta C_k) \tilde{X}_{\beta_k}^{t_k} (D_k + \delta D_k) = E_k + \delta E_k \quad k = 1, \dots, r,$$

with  $\|\delta S_k\|_2/\|S_k\|_2$  being small, for  $S = A, B, C, D, E$ . Unfortunately, algorithms of this family fail to be structurally backward stable even in the simplest case of a single Sylvester equation  $AX - XD = E$ , as shown in [22, §16.2] (see also the discussion in [9] for the case  $s = 1$ ).

Note that Theorem 24 is nevertheless sufficient to show that the residual of each equation  $R_k = \|A_k \tilde{X}_{\alpha_k}^{s_k} B_k - C_k \tilde{X}_{\beta_k}^{t_k} D_k - E_k\|_F$ , for  $k = 1, 2, \dots, r$ , is small. Indeed,  $\|M\tilde{\mathcal{X}} - \mathcal{E}\|_2 = \sqrt{\sum_{k=1}^r R_k^2}$  satisfies (see [23, Thm 7.1])

$$\frac{\|M\tilde{\mathcal{X}} - \mathcal{E}\|_2}{\|M\|_2\|\tilde{\mathcal{X}}\|_2 + \|\mathcal{E}\|_2} \leq \max\left(\frac{\|\delta M\|_2}{\|M\|_2}, \frac{\|\delta \mathcal{E}\|_2}{\|\mathcal{E}\|_2}\right)$$

In order to prove Theorem 24, we need the following technical results.

**Lemma 26.** *Let  $N \in \mathbb{C}^{m \times m}$  and  $x, y \in \mathbb{C}^m$ , with  $x, y \neq 0$ , be such that*

$$y = (N + \Delta N)x, \quad \frac{\|\Delta N\|_2}{\|N\|_2} \leq \varepsilon, \quad (39)$$

for some  $\varepsilon > 0$ . Let  $\delta y \in \mathbb{C}^m$  be such that

$$\frac{\|\delta y\|_2}{\|y\|_2} \leq \kappa, \quad (40)$$

for some  $\kappa > 0$ . Then  $y + \delta y = (N + \delta N)x$ , for some  $\delta N \in \mathbb{C}^{m \times m}$  with  $\frac{\|\delta N\|_2}{\|N\|_2} \leq \varepsilon + \kappa(1 + \varepsilon)$ .

*Proof.* From (39) and (40) we get

$$\|\delta y\|_2 \leq \kappa\|y\|_2 \leq \kappa(\|N\|_2 + \|\Delta N\|_2)\|x\|_2 \leq \kappa(1 + \varepsilon)\|N\|_2\|x\|_2. \quad (41)$$

Now, setting  $\tilde{N} := \|x\|_2^{-2} \cdot (\delta y)x^H$ , we have  $\tilde{N}x = \delta y$  and  $\|\tilde{N}\|_2 = \|\delta y\|_2/\|x\|_2$ , so  $\|\delta y\|_2 = \|\tilde{N}\|_2\|x\|_2$ . Then, by (41),

$$\|\tilde{N}\|_2 \leq \kappa(1 + \varepsilon)\|N\|_2. \quad (42)$$

Finally, taking  $\delta N := \Delta N + \tilde{N}$ , and using (42), we arrive at  $\|\delta N\|_2 \leq \|\Delta N\|_2 + \|\tilde{N}\|_2 \leq (\varepsilon + \kappa(1 + \varepsilon))\|N\|_2$ .  $\square$

**Lemma 27.** *Consider a square linear system of the form  $Fx = b - \sum_{k=1}^s N_k c_k$ , where  $F, N_k \in \mathbb{C}^{m \times m}$ , and  $b, c_k \in \mathbb{C}^m$  are given, for  $k = 1, \dots, s$ , and  $x$  is the unknown.*

*Forming the sum in the right-hand side, in floating point arithmetic, and then solving the linear system using an algorithm with backward error  $\varepsilon_F$ , produces a computed solution  $\tilde{x}$  which is the exact solution of a perturbed system*

$$(F + \delta F)\tilde{x} = b + \delta b - \sum_{k=1}^s (N_k + \delta N_k)c_k,$$

with

$$\frac{\|\delta F\|_2}{\|F\|_2} \leq \varepsilon_F, \quad \frac{\|\delta b\|_2}{\|b\|_2} \leq \gamma_s, \quad \frac{\|\delta N_k\|_2}{\|N_k\|_2} \leq m\gamma_m + \gamma_s(1 + m\gamma_m).$$

*Proof.* Let  $\tilde{d}_k = \text{fl}(N_k c_k)$ ,  $\tilde{f} = \text{fl}(b - \sum_{k=1}^s \tilde{d}_k)$ . By hypothesis,  $(F + \delta F)\tilde{x} = \tilde{f}$ , with  $\|\delta F\|_2/\|F\|_2 \leq \varepsilon_F$ . The usual backward error analysis of summation can be used to show that  $\tilde{f} = b +$

$\delta b - \sum_{k=1}^s (\tilde{d}_k + \delta \tilde{d}_k)$ , with  $|(\delta b)_i|/|b_i|, |(\delta \tilde{d}_k)_i|/|(\tilde{d}_k)_i| \leq \gamma_s$ , for  $i = 1, \dots, m$  (see [23, Section 4]). Now, by standard backward error analysis of matrix-vector multiplication, we know that  $\tilde{d}_k = (N_k + \Delta N_k)c_k$ , with  $\|\text{Col}_j(\Delta N_k)\|_2/\|\text{Col}_j(N_k)\|_2 \leq \gamma_m$ , for  $j = 1, \dots, m$  (see [23, Section 3.5]). Using (37), this implies  $\|\Delta N_k\|_2/\|N_k\|_2 \leq m\gamma_m$ . Now, we can apply Lemma 26, with  $y = \tilde{d}_k, \delta y = \delta \tilde{d}_k, x = c_k, N = N_k$  and  $\Delta N = \Delta N_k$ , to conclude that  $\tilde{d}_k + \delta \tilde{d}_k = (N_k + \delta N_k)c_k$ , with  $\|\delta N_k\|_2/\|N_k\|_2 \leq m\gamma_m + \gamma_s(1 + m\gamma_m)$ , as wanted.  $\square$

*Proof of Theorem 24.* We note that each step of the block back substitution corresponds to solving a linear system of the form (35). More precisely, this system is

$$M_{ij}\mathcal{X}_{ij} = \mathcal{E}_{ij} - \sum_{(s,t) \in \mathcal{U}_{ij}} N_{st}^{(ij)}\mathcal{X}_{st},$$

where  $\mathcal{U}_{ij} = \{(i', j') : \max\{i', j'\} \geq \max\{i, j\} \text{ and } \min\{i', j'\} \geq \min\{i, j\}\}$  and the matrices  $N_{st}^{(ij)}$  are given by writing (36) in matrix form. By Lemma 27, there are some matrices  $\delta M_{ij}$  and  $\delta N_{st}^{(ij)}$  such that  $(M_{ij} + \delta M_{ij})\tilde{\mathcal{X}}_{ij} = \mathcal{E}_{ij} + \delta \mathcal{E}_{ij} - \sum_{(s,t) \in \mathcal{U}_{ij}} (N_{st}^{(ij)} + \delta N_{st}^{(ij)})\tilde{\mathcal{X}}_{st}$ , where  $\tilde{\mathcal{X}}_{ij}$  are the computed solutions at the  $(i, j)$  step and  $\tilde{\mathcal{X}}_{st}$ , for  $s \geq i, t \geq j$ , with  $(s, t) \neq (i, j)$ , are the ones computed in the previous steps, and

$$\frac{\|\delta M_{ij}\|_2}{\|M_{ij}\|_2} \leq \varepsilon_{M_{ij}}, \quad \frac{\|\delta N_{st}^{(ij)}\|_2}{\|N_{st}^{(ij)}\|_2} \leq r\gamma_r + \gamma_n^2(1 + r\gamma_r), \quad \frac{\|\delta \mathcal{E}_{ij}\|_2}{\|\mathcal{E}_{ij}\|_2} \leq \gamma_n^2.$$

If the  $r \times r$  (or  $(2r) \times (2r)$ ) linear system is solved through the QR factorization of  $M_{ij}$ , then  $\varepsilon_{M_{ij}} \leq r\gamma_r$ , as mentioned before (see [23, Th. 19.10]).

This gives a backward error for each block-row of the matrix  $M$  and of the right-hand side  $\mathcal{E}$  in Theorem 24. Since these rows are never reused between equations, this defines a perturbation of  $M$  and  $\mathcal{E}$  which ensures (38).  $\square$

We note that Theorem 24 corresponds to Step 3 in the procedure described at the beginning of Section 8 for solving a general system (3). The remaining steps can be carried out also in a backward stable way, as we are going to explain.

Step 1 involves no computations, just relabeling of the equations, transpositions and conjugations (which are exact in floating point arithmetic).

Step 2 is backward stable since the periodic QZ algorithm relies on unitary transformations and the following change of variables is unitary.

In Step 4, the vectorization of (14) produces the linear system

$$(B_k^\top \otimes A_k) \text{vec}(X_{\alpha_k}^{s_k}) = \text{vec}(E_k) + (D_k^\top \otimes C_k) \text{vec}(X_{\beta_k}^{t_k}),$$

which is again in the form treated in Lemma 27, so we only have to ensure that the method used to solve this linear system of the form  $(B_k^\top \otimes A_k) \text{vec}(X) = \text{vec}(F)$  is backward stable. To solve this system, we first compute  $\tilde{Y} = \text{fl}(A_k^{-1}F)$  column by column, each time solving a linear system with  $A_k$ , and then similarly  $\tilde{X} = \text{fl}(\tilde{Y}B_k^{-1})$ , solving a linear system for each of its rows.

We assume that the linear systems with  $A_k$  are solved with a backward stable method, i.e.,

$$(A_k + \delta_j A_k) \text{Col}_j(\tilde{Y}) = \text{Col}_j(F), \quad \frac{\|\delta_j A_k\|_2}{\|A_k\|_2} \leq \varepsilon_{A_k},$$

(note that there is a different perturbation  $\delta_j A_k$  for each  $j$ ); hence we have

$$(\mathbb{A} + \delta \mathbb{A}) \text{vec}(\tilde{Y}) = \text{vec}(F), \quad \frac{\|\delta \mathbb{A}\|_2}{\|\mathbb{A}\|_2} \leq \varepsilon_{A_k},$$

where  $\mathbb{A} = I_n \otimes A_k$  and  $\delta\mathbb{A} = \text{diag}(\delta_1 A_k, \dots, \delta_n A_k)$ .

An analogous argument shows that

$$(\mathbb{B} + \delta\mathbb{B}) \text{vec}(\tilde{X}) = \text{vec}(\tilde{Y}), \quad \frac{\|\delta\mathbb{B}\|_2}{\|\mathbb{B}\|_2} \leq \varepsilon_{B_k},$$

where  $\mathbb{B} = B_k^\top \otimes I_n$ . Combining these two relations we have  $\text{vec}(F) = (\mathbb{A} + \delta\mathbb{A})(\mathbb{B} + \delta\mathbb{B}) \text{vec}(\tilde{X}) = (\mathbb{A}\mathbb{B} + \delta(\mathbb{A}\mathbb{B})) \text{vec}(\tilde{X})$ , with  $\delta(\mathbb{A}\mathbb{B}) = \delta\mathbb{A} \cdot \mathbb{B} + \mathbb{A} \cdot \delta\mathbb{B} + \delta\mathbb{A} \cdot \delta\mathbb{B}$ . Since  $\|\mathbb{A}\mathbb{B}\|_2 = \|\mathbb{A}\|_2 \|\mathbb{B}\|_2$  for our choice of  $\mathbb{A}$  and  $\mathbb{B}$  (thanks to the properties of the Kronecker product [24, p. 253]), we get

$$\begin{aligned} \frac{\|\delta(\mathbb{A}\mathbb{B})\|_2}{\|\mathbb{A}\mathbb{B}\|_2} &= \frac{\|\delta\mathbb{A} \cdot \mathbb{B} + \mathbb{A} \cdot \delta\mathbb{B} + \delta\mathbb{A} \cdot \delta\mathbb{B}\|_2}{\|\mathbb{A}\|_2 \|\mathbb{B}\|_2} \leq \frac{\|\delta\mathbb{A}\|_2 \|\mathbb{B}\|_2 + \|\mathbb{A}\|_2 \|\delta\mathbb{B}\|_2 + \|\delta\mathbb{A}\|_2 \|\delta\mathbb{B}\|_2}{\|\mathbb{A}\|_2 \|\mathbb{B}\|_2} \\ &\leq \varepsilon_{A_k} + \varepsilon_{B_k} + \varepsilon_{A_k} \varepsilon_{B_k}. \end{aligned}$$

As a consequence of these arguments, the procedure described at the beginning of Section 8 produces a backward stable algorithm for solving general systems of the form (3).

## 8.7 Numerical experiments

We have implemented the proposed algorithm for the solution in the case  $\star = \mathbb{T}$ . The case  $\star = \mathbb{H}$  can be obtained with minimal changes (from the algorithmic point of view), so we decided to avoid running the same experiments concerning stability and performance. We have run the tests on a server with a Xeon X5680 CPU and 24 GB of memory. Our implementation is available at <https://github.com/numpi/starsylv/>. The code has been compiled with GNU Fortran compiler and linked with the (single-threaded) BLAS reference implementation (`libblas.so`, <http://www.netlib.org/blas/>).

We have computed the CPU time required by our implementation as a function of the size of the matrices  $n$  and of the number of equations in the reduced system  $r$ , and we have compared it with the behavior predicted by our analysis. We have considered only systems with triangular factors. The general case requires the reduction to triangular factors through the periodic Schur form as described in Section 6, which has been already implemented in [5, subroutines MB03BD and MB03BZ] (see also [7, 26]).

The results are reported in Figure 1, on the left, for the CPU time required for the solution of a system of three equations with coefficients of variable size  $n$ , and on the right for a system of  $r$  equations of size 16. Both plots confirm the cubic and linear dependence of the CPU time on the parameters  $n$  and  $r$ , respectively, that we expect. The dashed lines in the two plots are obtained plotting the functions  $k_n n^3$  and  $k_r r$  for two appropriate constants  $k_n$  and  $k_r$ .

Beside timings, we have also tested the accuracy of the implementation. For each value of  $n$  and  $r$  we have generated several systems of  $\mathbb{T}$ -Sylvester equations (in the required triangular form), and we have computed the residuals  $R_k := \|A_k X_k B_k - C_k X_{k+1} D_k - E_k\|_F$  for  $k = 1, \dots, r-1$ , and  $R_r := \|A_r X_r B_r - C_r X_1^\top D_r - E_r\|_F$ . Then, the 2-norm of the residual of the linear system can be evaluated as  $R := \sqrt{R_1^2 + \dots + R_r^2}$ . In Figure 2 we have plotted an upper bound of the relative residuals  $R/\|M\|_2$ , obtained using the relation  $n\sqrt{r}\|M\|_2 \geq \|M\|_F$ , where  $M$  is the matrix of the ‘‘large’’ linear system, for different values of  $n$  and  $r$  (recall that  $M$  has size  $n^2 r$ ). Each value has been averaged over 100 runs. The Frobenius norm of  $M$  is easily computable recalling that, if two matrices  $M_1$  and  $M_2$  do not have non-zero entries in corresponding positions, then  $\|M_1 + M_2\|_F^2 = \|M_1\|_F^2 + \|M_2\|_F^2$ , and the relation  $\|A \otimes B\|_F = \|A\|_F \|B\|_F$ .

In these tests, the coefficients matrices  $A_k, B_k, C_k, D_k$  have been chosen with random entries with normal distribution, and with the correct triangular structure. We have then shifted  $A_k$  and  $B_k$  with  $\sqrt{n}I$  to avoid finding solutions with very large norms.

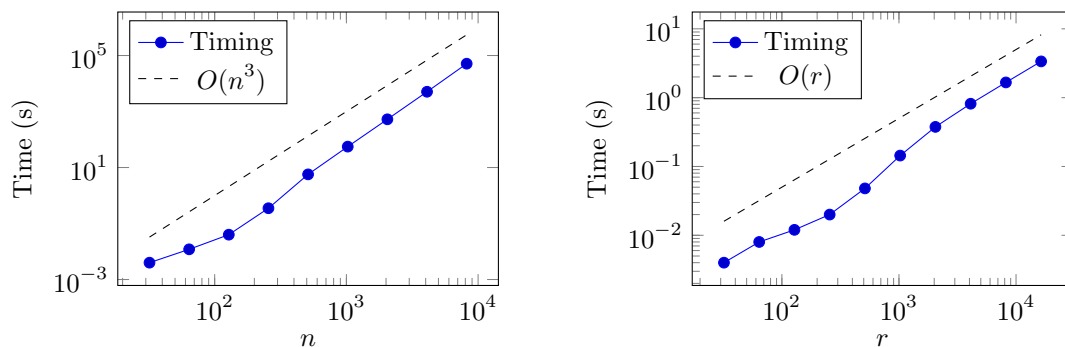


Figure 1: On the left, the CPU time required by the algorithm described in Section 8.6 for the  $\star = \text{T}$  case, as a function of  $n$ . The timings reported are for a system with 3 equations, already in the required triangular form. The problems tested have sizes ranging from  $n = 32$  to  $n = 8192$ . On the right, the CPU time required by the algorithm described in Section 8.6 for the  $\star = \text{T}$  case, as a function of  $r$ . The timings reported are for a system with  $r$  equations and coefficients of size  $16 \times 16$ , already in the required triangular form. The problems tested have sizes ranging from  $r = 32$  to  $r = 16384$ .

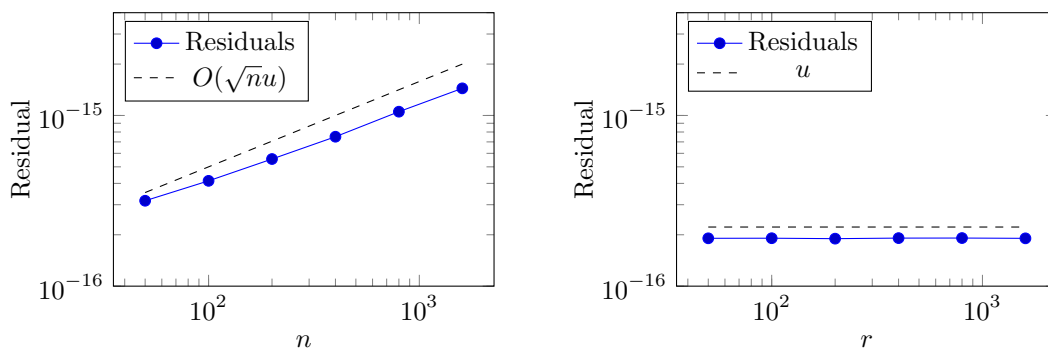


Figure 2: On the left, average residuals of 100 systems of T-Sylvester equations solved via the algorithm described in Section 8. The systems considered have 3 equations with a variable coefficient size  $n$ . On the right, average residuals of 100 systems of T-Sylvester equations solved via the algorithm described in Section 8. The systems considered have coefficients with size  $8 \times 8$ , and  $r$  equations.

From the tests performed so far, the algorithm behaves in a backward stable way, as predicted by our analysis. In fact, one can spot that the error growth with respect to  $n$  and  $r$  is even less than the upper bound proved in this section. The error seems to grow slightly less than  $\sqrt{n}$ , and to be independent of  $r$ . This behavior is often encountered in dense linear algebra algorithms, since on average the errors do not accumulate in the same direction (see e.g. [23, Section 4.5]).

## 9 Conclusions and future work

We have provided necessary and sufficient conditions for the nonsingularity of  $r$  coupled generalized Sylvester and  $\star$ -Sylvester equations (3), with square coefficients of the same size  $n \times n$ . We have shown that, in the nonsingular case, the problem can be reduced to periodic systems having at most one generalized  $\star$ -Sylvester equation. A characterization for the nonsingularity of periodic systems of just generalized Sylvester equations was obtained in an unpublished work by Byers and Rhee [9]. That characterization was given in terms of spectral properties of matrix



pencils constructed from the coefficients of the system. We have provided an analogous characterization for the nonsingularity of periodic systems with exactly one generalized  $\star$ -Sylvester equation. We have also provided a characterization for both types of periodic systems (namely, the one with exactly one generalized  $\star$ -Sylvester equation and the one with only generalized Sylvester equations) in terms of spectral properties of formal products constructed from the coefficients of the system. Finally, we have presented an  $O(n^3r)$  algorithm for computing the unique solution of a nonsingular system, which has been shown to be backward stable.

A future research line that naturally arises from this work is to get a characterization of nonsingularity in the more general setting of rectangular coefficients. Other possible generalizations, pointed out by the referees, include systems involving complex conjugation of the unknowns, like those considered in [17, 36], or systems of periodic type [6, 29].

## Acknowledgments

We wish to thank the anonymous referees for their comments that helped us to improve the presentation. This work does not have any conflicts of interest.

## References

- [1] P. Anderson, R. Granat, I. Jonsson, and B. K. gström. Parallel algorithms for triangular periodic Sylvester-type matrix equations. In *Lecture Notes in Computer Science*, pages 169–174. Euro-Par 2008–Parallel Processing, Springer, 2007.
- [2] J. L. Aurentz, T. Mach, L. Robol, R. Vandebril, and D. S. Watkins. *Core-Chasing Algorithms for the Eigenvalue Problem*, volume 13 of *Fundamentals of Algorithms*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2018.
- [3] J. L. Aurentz, T. Mach, L. Robol, R. Vandebril, and D. S. Watkins. Fast and backward stable computation of the eigenvalues of matrix polynomials. *Math. Comput.*, 2018.
- [4] Z. Bai and J. W. Demmel. On swapping diagonal blocks in real Schur form. *Linear Algebra Appl.*, 186:75–95, 1993.
- [5] P. Benner, V. Mehrmann, V. Sima, S. Van Huffel, and A. Varga. SLICOT — a subroutine library in systems and control theory. In B. N. Datta, editor, *Applied and Computational Control, Signals, and Circuits (1997)*, chapter 10, pages 499–539. Birkhäuser Boston, Boston, MA, 1997.
- [6] D. A. Bini, B. Iannazzo, and F. Poloni. A fast Newton’s method for a nonsymmetric algebraic Riccati equation. *SIAM J. Matrix Anal. Appl.*, 30(1):276–290, 2008.
- [7] A. W. Bojanczyk, G. H. Golub, and P. Van Dooren. Periodic Schur decomposition: algorithms and applications. In *Proc. SPIE Conference*, pages 31–42. International Society for Optics and Photonics, 1992.
- [8] R. Byers and D. Kressner. Structured condition numbers for invariant subspaces. *SIAM J. Matrix Anal. Appl.*, 28(2):326–347, 2006.
- [9] R. Byers and N. Rhee. Cyclic Schur and Hessenberg-Schur numerical methods for solving periodic Lyapunov and Sylvester equations. Technical report, Dept. of Mathematics, Univ. of Missouri at Kansas City, 1995.

- [10] C.-Y. Chiang, E. K.-W. Chu, and W.-W. Lin. On the  $\star$ -Sylvester equation  $AX \pm X^*B^* = C$ . *Appl. Math. Comput.*, 218:8393–8407, 2012.
- [11] K.-W. E. Chu. The solution of the matrix equations  $AXB - CXD = E$  and  $(YA - DZ, YC - BZ) = (E, F)$ . *Linear Algebra Appl.*, 93:93–105, 1987.
- [12] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to algorithms*. MIT Press, Cambridge, MA, third edition, 2009.
- [13] F. De Terán and F. M. Dopico. Consistency and efficient solution for the Sylvester equation for  $\star$ -congruence:  $AX + X^*B = C$ . *Electron. J. Linear Algebra*, 22:849–863, 2011.
- [14] F. De Terán, F. M. Dopico, N. Guillery, D. Montealegre, and N. Z. Reyes. The solution of the equation  $AX + X^*B = 0$ . *Linear Algebra Appl.*, 483:2817–2860, 2013.
- [15] F. De Terán and B. Iannazzo. Uniqueness of solution of a generalized  $\star$ -Sylvester matrix equation. *Linear Algebra Appl.*, 493:323–335, 2016.
- [16] F. De Terán, B. Iannazzo, F. Poloni, and L. Robol. Solvability and uniqueness criteria for generalized Sylvester-type equations. *Linear Algebra Appl.*, 542:501–521, 2018.
- [17] A. Dmytryshyn, V. Futorny, T. Klymchuk, and V. V. Sergeichuk. Generalization of Roth’s solvability criteria to systems of matrix equations. *Linear Algebra Appl.*, 527:294–302, 2017.
- [18] A. Dmytryshyn and B. Kågström. Coupled Sylvester-type matrix equations and block diagonalization. *SIAM J. Matrix Anal. Appl.*, 36(2):580–593, 2016.
- [19] F. M. Dopico, J. González, D. Kressner, and V. Simoncini. Projection methods for large-scale  $T$ -Sylvester equations. *Math. Comput.*, 85:2427–2455, 2016.
- [20] R. Granat and B. K. gström. Direct eigenvalue reordering in a product of matrices in periodic Schur form. *SIAM J. Matrix Anal. Appl.*, 28:285–300, 2006.
- [21] R. Granat, B. K. gström, and D. Kressner. Computing periodic deflating subspaces associated with a specified set of eigenvalues. *BIT*, 47:763–791, 2007.
- [22] N. J. Higham. Perturbation theory and backward error for  $AX - XB = C$ . *BIT*, 33(1):124–136, 1993.
- [23] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. SIAM, Philadelphia, PA, 1996.
- [24] R. A. Horn and C. R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, Cambridge, MA, 1994.
- [25] T. Košir. Kronecker bases for linear matrix equations, with application to two-parameter eigenvalue problems. *Linear Algebra Appl.*, 249:259–288, 1996.
- [26] D. Kressner. An efficient and reliable implementation of the periodic QZ algorithm. In *IFAC Workshop on Periodic Control Systems (PSYCO 2001), Como (Italy)*, pages 31–42. IFAC, 2001.
- [27] D. Kressner, E. Mengi, I. Nakić, and N. Truhar. Generalized eigenvalue problems with specified eigenvalues. *IMA J. Numer. Anal.*, 34:480–501, 2014.

- [28] D. Kressner, C. Schröder, and D. S. Watkins. Implicit QR algorithms for palindromic and even eigenvalue problems. *Numer. Algorithms*, 51(2):209–238, 2009.
- [29] I. Kuzmanović and N. Truhar. Sherman-Morrison-Woodbury formula for Sylvester and  $T$ -Sylvester equations with applications. *Int. J. Comput. Math.*, 90(2):306–324, 2013.
- [30] S. K. Mitra. The matrix equation  $AXB + CXD = E$ . *SIAM J. Appl. Math.*, 32:823–825, 1977.
- [31] T. Rees and J. Scott. A comparative study of null-space factorizations for sparse symmetric saddle point systems. *Numer. Linear Algebra Appl.*, 25(1):e2103, 2018.
- [32] V. Simoncini. Computational methods for linear matrix equations. *SIAM Rev.*, 58(3):377–441, 2016.
- [33] A. Varga and P. Van Dooren. Computational methods for periodic systems—an overview. In *Proc. IFAC Workshop on Periodic Control Systems*, pages 171–176. IFAC, 2001.
- [34] D. S. Watkins. Product eigenvalue problems. *SIAM Rev.*, 47(1):3–40, 2005.
- [35] M. Wedderburn. Note on the linear matrix equation. *Proc. Edinburgh Math. Soc.*, 22:49–53, 1904.
- [36] A.-G. Wu and Y. Zhang. *Complex conjugate matrix equations for systems and control*. Communications and Control Engineering Series. Springer, Singapore, 2017.