**RESEARCH ARTICLE** 

Revised: 22 February 2022

WILEY

# Preconditioners for fractional diffusion equations based on the spectral symbol

### Nikos Barakitis<sup>1</sup> | Sven-Erik Ekström<sup>1,2</sup> | Paris Vassalos<sup>1</sup>

<sup>1</sup>Department of Informatics, Athens University of Economics and Business, Athens, Greece

<sup>2</sup>Department of Information Technology, Division of Scientific Computing, Uppsala University, Uppsala, Sweden

#### Correspondence

Paris Vassalos, Department of Informatics, Athens University of Economics and Business, 76 Patission Str., GR-10434 Athens, Greece. Email: pvassal@aueb.gr

**Funding information** Athens University of Economics and Business

#### Abstract

It is well known that the discretization of fractional diffusion equations with fractional derivatives  $\alpha \in (1, 2)$ , using the so-called weighted and shifted Grünwald formula, leads to linear systems whose coefficient matrices show a Toeplitz-like structure. More precisely, in the case of variable coefficients, the related matrix sequences belong to the so-called generalized locally Toeplitz class. Conversely, when the given FDE has constant coefficients, using a suitable discretization, we encounter a Toeplitz structure associated to a nonnegative function  $\mathcal{F}_{\alpha}$ , called the spectral symbol, having a unique zero at zero of real positive order between one and two. For the fast solution of such systems by preconditioned Krylov methods, several preconditioning techniques have been proposed in both the one- and two-dimensional cases. In this article we propose a new preconditioner denoted by  $\mathcal{P}_{\mathcal{F}_{\alpha}}$  which belongs to the  $\tau$ -algebra and it is based on the spectral symbol  $\mathcal{F}_{\alpha}$ . Comparing with some of the previously proposed preconditioners, we show that although the low band structure preserving preconditioners are more effective in the one-dimensional case, the new preconditioner performs better in the more challenging multi-dimensional setting.

#### K E Y W O R D S

fractional differential equations, fractional order zero, GMRES, multi-level Toeplitz matrix, sine transform based preconditioner

### **1** | INTRODUCTION

Fractional calculus may be considered both an old and modern topic. Old, since it dates back to the letter from L'Hôpital to Leibniz in 1695, and a novel one, since it has been object of specialized conferences and treatises, for the last 40 years. In recent years considerable interest in fractional calculus has been stimulated by the applications that this calculus finds in numerical analysis and modeling. As an example, fractional diffusion equations (FDEs) are used to model anomalous diffusion or dispersion, where a particle plume spreads at a rate inconsistent with the classical Brownian motion model (e.g., see Reference 1 and the references therein). Such phenomena are ubiquitous in both natural and social sciences. In fact, many complex dynamical systems often contain anomalous diffusion. Fractional kinetic equations are usually an effective method for describing these complex systems, including diffusion type, diffusive convection type and Fokker–Planck type FDEs.<sup>2</sup> Since analytical solutions are rarely available, these kinds of equations are of numerical interest. When the order of fractional derivatives is  $\alpha = 1$ , we have the standard diffusion process. With  $0 < \alpha < 1$ , we describe a sub-diffusion process or dispersive, slow diffusion process with the anomalous diffusion index, while with  $\alpha > 1$ , an ultra-diffusion

<sup>2 of 22</sup> WILEY

process or increased, fast diffusion process. In Reference 3 it has been proved the strict relationship between the order of the fractional derivative and the order of the zero of the associated symbol of the coefficient matrix of the related system. In addition, it is well known (e.g., see References 4,5), that if the generating function  $f \ge 0$  of a Toeplitz matrix of size *n* has a unique zero of order  $\alpha \in (0, 1)$  and it is bounded, then  $T_n(f)$  has a condition number growing exactly as  $n^{\alpha}$ . Hence, the standard Conjugate Gradient method requires only  $O(n^{\alpha/2})$  iterations for the solution of a related linear systems up to the required accuracy (for a total of  $O(n^{\alpha/2} \log(n))$  arithmetic operations. When  $\omega f$  has all the range in the right complex plane for some  $\omega$  complex of modulus one, the generating function *f* is called weakly sectorial.<sup>6,7</sup> Then if *f* is weakly sectorial, essentially bounded and *f* has a unique zero of order  $\alpha \in (0, 1)$ , then again  $T_n(f)$  has a condition number growing exactly as  $n^{\alpha}$ . Hence, a good GMRES with (possibly) any standard circulant preconditioning is essentially satisfactory (e.g., see References 8-10). Thus, the case where the order of fractional derivative  $\alpha$  belongs to the interval (1, 2) is, computationally, more challenging. Moreover, in this article we are focus on the numerical solution of particular time-dependent space-fractional diffusion equation on rectangular domains in one and two dimensions using finite differences techniques. For numerical techniques concerning domains of general geometry or numerical schemes different from finite differences and multigrid techniques, the interested reader is referred, for example, to Reference 11-13, and the references therein.

Several definitions for the fractional derivative exist, and each of them approaches the definition of ordinary derivative in the integer order limit. In References 1,14 the authors proposed two unconditionally stable finite difference schemes, of first and second order accuracy, based on the shifted Grünwald–Letnikov definition of fractional derivatives. In Reference 15 it was shown that once one of these methods is chosen, the coefficient matrix of the generated system can be seen as the sum of two structures, each of them expressed as a diagonal matrix multiplied by a Toeplitz matrix. Since the efficient solution of such systems are of great interest many iterative solvers have been proposed. Representative examples are the multigrid method (MGM) scheme proposed by Noutsos and Vassalos,<sup>16</sup> the circulant-based preconditioners for the Conjugate Gradient Normal Residual (CGNR) method,<sup>17,18</sup> the splitting preconditioner,<sup>19</sup> and two structure-preserving preconditioners proposed in Reference 3. In the latter paper, the authors provide a detailed analysis, showing that the sequence of coefficient matrices belongs to the generalized locally Toeplitz (GLT) class and its spectral symbol, which describes the asymptotic singular and eigenvalue distribution, is explicitly derived. In Reference 20 the analysis is extended to the two-dimensional case and the authors compare the two-dimensional version of the structure preserving preconditioner based on a decomposition of the Laplacian<sup>3</sup> to a preconditioner based on an algebraic MGM.

In this work, based on the theoretical results presented in Reference 21 and motivated by an interest to study the effectiveness of suitable  $\tau$  preconditioners for ill-conditioned symmetric Toeplitz systems, we propose a new preconditioner for the solution of Toeplitz-like systems, stemming from the discretization of the considered FDEs. Specifically, in Reference 21 the authors proved the essential spectral equivalence between the matrix sequences  $\{T_n(f)\}_n$  and  $\{\tau_n(f)\}_n$ , where  $\{T_n(f)\}_n$  is the sequence of symmetric positive definite (SPD) Toeplitz matrices generated by an even, non-negative functions *f* with zeros of any positive order, that is,

$$[T_n(f)]_{kj} = [T_n(f)]_{k-j} = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-i(k-j)x} dx \qquad k, j = 1, 2, \dots, n, \quad i^2 = -1$$

and  $\{\tau_n(f)\}_n$  is the sequence of a specific  $\tau$  matrices, generated as

$$\tau_n(f) = \mathbb{S}_n \operatorname{diag}(f(\theta))\mathbb{S}_n, \qquad \theta = [\theta_1, \theta_2, \dots, \theta_n], \qquad \theta_j = \frac{j\pi}{n+1} = j\pi h, \qquad j = 1, \dots, n,$$

and

$$[\mathbb{S}_n]_{i,j} = \sqrt{\frac{2}{n+1}} \sin\left(i\theta_j\right), \qquad i,j=1,\ldots,n.$$
(1)

We recall here that  $S_n$  is symmetric and orthogonal and so it coincides with its inverse and that "essential spectral equivalence" means that all the eigenvalues of  $\{\tau_n^{-1}(f)T_n(f)\}_n$  belong to an interval [c, C] except possible *m* outliers, not converging to zero as the matrix size tends to infinity. In the case of generating functions with the order of their zero lying in the interval [0, 3] it is worth noticing that there are no outliers.

According to the analysis given in the aforementioned works, the coefficient matrix of the corresponding linear system depends on the diffusion coefficients of the FDE. In the simplest case where they are constant and equal, the related matrix is a diagonal matrix multiplied by a real SPD Toeplitz matrix with its generating function  $\mathcal{F}_{\alpha}$  being even, positive, and real, having a zero at zero of real positive order between one and two, plus a positive diagonal with constant entries that asymptotically tend to zero. Analysis shows that this matrix is present in the more general case where the diffusion coefficients are not constant and not equal to each other. In this case, a diagonal times skew-symmetric real Toeplitz matrix is then added to the coefficient matrix. Taking advantage of this fact, we propose the preconditioner  $\mathcal{P}_{F_{\alpha}} = D_n \tau_n(\mathcal{F}_{\alpha})$ , where  $D_n$  is a suitable diagonal matrix defined in Section 3. We show that this preconditioner can effectively keep the real part of the eigenvalues away from zero, while the sine transform keeps the cost per iteration  $\mathcal{O}(n \log n)$ , using a specific real algorithm or using the fast Fourier transform (FFT). It turns out that this preconditioner is very efficient and performs better, especially in multi-dimensional case, than the proposed preconditioners in References 3 and 20.

The article is organized as follows. In Sections 1.1–1.4, we present the one and two-dimensional FDE problems and the respective discretizations. Then, in Section 2 we summarize the spectral analysis performed in References 3,20, which turns out to be necessary for the definition of the new preconditioner. In Section 3, we also define the proposed preconditioners for the one and two-dimensional cases. In Section 4, we report numerical experiments and results that confirm the efficiency of the proposed preconditioner. Finally, in Section 5 we discuss the advantages and disadvantages of the proposed preconditioners and possible future research directions.

### **1.1** | Fractional diffusion equations

Consider the two dimensional initial-boundary value problem

$$\begin{cases} \frac{\partial u(x,y,t)}{\partial t} = d_{+}(x,y,t) \frac{\partial^{a} u(x,y,t)}{\partial_{+}x^{a}} + d_{-}(x,y,t) \frac{\partial^{a} u(x,y,t)}{\partial_{-}x^{a}} + \\ + e_{+}(x,y,t) \frac{\partial^{b} u(x,y,t)}{\partial_{+}y^{\beta}} + e_{-}(x,y,t) \frac{\partial^{b} u(x,y,t)}{\partial_{-}y^{\beta}} + f(x,y,t), \qquad (x,y,t) \in \Omega \times (0,T), \\ u(x,y,t) = 0, \qquad (x,y,t) \in \mathbb{R}^{2} \setminus \Omega \times [0,T], \\ u(x,y,0) = u_{0}(x,y), \qquad (x,y) \in \overline{\Omega}, \end{cases}$$
(2)

where  $\Omega = (L_1, R_1) \times (L_2, R_2), \alpha, \beta \in (1, 2)$  is the fractional derivative order, f(x, y, t) is the source term and the nonnegative functions  $d_{\pm}(x, y, t)$  and  $e_{\pm}(x, y, t)$  are the diffusion coefficients. Accordingly, in the one-dimensional setting we drop the dependency on *y*, while the terms including  $e_{\pm}(x, y, t)$  are not present.

The left-handed  $(\partial_+)$  and the right-handed  $(\partial_-)$  fractional derivatives in (2) are defined in Riemann–Liouville form as follows:

$$\frac{\partial^{\alpha}u(x,y,t)}{\partial_{+}x^{\alpha}} = \frac{1}{\Gamma(2-\alpha)}\frac{\partial^{2}}{\partial x^{2}}\int_{L_{1}}^{x}\frac{u(\xi,y,t)}{(x-\xi)^{\alpha-1}}d\xi, \qquad \qquad \frac{\partial^{\alpha}u(x,y,t)}{\partial_{-}x^{\alpha}} = \frac{1}{\Gamma(2-\alpha)}\frac{\partial^{2}}{\partial x^{2}}\int_{x}^{R_{1}}\frac{u(\xi,y,t)}{(\xi-x)^{\alpha-1}}d\xi, \\ \frac{\partial^{\beta}u(x,y,t)}{\partial_{+}y^{\beta}} = \frac{1}{\Gamma(2-\beta)}\frac{\partial^{2}}{\partial y^{2}}\int_{L_{2}}^{y}\frac{u(x,\eta,t)}{(y-\eta)^{\beta-1}}d\eta, \qquad \qquad \frac{\partial^{\beta}u(x,y,t)}{\partial_{-}y^{\beta}} = \frac{1}{\Gamma(2-\beta)}\frac{\partial^{2}}{\partial y^{2}}\int_{y}^{R_{2}}\frac{u(x,\eta,t)}{(\eta-y)^{\beta-1}}d\eta.$$

### **1.2** | First-order finite difference discretization

In this section, we consider the one-dimensional version of (2) (for two-dimensional derivation see Section 1.4 and Reference 20). Applying the shifted Grünwald formulas we can approximate the left and right fractional derivatives by

$$\begin{split} \frac{\partial^{\alpha} u(x,t)}{\partial_{+}x^{\alpha}} &= \frac{1}{h_{x}^{\alpha}} \sum_{k=0}^{\lfloor (x-L_{1})/h_{x} \rfloor} g_{k}^{(\alpha)} u(x-(k-1)h_{x},t) + \mathcal{O}(h_{x}), \\ \frac{\partial^{\alpha} u(x,t)}{\partial_{-}x^{\alpha}} &= \frac{1}{h_{x}^{\alpha}} \sum_{k=0}^{\lfloor (R_{1}-x)/h_{x} \rfloor} g_{k}^{(\alpha)} u(x+(k-1)h_{x},t) + \mathcal{O}(h_{x}), \end{split}$$

4 of 22 WILEY

where  $\lfloor \cdot \rfloor$  is the floor function,  $n_1$  is the discretization parameter giving  $h_x = (R_1 - L_1)/(n_1 + 1) = (R_1 - L_1)h_1$ , and  $g_k^{(\alpha)}$  are the alternating fractional binomial coefficients defined as

$$g_k^{(\alpha)} = (-1)^k \binom{\alpha}{k} = \frac{(-1)^k}{k!} \alpha(\alpha - 1) \cdots (\alpha - k + 1), \quad k = 0, 1, \dots,$$
(3)

where  $\binom{\alpha}{0} = 1$ . Using the implicit Euler method for time discretization, we define the number of time steps (index *m*) to be *M*, and thus  $h_t = T/M$ , and

$$\frac{u_i^{(m)} - u_i^{(m-1)}}{h_t} = \frac{d_{+,i}^{(m)}}{h_x^{\alpha}} \sum_{k=0}^{i+1} g_k^{(\alpha)} u_{i-k+1}^{(m)} + \frac{d_{-,i}^{(m)}}{h_x^{\alpha}} \sum_{k=0}^{n_i-i+2} g_k^{(\alpha)} u_{i+k-1}^{(m)} + f_i^{(m)},$$

where  $d_{\pm,i}^{(m)} = d_{\pm}(x_i, t_m)$ ,  $u_i^{(m)} = u(x_i, t_m)$ , and  $f_i^{(m)} = f(x_i, t_m)$ , where  $x_i = L_1 + ih_x$  and  $t_m = mh_t$ . After rearranging terms, we find

$$\frac{h_x^{\alpha}}{h_t}u_i^{(m)} - d_{+,i}^{(m)}\sum_{k=0}^{i+1}g_k^{(\alpha)}u_{i-k+1}^{(m)} - d_{-,i}^{(m)}\sum_{k=0}^{n_1-i+2}g_k^{(\alpha)}u_{i+k-1}^{(m)} = \frac{h_x^{\alpha}}{h_t}u_i^{(m-1)} + h_x^{\alpha}f_i^{(m)}$$

or in matrix form, the linear systems

$$\left(\nu_{M,n_{1}}\mathbb{I}_{n_{1}}+D_{+}^{(m)}T_{\alpha,n_{1}}+D_{-}^{(m)}T_{\alpha,n_{1}}^{\mathrm{T}}\right)\mathbf{u}^{(m)}=\nu_{M,n_{1}}\mathbf{u}^{(m-1)}+h_{x}^{\alpha}\mathbf{f}^{(m)},\tag{4}$$

where

 $\mathbb{I}_{n_1}: \text{The identity matrix of size } n_1, \tag{5}$ 

$$v_{M,n_{1}} = \frac{h_{x}^{\alpha}}{h_{t}},$$

$$\mathbf{u}^{(m)} = \begin{bmatrix} u_{1}^{(m)}, u_{2}^{(m)}, \dots, u_{n_{1}}^{(m)} \end{bmatrix}^{\mathrm{T}},$$

$$\mathbf{f}^{(m)} = \begin{bmatrix} f_{1}^{(m)}, f_{2}^{(m)}, \dots, f_{n_{1}}^{(m)} \end{bmatrix}^{\mathrm{T}},$$

$$[D_{\pm}^{(m)}]_{i,i} = d_{\pm}^{(m)}(x_{i}, t_{m}), \quad i = 1, \dots, n_{1},$$
(6)

and

with the coefficients  $g_k^{(\alpha)}$  given in (3).

Now define

$$\mathcal{M}_{\alpha,n_{1}}^{(m)} = \left( \nu_{M,n_{1}} \mathbb{I}_{n_{1}} + D_{+}^{(m)} T_{\alpha,n_{1}} + D_{-}^{(m)} T_{\alpha,n_{1}}^{\mathrm{T}} \right),$$
  
$$\mathbf{b}^{(m)} = \mathbf{v}_{M,n_{1}} u^{(m-1)} + h_{x}^{\alpha} \mathbf{f}^{(m)}.$$
 (8)

Then, for each time step *m*, we solve the system

$$\mathcal{M}_{\alpha,n_1}^{(m)} \mathbf{u}^{(m)} = \mathbf{b}^{(m)}.$$
(9)

#### **1.3** | Second-order finite difference discretization

For the second order finite difference discretization in space, we can just exchange the matrix  $T_{\alpha,n_1}$  in (4) with a matrix  $S_{\alpha,n_1}$  defined by

$$S_{\alpha,n_{1}} = - \begin{bmatrix} w_{1}^{(\alpha)} & w_{0}^{(\alpha)} & & & \\ w_{2}^{(\alpha)} & w_{1}^{(\alpha)} & w_{0}^{(\alpha)} & & \\ w_{3}^{(\alpha)} & w_{2}^{(\alpha)} & w_{1}^{(\alpha)} & w_{0}^{(\alpha)} & \\ \vdots & \ddots & \ddots & \ddots & \ddots & \\ \vdots & \ddots & \ddots & \ddots & \ddots & \\ w_{n_{1}-1}^{(\alpha)} & w_{n_{1}-2}^{(\alpha)} & \cdots & \ddots & w_{2}^{(\alpha)} & w_{1}^{(\alpha)} \\ w_{n_{1}}^{(\alpha)} & w_{n_{1}-1}^{(\alpha)} & \cdots & \cdots & w_{3}^{(\alpha)} & w_{2}^{(\alpha)} & w_{1}^{(\alpha)} \end{bmatrix},$$
(10)

where

$$\begin{split} & w_0^{(\alpha)} = \frac{\alpha}{2} g_0^{(\alpha)}, \\ & w_k^{(\alpha)} = \frac{\alpha}{2} g_k^{(\alpha)} + \frac{2 - \alpha}{2} g_{k-1}^{(\alpha)}, \quad k \ge 1, \end{split}$$

and the coefficients  $g_k^{(\alpha)}$  are expressed as in relation (3).

#### 1.4 | Two-dimensional case

Similarly to 1D case, we can extend the discretization scheme to the two-dimensional setting. In the next paragraph we summarize the main points of the numerical procedure, referring the reader in Reference 20 for further details. Define

$$h_x = \frac{R_1 - L_1}{n_1 + 1} = (R_1 - L_1)h_1, \qquad x_i = L_1 + ih_x, \quad i = 1, \dots, n_1,$$
  
$$h_y = \frac{R_2 - L_2}{n_2 + 1} = (R_2 - L_2)h_2, \qquad y_i = L_2 + ih_y, \quad i = 1, \dots, n_2,$$

and  $N = n_1 n_2$ . The solution u(x, y, t) is discretized as  $u_{i,j}^{(m)} = u(x_i, y_j, t^{(m)})$ ,

$$\mathbf{u}^{(m)} = [u_{1,1}^{(m)}, \ldots, u_{n_1,1}^{(m)}, u_{1,2}^{(m)}, \ldots, u_{n_1,2}^{(m)}, \ldots, u_{1,n_2}^{(m)}, \ldots, u_{n_1,n_2}^{(m)}]^{\mathrm{T}},$$

and the four diffusion function  $d_+(x, y, t)$ ,  $d_-(x, y, t)$ ,  $e_+(x, y, t)$ ,  $e_-(x, y, t)$  are discretized as  $d_{i,j}^{\pm,(m)} = d_{\pm}(x_i, y_j, t^{(m)})$  and  $e_{i,j}^{\pm,(m)} = e_{\pm}(x_i, y_j, t^{(m)})$ ,

$$\mathbf{d}_{\pm}^{(m)} = \left[ d_{1,1}^{\pm,(m)}, \dots, d_{n_{1},1}^{\pm,(m)}, d_{1,2}^{\pm,(m)}, \dots, d_{n_{1},2}^{\pm,(m)}, \dots, d_{1,n_{2}}^{\pm,(m)}, \dots, d_{n_{1},n_{2}}^{\pm,(m)} \right]^{\mathrm{I}}, \\ \mathbf{e}_{\pm}^{(m)} = \left[ e_{1,1}^{\pm,(m)}, \dots, e_{n_{1},1}^{\pm,(m)}, e_{1,2}^{\pm,(m)}, \dots, e_{n_{1},2}^{\pm,(m)}, \dots, e_{1,n_{2}}^{\pm,(m)}, \dots, e_{n_{1},n_{2}}^{\pm,(m)} \right]^{\mathrm{T}}.$$

The source term f(x, y, t) is discretized as  $f_{i,j}^{(m)} = f(x_i, y_j, t^{(m)})$ ,

$$\mathbf{v}^{(m-1/2)} = \left[ f_{1,1}^{(m-1/2)}, \dots, f_{n_1,1}^{(m-1/2)}, f_{1,2}^{(m-1/2)}, \dots, f_{n_1,2}^{(m-1/2)}, \dots, f_{1,n_2}^{(m-1/2)}, \dots, f_{n_1,n_2}^{(m-1/2)} \right]^{\mathrm{T}}.$$

We also define the four matrices  $D_{\pm}^{(m)} = \text{diag}(\mathbf{d}_{\pm}^{(m)})$  and  $E_{\pm}^{(m)} = \text{diag}(\mathbf{e}_{\pm}^{(m)})$ .

WILEY <u>5 of 22</u>

### 6 of 22 WILEY

If we have two fractional derivatives,  $\alpha$  and  $\beta$ , in each spatial direction we define the two matrices  $S_{\alpha,n_1}$  and  $S_{\beta,n_2}$  (or  $T_{\alpha,n_1}$  and  $T_{\beta,n_2}$  for the considered first-order discretization).

We also define the two  $N\times N$  matrices

$$\begin{split} A_x^{(m)} &= D_+^{(m)}(\mathbb{I}_{n_2} \otimes S_{\alpha,n_1}) + D_-^{(m)}(\mathbb{I}_{n_2} \otimes S_{\alpha,n_1}^{\mathrm{T}}), \\ A_y^{(m)} &= E_+^{(m)}(S_{\beta,n_2} \otimes \mathbb{I}_{n_1}) + E_-^{(m)}(S_{\beta,n_2}^{\mathrm{T}} \otimes \mathbb{I}_{n_1}), \end{split}$$

where  $\mathbb{I}_n$  denotes the identity matrix of size *n*, and  $\otimes$  is the Kronecker product. Using Crank–Nicolson approach for time discretization (e.g., see Reference 20) we obtain the system

$$\left(\frac{1}{r}\mathbb{I}_{N}+A_{x}^{(m)}+\frac{s}{r}A_{y}^{(m)}\right)\mathbf{u}^{(m)}=\left(\frac{1}{r}\mathbb{I}_{N}-A_{x}^{(m-1)}-\frac{s}{r}A_{y}^{(m-1)}\right)\mathbf{u}^{(m-1)}+2h_{x}^{\alpha}\mathbf{v}^{(m-1/2)},$$

where  $r = \frac{h_t}{2h_x^{\alpha}}$ ,  $s = \frac{h_t}{2h_y^{\theta}}$ . In compact form we have

$$\mathcal{M}_{(\alpha,\beta),N}^{(m)}\mathbf{u}^{(m)}=\mathbf{b}^{(m)},$$

where

$$\mathcal{M}_{(\alpha,\beta),N}^{(m)} = \frac{1}{r} \mathbb{I}_N + A_x^{(m)} + \frac{s}{r} A_y^{(m)},$$
  
$$\mathbf{b}^{(m)} = \left(\frac{1}{r} \mathbb{I}_N - A_x^{(m-1)} - \frac{s}{r} A_y^{(m-1)}\right) \mathbf{u}^{(m-1)} + 2h_x^{\alpha} \mathbf{v}^{(m-1/2)}.$$

#### 2 | SPECTRAL ANALYSIS

In this section we provide some definitions that are used in the analysis. We also employ the theory of GLT matrix sequences to study the spectral properties of  $\mathcal{M}_{\alpha,n_1}^{(m)}$  of (9) (for both the first and second order version) as the matrix dimension tends to infinity. We refer the reader to Reference 22 for an introduction to the theory of GLT matrix sequences. Here, we only list some basic properties that are used in the analysis that follows. The results reported in Sections 2.1 and 2.2 are taken from References 3,20.

**Definition 1.** Let  $\{A_n\}_n$  be a matrix sequence and  $f : D \to \mathbb{C}$  be a measurable function defined on a measurable set  $D \subset \mathbb{R}^k$  with  $0 < \mu(D) < \infty$ .

• We say that the sequence  $\{A_n\}_n$  has an asymptotic singular value distribution described by f, and we write  $\{A_n\}_n \sim_{\sigma} f$  if,

$$\lim_{n\to\infty}\frac{1}{n}\sum_{j=1}^{n}F(\sigma_{j}(A_{n}))=\frac{1}{\mu(D)}\int_{D}F(|f(x)|)dx,\quad\forall F\in C_{c}(\mathbb{R}),$$

where  $C_c(\mathbb{R})$  is the set of continuous functions with compact support over  $\mathbb{R}$ .

• We say that  $\{A_n\}_n$  has an asymptotic eigenvalue distribution described by f, and write  $\{A_n\}_n \sim_{\lambda} f$  if

$$\lim_{n\to\infty}\frac{1}{n}\sum_{j=1}^{n}F(\lambda_{j}(A_{n}))=\frac{1}{\mu(D)}\int_{D}F(f(x))dx,\quad\forall F\in C_{c}(\mathbb{C}),$$

where  $C_c(\mathbb{R})$  is the set of continuous functions with compact support over  $\mathbb{R}$ .

**Definition 2.** Let  $f \in L^1([-\pi, \pi])$  and  $\{f_k\}_{k \in \mathbb{Z}}$  its Fourier coefficients defined as

$$f_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(\theta) e^{-\mathbf{i}k\theta} \ d\theta, \quad k = 0, \pm 1, \pm 2, \dots$$

The sequence of matrices  $\{T_n(f)\}_{n \in \mathbb{N}}, T_n(f) = [f_{i-j}]_{i,j=1}^n$ , is called a Toeplitz sequence generated by f.

The eigenvalue and singular value distribution of Toeplitz sequences generated by  $f \in L^1([-\pi, \pi])$  is given by generalized Szegő theorem:<sup>23</sup>

**Theorem 1.** Let  $f \in L^1([-\pi, \pi])$  and  $T_n(f)$  be the Toeplitz matrix generated by f. Then f is the spectral symbol of the sequence, that is

$$\{T_n(f)\}_n \sim_\sigma f$$

If, moreover, f is real-valued, then

$$\{T_n(f)\}_n \sim_{\lambda} f.$$

The basic properties of the GLT class follow.

- GLT1 Each GLT sequence  $\{A_n\}_n$  has a singular value symbol  $\tilde{f} : [0, 1] \times [-\pi, \pi] \to \mathbb{C}$ . If all the matrices of the sequence are Hermitian, then the distribution also holds in the eigenvalue sense. We call  $\tilde{f}(x, \theta)$  the GLT symbol of  $\{A_n\}_n$  and we write  $\{A_n\}_n \sim_{\text{GLT}} \tilde{f}$ .
- GLT2 The set of GLT sequences is closed under linear combinations, products, inversion (whenever the symbol is zero in at most a set of zero Lebesgue measure) and conjugation. The sequence obtained via algebraic operations on a finite set of given GLT sequences is still a GLT sequence and its symbol is obtained by performing the same algebraic manipulations on the corresponding symbols of the input GLT sequences.
- GLT3 Every Toeplitz sequence generated by a function  $f \in L^1([-\pi, \pi])$  is a GLT sequence and its symbol is  $\tilde{f}(x, \theta) = f(\theta)$ . If  $\alpha : [0, 1] \to \mathbb{C}$  is a Riemann integrable function, the diagonal matrix sequence of the form  $\{D_n(\alpha)\}_n, n \in \mathbb{N}, D_n(\alpha) = \text{diag}_{j=1,...,n}(\alpha(\frac{j}{n}))$  is a GLT sequence with spectral symbol  $\tilde{f}(x, \theta) = \alpha(x)$ .

### 2.1 | Spectral analysis: Matrices $T_{\alpha,n}$ and $S_{\alpha,n}$

From Reference 3 we know that  $T_{\alpha,n_1}$  in (7) is a Toeplitz sequence with spectral symbol

$$g_{\alpha}(\theta) = -e^{-\mathbf{i}\theta} \left(1 - e^{\mathbf{i}\theta}\right)^{\alpha},$$

and thus from Theorem 1 and GLT3

$$\{T_{\alpha,n}\}_n = \{T_n(g_\alpha)\}_n \sim_{\text{GLT},\sigma} g_\alpha.$$
<sup>(11)</sup>

Furthermore, as shown in Reference 20,  $S_{a,n}$  in (10) is a Toeplitz sequence with spectral symbol

$$w_{\alpha}(\theta) = -\left(\frac{2-\alpha(1-e^{-\mathbf{i}\theta})}{2}\right)\left(1-e^{\mathbf{i}\theta}\right)^{\alpha},\tag{12}$$

and thus from Theorem 1 and GLT3

$$\{S_{\alpha,n}\}_n = \{T_n(w_\alpha)\}_n \sim_{\text{GLT},\sigma} w_\alpha.$$
<sup>(13)</sup>

#### 2.2 | Spectral analysis: Constant coefficient case

**Theorem 2.** Assuming  $d_{\pm}(x,t) = d > 0$  and that  $v_{M,n} = o(1)$ , we have for the first order spatial discretization

$$\left\{\mathcal{M}_{\alpha,n}^{(m)}\right\}_n\sim_{\mathrm{GLT},\sigma,\lambda}d\cdot p_\alpha(\theta),$$

where

$$p_{\alpha}(\theta) = g_{\alpha}(\theta) + g_{\alpha}(-\theta).$$
(14)

For the second order spatial discretization we have

$$\left\{\mathcal{M}_{\alpha,n}^{(m)}\right\}_n \sim_{\mathrm{GLT},\sigma,\lambda} d \cdot q_\alpha(\theta),$$

where

$$q_{\alpha}(\theta) = w_{\alpha}(\theta) + w_{\alpha}(-\theta). \tag{15}$$

In the subsequent analysis, whenever either symbol  $p_{\alpha}(\theta)$  or  $q_{\alpha}(\theta)$  is applicable, we denote both symbols by  $\mathcal{F}_{\alpha}(\theta)$ .

**Proposition 1.** Let  $\alpha \in (1, 2)$ , then the function  $p_{\alpha}(\theta)$  has a zero of order  $\alpha$  at 0.

Moreover, in connection with Proposition 1, it is worth noticing the following: if *f* is nonnegative with a unique zero of order  $\alpha > 0$ , then the matrix  $T_n(f)$  is positive definite for any *n* its minimal eigenvalue tends to zero as *n* tends to infinity as  $n^{-\alpha}$ ; furthermore, if *f* is also bounded then the condition number of  $T_n(f)$  grows asymptotically as  $n^{\alpha}$  (e.g., see References 5,6).

### **3** | MAIN RESULTS

In this section we propose two new preconditioners, based on the spectral symbol, for the one and two-dimensional problems.

#### 3.1 | Proposed preconditioner: One dimension

To be consistent with Reference 3, so that results can be compared, we use the first-order spatial discretization for the one dimensional case. We also omit the time dependency mark to simplify the notation. Thus, let  $T_n = T_{\alpha,n1}$  be defined as in (7) and let  $\mathcal{M}_n = \mathcal{M}_{\alpha,n_1}$  be defined as in (8).

As previously mentioned in Section 1, the proposed preconditioner is similar to a diagonal matrix  $D_n$  times a specific  $\tau$  matrix, that is,  $\mathcal{P}_{\mathcal{F}_\alpha} = D_n \tau_n(\mathcal{F}_\alpha(\theta))$ , where both these parts will be clarified through this paragraph.

The product of two or more matrices as preconditioner is not a new proposal (see, e.g., Reference 24). The coefficient matrix of the system  $\mathcal{M}_n = v_{\mathcal{M},n} \mathbb{I}_n + D_+ T_n + D_- T_n^{\mathrm{T}}$  suggests the following candidate for the diagonal matrix

$$D_n = \frac{1}{2} (D_+ + D_-),$$
  

$$[D_n]_{i,i} = \frac{d_{+,i} + d_{-,i}}{2},$$
(16)

that has been used in other preconditioning strategies (see e.g., Reference 3). Then, assuming that  $d_{\pm}$  do not have a common zero at  $x_0 \in [L, R]$ , we deduce that  $D_n^{-1}$  is uniformly bounded and

$$D_n^{-1}\mathcal{M}_n = v_{M,n}D_n^{-1} + D_n^{-1}D_+T_n + D_n^{-1}D_-T_n^{\mathrm{T}}.$$

Defining  $\delta(x) = \frac{d_{\pm}(x)}{d_{\pm}(x)+d_{\pm}(x)}$ ,  $\delta_i = \delta(x_i)$ ,  $\delta = [\delta_1, \delta_2, \dots, \delta_n]$ , and  $G_n = \text{diag}(\delta)$ , taking into account that  $d_{\pm}$  are non negative functions, we have that  $0 \le \delta(x) \le 1$  and also

$$D_n^{-1}D_+ = 2G_n,$$
  
 $D_n^{-1}D_- = 2(\mathbb{I}_n - G_n).$ 

Hence,  $D_n^{-1}\mathcal{M}_n$  can be written as

$$D_n^{-1}\mathcal{M}_n = v_{M,n}D_n^{-1} + D_n^{-1}D_+T_n + D_n^{-1}D_-T_n^{\mathrm{T}}$$
  
=  $v_{M,n}D_n^{-1} + 2G_nT_n + 2(\mathbb{I}_n - G_n)T_n^{\mathrm{T}}$   
=  $v_{M,n}D_n^{-1} + (T_n + T_n^{\mathrm{T}}) + (2G_n - \mathbb{I}_n)(T_n - T_n^{\mathrm{T}}).$ 

WILEY 9 of 22

Since, from (11), 
$$T_n := T_n (-e^{-i\theta} (1 - e^{i\theta})^a) = T_n (g_\alpha(\theta)) \text{ and } T_n^T := T_n (-e^{i\theta} (1 - e^{-i\theta})^a) = T_n (g_\alpha(-\theta)) \text{ we have}$$
  

$$D_n^{-1} \mathcal{M}_n = v_{M,n} D_n^{-1} + (T_n + T_n^T) + (2G_n - \mathbb{I}_n)(T_n - T_n^T)$$

$$= v_{M,n} D_n^{-1} + T_n (g_\alpha(\theta) + g_\alpha(-\theta)) + (2G_n - \mathbb{I}_n) T_n (g_\alpha(\theta) - g_\alpha(-\theta))$$

$$= v_{M,n} D_n^{-1} + T_n (p_\alpha(\theta)) + (2G_n - \mathbb{I}_n) T_n (2\mathbf{i}\mathfrak{F} \{g_\alpha(\theta)\}), \qquad (17)$$

where  $p_{\alpha}(\theta)$ , defined in (14), is real, positive and even. With  $\mathfrak{F}$  we denote the imaginary part of a function. The above derivation of the  $D_n^{-1}\mathcal{M}_n$  matrix is of interest since it makes clear why it is reasonable to use the  $\tau$  preconditioner. The first term of the above matrix,  $v_{M,n}D_n^{-1}$ , is diagonal with positive and o(1) entries, since we have supposed that the  $d_{\pm}$  functions do not have zeros at the same point in the domain [L, R] and  $v_{M,n} = o(1)$ . We mention here that although the entries of this term are o(1), its effect on the eigenvalues of the preconditioned matrix can be significant. The reason is explained in the end of this section. The third term in (17) is a diagonal matrix with entries in [-1, 1] times a skew-symmetric Toeplitz matrix with generating function  $2\mathbf{i}\mathfrak{F}\{g_{\alpha}(\theta)\}$ . If  $d_+ = d_-$  this term is vanishing while if the  $d_{\pm}$  are constant but not equal it is a pure skew-symmetric Toeplitz (in that case  $(2G_n - \mathbb{I}_n) = c\mathbb{I}_n$  for some constant c).

The term in (17) which is mainly responsible for the dispersion of the real part of the spectrum, is the second term, that is,  $T_n(p_\alpha(\theta))$ . The  $\tau$  preconditioner will effectively cluster the eigenvalues of this matrix, and consequently the eigenvalues of the whole matrix  $D_n^{-1}\mathcal{M}_n$ . Hence, taking advantage of the essential spectral equivalence between the matrix sequences  $\{\tau_n(f)\}_n$  and  $\{T_n(f)\}_n$  proven in Reference 21, we propose a preconditioner expressed as

$$\mathcal{P}_{\mathcal{F}_{\alpha},n} = D_n \tau_n(p_{\alpha}(\theta)) = D_n \mathbb{S}_n F_n \mathbb{S}_n, \tag{18}$$

where

$$F_n = \operatorname{diag}(p_\alpha(\theta)), \qquad \theta = [\theta_1, \theta_2, \dots, \theta_n], \qquad \theta_j = \frac{j\pi}{n+1} = j\pi h, \qquad j = 1, \dots, n$$

with  $D_n$  defined in (16) and  $\mathbb{S}_n$  being the sine transform matrix reported in (1).

#### 3.1.1 | Case I: $d_{\pm}$ are constants

In the case where the diffusion coefficient functions are constants, the (17) becomes:

$$\left(2\frac{\nu_{M,n}}{d_++d_-}\right)\mathbb{I}_n+T_n\left(p_\alpha(\theta)\right)+\left(\frac{d_+-d_-}{d_++d_-}\right)T_n\left(2\mathbf{i}\Im\left\{g_\alpha(\theta)\right\}\right)=T_n\left(2\frac{\nu_{M,n}}{d_++d_-}+p_\alpha(\theta)\right)+T_n\left(2\left(\frac{d_+-d_-}{d_++d_-}\right)\mathbf{i}\Im\left\{g_\alpha(\theta)\right\}\right),$$

that is, is exactly the sum of a symmetric and a skew-symmetric Toeplitz matrix. It is worth noticing that according to the GLT machinery, the term  $\frac{2 \cdot v_{M,n}}{d_++d_-}$  which is added to the symbol of the first Toeplitz matrix sequence does not change the symbol of the sequence since is of order o(1). However it affects the speed in which the minimum eigenvalue of the sequence approaches zero as the dimension of the matrix tends to infinity. Thus, in this special case, the  $\tau$  part of preconditioner is defined as

$$\tau_{M,n}\left(p_{\alpha}(\theta) + \frac{2 \cdot v_{M,n}}{d_{+} + d_{-}}\right) = \mathbb{S}_{n} \operatorname{diag}\left(p_{\alpha}(\theta) + \frac{2 \cdot v_{M,n}}{d_{+} + d_{-}}\right) \mathbb{S}_{n} = \mathbb{S}_{n} \hat{F}_{n} \mathbb{S}_{n}.$$

Then,

$$\begin{aligned} \tau_{M,n}^{-1} \left( p_{\alpha}(\theta) + \frac{2 \cdot v_{M,n}}{d_{+} + d_{-}} \right) \left[ T_{n} \left( \frac{2 \cdot v_{M,n}}{d_{+} + d_{-}} + p_{\alpha}(\theta) \right) + T_{n} \left( 2 \frac{d_{+} - d_{-}}{d_{+} + d_{-}} \mathbf{i} \mathfrak{F} \left\{ g_{\alpha}(\theta) \right\} \right) \right] \\ &\sim \hat{F}_{n}^{-\frac{1}{2}} \mathbb{S}_{n} \left[ T_{n} \left( \frac{2 \cdot v_{M,n}}{d_{+} + d_{-}} + p_{\alpha}(\theta) \right) + T_{n} \left( 2 \frac{d_{+} - d_{-}}{d_{+} + d_{-}} \mathbf{i} \mathfrak{F} \left\{ g_{\alpha}(\theta) \right\} \right) \right] \mathbb{S}_{n} \hat{F}_{n}^{-\frac{1}{2}} \\ &= \hat{F}_{n}^{-\frac{1}{2}} \mathbb{S}_{n} T_{n} \left( \frac{2 \cdot v_{M,n}}{d_{+} + d_{-}} + p_{\alpha}(\theta) \right) \mathbb{S}_{n} \hat{F}_{n}^{-\frac{1}{2}} + \hat{F}_{n}^{-\frac{1}{2}} \mathbb{S}_{n} T_{n} \left( 2 \frac{d_{+} - d_{-}}{d_{+} + d_{-}} \mathbf{i} \mathfrak{F} \left\{ g_{\alpha}(\theta) \right\} \right) \mathbb{S}_{n} \hat{F}_{n}^{-\frac{1}{2}} \end{aligned}$$

The first term in the above sum is symmetric and its eigenvalues are strongly clustered at 1 since the conditions of the main theoretical result of Reference 21 are fulfilled concerning the spectral equivalence between a  $\tau$  matrix and a Toeplitz one.

### 10 of 22 | WILFY

The second term is skew-symmetric and it does not affect the real part of the eigenvalues of the whole matrix. Moreover, it is absent whenever  $d_+ = d_-$ . Hence, the real parts of the eigenvalues of the preconditioned matrix are strongly clustered at 1 and are bounded by constants c, C with  $0 < c \le 1 \le C < \infty$ .

3.1.2 | Case II: 
$$d_{-}(x) = d_{+}(x) > 0$$

In this case, the term  $2G_n - \mathbb{I}_n = \mathbf{0}$  in 17 is equal to zero and the preconditioned matrix becomes  $\tau_n^{-1}(p_\alpha(\theta))(v_{M,n}D_n^{-1} + T_n(p_\alpha(\theta)))$  which is similar to the SPD

$$\tau_n^{-1}(p_{\alpha}(\theta))(\nu_{M,n}D_n^{-1} + T_n(p_{\alpha}(\theta))) \sim F_n^{-1/2} \mathbb{S}_n(\nu_{M,n}D_n^{-1} + T_n(p_{\alpha}(\theta))) \mathbb{S}_n F_n^{-1/2}$$
  
=  $\nu_{M,n}F_n^{-1/2} \mathbb{S}_n(D_n^{-1}) \mathbb{S}_n F_n^{-1/2} + F_n^{-1/2} \mathbb{S}_n(T_n(p_{\alpha}(\theta))) \mathbb{S}_n F_n^{-1/2}.$  (19)

In the above splitting in positive symmetric terms, the first one has o(n) eigenvalues tending to infinity while the second one fulfills the main theoretical result of Reference 21 and thus, for every n, it has eigenvalues belonging to an interval [c,C] with c, C constants and  $0 < c \le 1 \le C < \infty$ . The claim about the spectrum of the first term can be proved if we equivalently show that the inverse of it, that is,  $F_n(\mathbb{S}_n D_n \mathbb{S}_n)$  has at most o(n) eigenvalues tending to 0 as  $n \to \infty$ . Since  $F_n$  is the diagonal matrix formed by the values  $p_\alpha(j\pi h)$ , j = 1, ..., n, which has a unique zero at zero of order  $\alpha$ , there will be an index  $\hat{j}$  with  $\hat{j}$  of order o(n) such that  $p_\alpha(j\pi h)$  being of order o(1) for all  $j \le \hat{j}$ . Thus, at most o(n) eigenvalues of  $F_n$  can tend to zero. Using Rayleigh quotient and taking into account that the matrix  $D_n$  is a diagonal matrix with entries bounded from above end below by positive universal constants, our claim is proved. Consequently, using the Weyl's theorem on (19) we obtain that

$$\lambda_k \left( \nu_{M,n} F_n^{-1} \mathbb{S}_n (D_n^{-1}) \mathbb{S}_n + F_n^{-1} \mathbb{S}_n (T_n(p_\alpha(\theta))) \mathbb{S}_n \right) \leq \nu_{M,n} \lambda_k (F_n^{-1} \mathbb{S}_n (D_n^{-1}) \mathbb{S}_n) + \lambda_n \left( F_n^{-1} \mathbb{S}_n (T_n(p_\alpha(\theta))) \mathbb{S}_n \right)$$

Accordingly, at most o(n) eigenvalues of  $\tau_n^{-1}(p_\alpha(\theta))(v_{M,n}D_n^{-1} + T_n(p_\alpha(\theta)))$  can tend to infinity. Clearly the term  $v_{M,n}$  which in general tends to zero as  $O(n^{1-\alpha})$ , can further reduce the number of eigenvalues tending to infinity.

We remark that as in the semi elliptic case (see Reference 25 and especially the numerical experiments therein), if the equal functions  $d_{\pm}$  have a root then we expect an unpredictable asymptotical behavior of the eigenvalues of coefficient matrix  $\mathcal{M}_{\alpha}$ .

### 3.1.3 | Case III: General case

In the case where  $d_+ \neq d_-$  the term  $(2G_n - \mathbb{I}_n)(T_n - T_n^T)$  is nonzero and it affects the spectrum of the preconditioned matrix. Specifically,

$$\begin{aligned} \tau_n^{-1}(p_{\alpha}(\theta))(v_{M,n}D_n^{-1} + T_n(p_{\alpha}(\theta)) + (2G_n - \mathbb{I}_n)T_n(2\mathbf{i}\Im \{g_{\alpha}(\theta)\})) \\ &\sim F_n^{-1/2} \mathbb{S}_n(v_{M,n}D_n^{-1} + T_n(p_{\alpha}(\theta)) + (2G_n - \mathbb{I}_n)T_n(2\mathbf{i}\Im \{g_{\alpha}(\theta)\})) \mathbb{S}_n F_n^{-1/2} \\ &= F_n^{-1/2} \mathbb{S}_n(v_{M,n}D_n^{-1}) \mathbb{S}_n F_n^{-1/2} + F_n^{-1/2} \mathbb{S}_n(T_n(p_{\alpha}(\theta))) \mathbb{S}_n F_n^{-1/2} + F_n^{-1/2} \mathbb{S}_n(2G_n - \mathbb{I}_n)T_n(2\mathbf{i}\Im \{g_{\alpha}(\theta)\}) \mathbb{S}_n F_n^{-1/2}, \end{aligned}$$

where only the, new, third term can add imaginary quantity on the eigenvalues. However, we have observed through experimentation, that the effect of this third term on the real part of the eigenvalues is negligible. In this sense, we chose all the numerical experiments given in Section 4 belong to this case mainly for showing the performance of our proposal there were our spectral analysis do not explicitly and in depth cover the topic.

### 3.2 | Proposed preconditioner: Two dimensions

In the two-dimensional case we use the second order spatial discretization, in order to be consistent with Reference 20 and be able to readily compare the results. In this case, as reported in Section 1.4, the coefficient matrix of the system is defined as

$$\mathcal{M}_{(\alpha,\beta),N}^{(m)} = \frac{1}{r} \mathbb{I}_{N} + D_{+}^{(m)} (\mathbb{I}_{n_{2}} \otimes S_{\alpha,n_{1}}) + D_{-}^{(m)} (\mathbb{I}_{n_{2}} \otimes S_{\alpha,n_{1}}^{\mathrm{T}}) + \frac{s}{r} \left( E_{+}^{(m)} (S_{\beta,n_{2}} \otimes \mathbb{I}_{n_{1}}) + E_{-}^{(m)} (S_{\beta,n_{2}}^{\mathrm{T}} \otimes \mathbb{I}_{n_{1}}) \right).$$
(20)

VII FV <u>11 of 22</u>

We recall that  $S_{\alpha,n_1} = T_{n_1}(w_{\alpha}(\theta))$  and  $S_{\beta,n_2} = T_{n_2}(w_{\beta}(\theta))$  (see 12 and 13). Again, for simplicity we here omit the time dependency in the notation.

Now let  $\mathcal{F}_{(\alpha,\beta)}(\theta_1, \theta_2) = q_{\alpha}(\theta_1) + \frac{s}{r}q_{\beta}(\theta_2)$  where *q* is the real, nonnegative and even function defined in (15),  $\theta_1, \theta_2 \in [-\pi, \pi]$ , and  $n_1, n_2$  the two integers used for the discretization of the domain  $[L_x, R_x] \times [L_y, R_y]$ . Using the grid in (1) we define the diagonal matrices

$$F_{n_1,j} = \operatorname{diag}(\mathcal{F}_{(\alpha,\beta)}(\theta_{i,n_1},\theta_{j,n_2}), i = 1, \dots, n_1),$$

for each  $j = 1, ..., n_2$ . Then, the  $N \times N$  diagonal matrix is expressed as

$$F_{N} = \begin{bmatrix} F_{n_{1},1} & & & \\ & F_{n_{1},2} & & \\ & & \ddots & \\ & & & \ddots & \\ & & & & F_{n_{1},n_{2}} \end{bmatrix}.$$
 (21)

Let  $\mathbb{S}_{n_1}$  and  $\mathbb{S}_{n_2}$  be the discrete sine transform matrices of sizes  $n_1$  and  $n_2$ , respectively, as they defined in (1). Then, generalizing the idea of (18), our proposed preconditioner for this case is

$$\mathcal{P}_{\mathcal{F}_{(a,\delta)},N} = D_N\left(\mathbb{S}_{n_2} \otimes \mathbb{S}_{n_1}\right) F_N\left(\mathbb{S}_{n_2} \otimes \mathbb{S}_{n_1}\right),\tag{22}$$

where

$$D_N = (D_+ + D_- + E_+ + E_-)/4.$$

The motivation of the above construction is to create a preconditioner that properly acts on the different sources affecting the spectrum of  $\mathcal{M}_{(\alpha,\beta),N}$ . Specifically, the diagonal part operates on the spatial space treating the influence that the coefficients of the equation have on the matrix, while the  $\tau$  matrix focuses on the spectral space and the ill-conditioning generated by the discretization of the fractional differential operator. This observation is a direct result of the GLT symbol associated to  $\mathcal{M}_{(\alpha,\beta),N}$  and has been extensively studied in References 25 and 26, for the case of semi elliptic differential equations. In the simplest, but not unusual in applications, case where  $d_{\pm} = d$ ,  $e_{\pm} = e$ , we can counterbalance the influence of the term  $\frac{1}{r}$  in the spectrum of  $\mathcal{M}_{(\alpha,\beta),N}$  incorporating it into the  $\tau$  part of the preconditioner. Particularly, we define  $\hat{\mathcal{F}}_{(\alpha,\beta)}(\theta_1, \theta_2) = \frac{1}{r} + d \cdot q_{\alpha}(\theta_1) + \frac{s}{r}e \cdot q_{\beta}(\theta_2)$  replacing the sampling of  $\mathcal{F}_{(\alpha,\beta)}$  with that of  $\hat{\mathcal{F}}_{(\alpha,\beta)}$  for the construction of  $\hat{F}_N$  instead of  $F_N$  in (21). Accordingly, the new corresponding preconditioner  $\mathcal{P}_{\hat{f}_{(\alpha,\beta)},N}$  is defined as

$$\mathcal{P}_{\hat{\mathcal{F}}_{(q,\delta)}N} = \left(\mathbb{S}_{n_2} \otimes \mathbb{S}_{n_1}\right) \hat{\mathcal{F}}_N \left(\mathbb{S}_{n_2} \otimes \mathbb{S}_{n_1}\right).$$
<sup>(23)</sup>

The following theorem shows that in this case, the spectrum of the preconditioned matrix is bounded by positive constants independent of the size of the matrix.

**Theorem 3.** Assume that  $d_{\pm} = d > 0$ ,  $e_{\pm} = e > 0$ . In this case the coefficient matrix of the system becomes

$$A_{N} = \frac{1}{r} \mathbb{I}_{N} + (\mathbb{I}_{n_{2}} \otimes \hat{A}_{n_{1}}^{\alpha}) + (A_{n_{2}}^{\beta} \otimes \mathbb{I}_{n_{1}}) = \left(\mathbb{I}_{n_{2}} \otimes \left(\frac{1}{r} \mathbb{I}_{n_{1}} + \hat{A}_{n_{1}}^{\alpha}\right)\right) + (A_{n_{2}}^{\beta} \otimes \mathbb{I}_{n_{1}}) = \mathbb{I}_{n_{2}} \otimes A_{n_{1}}^{\alpha} + A_{n_{2}}^{\beta} \otimes \mathbb{I}_{n_{1}},$$
(24)

where

$$A_{n_1}^{\alpha} = \frac{1}{r} \mathbb{I}_N + T_{n_1} \left( d \cdot \left( w_{\alpha}(\theta) + w_{\alpha}(-\theta) \right) \right) = T_{n_1} \left( \frac{1}{r} + d \cdot q_{\alpha}(\theta) \right)$$
$$A_{n_2}^{\beta} = T_{n_2} \left( e \frac{s}{r} \cdot \left( w_{\beta}(\theta) + w_{\beta}(-\theta) \right) \right) = T_{n_2} \left( e \frac{s}{r} \cdot q_{\beta}(\theta) \right).$$

Then, the spectrum of the preconditioned matrix sequence  $\left\{ \mathcal{P}_{\hat{\mathcal{F}}_{(\alpha,\beta)},N}^{-1} A_N \right\}_N$  is bounded by positive constants c, C independent of N.

# 12 of 22 | WILEY

Proof. We recall that

$$h_x = (R_x - L_x)h_1, \quad h_y = (R_y - L_y)h_2,$$
  
 $r = \frac{h_t}{2h_x^{\alpha}}, \quad s = \frac{h_t}{2h_y^{\beta}},$ 

and

$$\hat{F}_N = \mathbb{I}_{n_2} \otimes F_{n_1}^{\alpha} + F_{n_2}^{\beta} \otimes \mathbb{I}_{n_1}, \tag{25}$$

where  $\mathbb{I}_n$  is the identity matrix of order *n* and

$$F_{n_1}^{\alpha} = \operatorname{diag}\left(d \cdot \mathcal{F}_{\alpha}(\theta_{i,n_1}) + \frac{1}{r}\right), \quad i = 1, \dots, n_1,$$
(26)

$$F_{n_2}^{\beta} = \operatorname{diag}\left(e\frac{s}{r} \cdot \mathcal{F}_{\beta}(\theta_{j,n_2})\right), \quad j = 1, \dots, n_2.$$

$$(27)$$

The matrix  $A_N$  is SPD since each of its terms is a Kronecker product of a diagonal with a SPD Toeplitz matrix. Hence,

$$\mathcal{P}_{N}^{-1}A_{N} = \left(\mathbb{S}_{n_{2}}\otimes\mathbb{S}_{n_{1}}\right)\hat{F}_{N}^{-1}\left(\mathbb{S}_{n_{2}}\otimes\mathbb{S}_{n_{1}}\right)A_{N},$$

which is similar to the matrix

$$\hat{F}_{N}^{-1/2}\left(\mathbb{S}_{n_{2}}\otimes\mathbb{S}_{n_{1}}\right)A_{N}\left(\mathbb{S}_{n_{2}}\otimes\mathbb{S}_{n_{1}}\right)\hat{F}_{N}^{-1/2}.$$

Thus,

$$\begin{split} \hat{F}_{N}^{-1/2}(S_{n_{2}}\otimes\mathbb{S}_{n_{1}})\left(\left(\mathbb{I}_{n_{2}}\otimes A_{n_{1}}^{\alpha}\right)+\left(A_{n_{2}}^{\beta}\otimes\mathbb{I}_{n_{1}}\right)\right)\left(\mathbb{S}_{n_{2}}\otimes\mathbb{S}_{n_{1}}\hat{F}_{N}^{-1/2}\right) \\ &=\hat{F}_{N}^{-1/2}\left(\left(\mathbb{S}_{n_{2}}\otimes\mathbb{S}_{n_{1}}\right)\left(\mathbb{I}_{n_{2}}\otimes A_{n_{1}}^{\alpha}\right)\left(\mathbb{S}_{n_{2}}\otimes\mathbb{S}_{n_{1}}\right)+\left(\mathbb{S}_{n_{2}}\otimes\mathbb{S}_{n_{1}}\right)\left(A_{n_{2}}^{\beta}\otimes\mathbb{I}_{n_{1}}\right)\left(\mathbb{S}_{n_{2}}\otimes\mathbb{S}_{n_{1}}\right)\right)\hat{F}_{N}^{-1/2} \\ &=\hat{F}_{N}^{-1/2}\left(\mathbb{I}_{n_{2}}\otimes\mathbb{S}_{n_{1}}A_{n_{1}}^{\alpha}S_{n_{1}}+\mathbb{S}_{n_{2}}A_{n_{2}}^{\beta}\mathbb{S}_{n_{2}}\otimes\mathbb{I}_{n_{1}}\right)\hat{F}_{N}^{-1/2} \\ &=\hat{F}_{N}^{-1/2}\left(\mathbb{I}_{n_{2}}\otimes(F_{n_{1}}^{\alpha})^{1/2}(F_{n_{1}}^{\alpha})^{-1/2}\mathbb{S}_{n_{1}}A_{n_{1}}^{\alpha}\mathbb{S}_{n_{1}}(F_{n_{1}}^{\alpha})^{-1/2}(\mathbb{F}_{n_{2}}^{\alpha})^{1/2}+\right.\\ &+\left.\left(F_{n_{2}}^{\beta}\right)^{1/2}(F_{n_{2}}^{\beta})^{-1/2}\mathbb{S}_{n_{2}}A_{n_{2}}^{\beta}\mathbb{S}_{n_{2}}(F_{n_{2}}^{\beta})^{-1/2}\mathbb{S}_{n_{1}}A_{n_{1}}^{\alpha}\mathbb{S}_{n_{1}}(F_{n_{2}}^{\alpha})^{-1/2}\right)(\mathbb{I}_{n_{2}}\otimes(F_{n_{1}}^{\alpha})^{1/2})+\right.\\ &+\left.\left.\left(F_{n_{2}}^{\beta}\right)^{1/2}\otimes\mathbb{I}_{n_{1}}\right)\left((F_{n_{2}}^{\beta})^{-1/2}\mathbb{S}_{n_{2}}A_{n_{2}}^{\beta}\mathbb{S}_{n_{2}}(F_{n_{2}}^{\beta})^{-1/2}\right)\mathbb{S}_{n_{1}}A_{n_{1}}^{\alpha}\mathbb{S}_{n_{1}}(F_{n_{1}}^{\alpha})^{-1/2}\right)(\mathbb{I}_{n_{2}}\otimes(F_{n_{1}}^{\alpha})^{1/2})+\right.\\ &+\left.\left.\left(F_{n_{2}}^{\beta}\right)^{1/2}\otimes\mathbb{I}_{n_{1}}\right)\left((F_{n_{2}}^{\beta})^{-1/2}\mathbb{S}_{n_{2}}A_{n_{2}}^{\beta}\mathbb{S}_{n_{2}}(F_{n_{2}}^{\beta})^{-1/2}\right)\mathbb{S}_{n_{1}}A_{n_{1}}^{\alpha}\mathbb{S}_{n_{1}}(F_{n_{2}}^{\alpha})^{-1/2}}\right)\mathbb{I}_{n_{1}}(F_{n_{2}}^{\beta})^{1/2}\otimes\mathbb{I}_{n_{1}}\right)\right)\hat{F}_{N}^{-1/2}\\ &=\left.\hat{F}_{N}^{-1/2}\left(\mathbb{I}_{n_{2}}\otimes(F_{n_{1}}^{\alpha})^{1/2}\right)L(\mathbb{I}_{n_{2}}\otimes(F_{n_{1}}^{\alpha})^{1/2})\hat{F}_{N}^{-1/2}}+\frac{\hat{F}_{N}^{-1/2}((F_{n_{2}}^{\beta})^{1/2}\otimes\mathbb{I}_{n_{1}})R((F_{n_{2}}^{\beta})^{1/2}\otimes\mathbb{I}_{n_{1}})\hat{F}_{N}^{-1/2}}\right)\\ &=\left.\hat{F}_{n_{2}}^{-1/2}\left(\mathbb{I}_{n_{2}}\otimes(F_{n_{1}}^{\alpha})^{1/2}\right)L(\mathbb{I}_{n_{2}}\otimes(F_{n_{1}}^{\alpha})^{1/2})\hat{F}_{N}^{-1/2}}+\frac{\hat{F}_{N}^{-1/2}((F_{n_{2}}^{\beta})^{1/2}\otimes\mathbb{I}_{n_{1}})R((F_{n_{2}}^{\beta})^{1/2}\otimes\mathbb{I}_{n_{1}})\hat{F}_{N}^{-1/2}}\right). \end{split}$$

Let

$$P_{n_1}^{\alpha} = \mathbb{S}_{n_1} F_{n_1}^{\alpha} \mathbb{S}_{n_1},$$
$$P_{n_2}^{\beta} = \mathbb{S}_{n_2} F_{n_2}^{\beta} \mathbb{S}_{n_2}.$$

Then, (see Reference 21), there exist positive constants c and C independent of  $n_1, n_2$ , such that

$$c < \sigma\left(\left(P_{n_1}^{\alpha}\right)^{-1}A_{n_1}^{\alpha}\right) < C \Rightarrow c < \sigma\left((F_{n_1}^{\alpha})^{-1/2}\mathbb{S}_{n_1}A_{n_1}^{\alpha}\mathbb{S}_{n_1}(F_{n_1}^{\alpha})^{-1/2}\right) < C,$$

and

$$c < \sigma \left( \left( P_{n_2}^{\beta} \right)^{-1} A_{n_2}^{\beta} \right) < C \Rightarrow c < (F_{n_2}^{\beta})^{-1/2} \mathbb{S}_{n_2} A_{n_2}^{\beta} \mathbb{S}_{n_2} (F_{n_2}^{\beta})^{-1/2} < C.$$

Consequently, for every normalized vector  $x \in \mathbb{R}^N$  we find that:

 $c < x^{\mathrm{T}}Lx < C, \qquad c < x^{\mathrm{T}}Rx < C.$ 

Since the matrices  $A_L$ ,  $A_R$  that form (28) are SPD, we recall some properties concerning such kind of matrices. Specifically, we use the inequality A > B for A, B SPD matrices if A - B > 0 is positive definite. In addition if A, B, C, D, and E are SPD, then

$$A > B \Leftrightarrow EAE > EBE, \tag{29}$$

WILEY <u>13 of 22</u>

$$A > B$$
 and  $C > D \Leftrightarrow A + C > B + D.$  (30)

Therefore, we infer

$$\begin{cases} c\mathbb{I}_N < L < C\mathbb{I}_N, \\ c\mathbb{I}_N < R < C\mathbb{I}_N, \end{cases}$$

and, using (29) and (30), we deduce

$$\begin{cases} c\hat{F}_{N}^{-1}(\mathbb{I}_{n_{2}}\otimes F_{n_{1}}^{\alpha}) < A_{L} < C\hat{F}_{N}^{-1}(\mathbb{I}_{n_{2}}\otimes F_{n_{1}}^{\alpha}), \\ c\hat{F}_{N}^{-1}(F_{n_{2}}^{\beta}\otimes \mathbb{I}_{n_{1}}) < A_{R} < C\hat{F}_{N}^{-1}(F_{n_{2}}^{\beta}\otimes \mathbb{I}_{n_{1}}). \end{cases}$$
(31)

Using again (29) and (30), taking into account the two inequalities of (31), and (25), we have

$$c\hat{F}_{N}^{-1}(\mathbb{I}_{n_{2}}\otimes F_{n_{1}}^{\alpha}) + c\hat{F}_{N}^{-1}(F_{n_{2}}^{\beta}\otimes\mathbb{I}_{n_{1}}) = c\hat{F}_{N}^{-1}\hat{F}_{N} = c\mathbb{I}_{N},$$
  
$$c\hat{F}_{N}^{-1}(\mathbb{I}_{n_{2}}\otimes F_{n_{1}}^{\alpha}) + c\hat{F}_{N}^{-1}(F_{n_{2}}^{\beta}\otimes\mathbb{I}_{n_{1}}) = c\hat{F}_{N}^{-1}\hat{F}_{N} = C\mathbb{I}_{N}.$$

Consequently, we conclude that

$$c\mathbb{I}_N \leq F_N^{-1/2}(\mathbb{S}_{n_1} \otimes \mathbb{S}_{n_2})A_N(\mathbb{S}_{n_1} \otimes \mathbb{S}_{n_2})F_N^{-1/2} \leq C\mathbb{I}_N.$$

Therefore, the spectrum of the preconditioned matrix, which is similar to the  $F_N^{-1/2}(\mathbb{S}_{n_1} \otimes \mathbb{S}_{n_2})A_N(\mathbb{S}_{n_1} \otimes \mathbb{S}_{n_2})F_N^{-1/2}$ , lies in [c, C]. Moreover, from Reference 21 we expect all the eigenvalues to be clustered at 1, something that is numerically confirmed in the next section.

**Corollary 1.** Let the functions  $d_+(x, y, t)$ ,  $d_-(x, y, t)$ ,  $e_+(x, y, t)$ ,  $e_-(x, y, t)$  being strictly positive functions on  $\Omega$ , with  $d_+(x, y, t) = d_-(x, y, t) = e_+(x, y, t) = e_-(x, y, t)$ . Then, the preconditioned matrix sequence  $\left\{ \mathcal{P}_{\hat{\mathcal{F}}_{(\alpha,\beta),N}}^{-1} \mathcal{M}_{(\alpha,\beta),N}^{(m)} \right\}_N$  is bounded by positive constants *c*, *C* independent of *N*.

*Proof.* The proof can be easily obtained from the results of Theorem 3 and the observation that the coefficient matrix in (20) can be bounded by

$$A_N^c \le \mathcal{M}_{(\alpha,\beta),N}^{(m)} \le A_N^C,$$

where

$$A_{N}^{c} = \frac{1}{r} \mathbb{I}_{N} + c(\mathbb{I}_{n_{2}} \otimes S_{\alpha,n_{1}}) + c(\mathbb{I}_{n_{2}} \otimes S_{\alpha,n_{1}}^{T}) + \frac{s \cdot c}{r} \left( (S_{\beta,n_{2}} \otimes \mathbb{I}_{n_{1}}) + (S_{\beta,n_{2}}^{T} \otimes \mathbb{I}_{n_{1}}) \right),$$
  
$$A_{N}^{C} = \frac{1}{r} \mathbb{I}_{N} + C(\mathbb{I}_{n_{2}} \otimes S_{\alpha,n_{1}}) + C(\mathbb{I}_{n_{2}} \otimes S_{\alpha,n_{1}}^{T}) + \frac{s \cdot C}{r} \left( (S_{\beta,n_{2}} \otimes \mathbb{I}_{n_{1}}) + (S_{\beta,n_{2}}^{T} \otimes \mathbb{I}_{n_{1}}) \right)$$

14 of 22 | WILEY-

and

$$\begin{split} c &= \min_{(x,y,t)\in\Omega} \{ d_+(x,y,t), d_-(x,y,t), e_+(x,y,t), e_-(x,y,t) \}, \\ C &= \max_{(x,y,t)\in\Omega} \{ d_+(x,y,t), d_-(x,y,t), e_+(x,y,t), e_-(x,y,t) \}. \end{split}$$

Then, using Rayleigh quotient we obtain

$$\mathcal{P}_{\hat{\mathcal{F}}_N}^{-1} A_N^c \leq \mathcal{P}_{\hat{\mathcal{F}}_N}^{-1} \mathcal{M}_{(\alpha,\beta),N}^{(m)} \leq \mathcal{P}_{\hat{\mathcal{F}}_N}^{-1} A_N^C$$
$$\lambda_1(\mathcal{P}_{\hat{\mathcal{F}}_N}^{-1} A_N^c) \leq \lambda_1(\mathcal{P}_{\hat{\mathcal{F}}_N}^{-1} \mathcal{M}_{(\alpha,\beta),N}^{(m)}) \leq \lambda_N(\mathcal{P}_{\hat{\mathcal{F}}_N}^{-1} \mathcal{M}_{(\alpha,\beta),N}^{(m)}) \leq \lambda_N(\mathcal{P}_{\hat{\mathcal{F}}_N}^{-1} A_N^C),$$

and the proof is completed.

In the subsequent section we report several numerical experiments which numerically confirm that a similar spectral behavior of the preconditioned matrix is expected in the more general case where the coefficients functions of the equation are all different to each other.

### 4 | NUMERICAL EXAMPLES

In this section we present three numerical examples to show the efficiency of the proposed preconditioners, compared with preconditioners discussed in Reference 3 (one dimension) and Reference 20 (two dimensions). We have chosen to compare our work with these works since they are the most recent and have shown their superiority against the other proposals in the literature.

- Example 1 is a one-dimensional problem, taken from Reference 3 Example 1. We compare and discuss the preconditioners therein with the proposed  $\mathcal{P}_{\mathcal{F}_{\alpha},n}$ , and a few variations based on the spectral symbol. The fractional derivatives are of order  $\alpha \in \{1.2, 1.5, 1.8\}$ .
- Example 2 is a two-dimensional problem, taken from Reference 20 Example 1. We compare and discuss the preconditioners therein with the proposed  $\mathcal{P}_{\mathcal{F}_{(\alpha,\beta)},N}$ . The fractional derivatives are  $\alpha = 1.8$  and  $\beta = 1.6$ .
- Example 3 is the same experiment as Example 2, but with the fractional derivatives  $\alpha = 1.8$  and  $\beta = 1.2$ .

The numerical experiments presented in Tables 1–4 were implemented in JULIA v1.1.0, using GMRES from the package ITERATIVESOLVERS.JL (GMRES tolerance is set to  $10^{-7}$ ) and the FFTW.JL package. Benchmarking is done with BENCHMARKTOOLS.JL with 100 samplings and minimum time is presented in milliseconds. Experiments were run, in serial, on a computer with dual Intel Xeon E5 2630 v4 2.20 GHz (10 cores each) CPUs, and with 128 GB of RAM.

Figures 1–3,5, and 6 show the scaled spectra of the preconditioned coefficient matrix  $\mathcal{P}^{-1}\mathcal{M}_{\alpha,n_1}$  (and  $\mathcal{P}^{-1}\mathcal{M}_{(\alpha,\beta),N}$ ) for different preconditioners  $\mathcal{P}$ , fractional derivatives  $\alpha$ , and matrix orders  $n_1$  (and  $\beta$ ,  $N = n_1, n_2$ ). The scaling by a constant  $c_0$  is performed as follows: we find the smallest disk enclosing all the eigenvalues of the considered matrix A. The center is denoted  $c_0$  and the radius is r. Then, the spectrum is scaled as  $\lambda_j(A)/c_0$  and the circle scaled and centered in (1, 0). The Julia package BOUNDINGSPHERE.JL was used to compute  $c_0$  and r for all figures. The current scaling of the eigenvalues of preconditioned coefficient matrices is a visualization of the important effect for the convergence rate of GMRES of both the clustering and of the shape of the clustering.

In Tables 1–4, for each preconditioner, we present the number of iterations (it), minimal timing (ms), and the condition number of the preconditioned matrix  $\kappa$ . The best results are highlighted in bold.

#### 4.1 | Example 1

We compare the proposed preconditioner  $\mathcal{P}_{\mathcal{F}_{a},n}$  with the ones presented in Example 1 from Reference 3 (and two alternative symbol-based preconditioners). We consider the one-dimensional form of (2) in the domain  $[L_1, R_1] \times [t_0, T] = [0, 2] \times [0, 1]$ , where the diffusion coefficients

$$\begin{split} d_+(x) &= \Gamma(3-\alpha) x^{\alpha}, \\ d_-(x) &= \Gamma(3-\alpha) (2-x)^{\alpha}. \end{split}$$

		I <sub>n1</sub>		$\mathcal{P}_{\mathbf{C},n_1}$	$\mathcal{P}_{C,n_1}$			$\mathcal{P}_{\mathrm{FULL},n_1}$			$\mathcal{P}_{\mathcal{F}_{\alpha},n_1}$		
α	$n_1 + 1$	(it)	(ms)	ĸ	(it)	(ms)	к	(it)	(ms)	к	(it)	(ms)	ĸ
1.2	2 <sup>6</sup>	28.0	1.7	9.6	13.0	9.6	3.3	14.0	3.8	1.6	7.2	2.3	30.8
	27	39.0	24.3	11.5	14.0	53.5	3.6	14.0	17.6	1.8	8.6	13.3	63.7
	2 <sup>8</sup>	46.0	114.9	13.4	13.0	119.8	3.8	14.0	68.8	2.0	9.9	58.2	132.2
	2 <sup>9</sup>	51.0	594.5	15.5	12.0	574.0	4.2	13.0	312.7	2.2	9.9	285.2	274.7
	$2^{10}$	54.0	2882.0	17.9	11.0	1927.0	4.5	12.0	1415.0	2.4	10.9	1450.0	571.4
	$2^{11}$	56.0	18569.0	20.5	10.0	11749.0	4.9	11.0	8840.0	2.5	12.8	9773.0	1189.7
1.5	2 <sup>6</sup>	32.0	2.0	33.4	12.0	8.8	7.1	13.0	3.2	1.8	6.7	2.2	16.1
	27	60.0	37.2	51.2	12.0	46.7	9.2	13.0	16.4	2.1	8.0	12.5	33.3
	2 <sup>8</sup>	89.0	213.1	75.8	12.0	111.3	12.0	13.0	64.5	2.3	8.5	52.6	70.9
	2 <sup>9</sup>	122.0	1389.0	109.9	12.0	544.2	15.8	12.0	288.9	2.6	10.0	280.2	152.7
	$2^{10}$	158.0	8007.0	157.7	11.0	1779.0	21.2	11.0	1366.0	2.9	10.0	1386.0	331.8
	$2^{11}$	195.0	56266.0	224.7	10.0	11538.0	28.6	10.0	8551.0	3.2	11.0	9142.0	724.3
1.8	2 <sup>6</sup>	32.0	2.1	136.5	9.0	6.6	23.0	10.0	2.6	2.6	6.1	2.2	9.7
	2 <sup>7</sup>	67.0	42.2	266.3	9.0	36.1	37.8	11.0	14.5	2.8	6.8	11.2	19.5
	2 <sup>8</sup>	131.0	332.3	494.8	9.0	89.8	63.0	10.0	53.6	2.9	7.0	47.2	40.8
	2 <sup>9</sup>	231.2	3085.0	893.8	9.0	446.8	106.3	9.0	257.9	2.9	8.6	262.8	86.9
	$2^{10}$	341.0	20620.0	1589.3	8.0	1503.0	180.5	8.0	1191.0	3.0	10.0	1370.0	187.5
	$2^{11}$	470.0	163700.0	2800.9	8.0	10197.0	308.3	7.0	7759.0	3.0	11.0	9125.0	408.1

**TABLE 1** Example 1: 1D,  $\alpha = \{1.2, 1.5, 1.8\}$ : Numerical experiments with GMRES and different preconditioners

*Note*: For each preconditioner we present: average number of iterations for one time step (it), total timing in milliseconds (ms) to attain the approximate solution at time *T*, and the condition number  $\kappa$  of the preconditioned matrix,  $\mathcal{P}^{-1}\mathcal{M}_{\alpha,n_1}$ . The best results are highlighted in bold.



**FIGURE 1** Example 1: 1D,  $\alpha = \{1.2, 1.5, 1.8\}$ : Scaled spectra of the resulting matrices when the preconditioners  $\mathbb{I}_{n_1}$ ,  $\mathcal{P}_{C,n_1}$ , and  $\mathcal{P}_{\text{FULL},n_1}$  are applied to the coefficient matrices  $\mathcal{M}_{\alpha,n_1}$  and  $n_1 = 2^6 - 1$ . Left:  $\alpha = 1.2$ . Middle:  $\alpha = 1.5$ . Right:  $\alpha = 1.8$ 

are non-constant in space. Furthermore, the source term is

$$f(x,t) = -32e^{-t}\left(x^2 + \frac{(2-x)^2(8+x^2)}{8} - \frac{3(x^3 + (2-x)^3)}{3-\alpha} + \frac{3(x^4 + (2-x)^4)}{(4-\alpha)(3-\alpha)}\right)$$

and the initial condition is

$$u(x,0) = 4x^2(2-x)^2,$$

		$\mathcal{P}_{1,n_1}$			$\mathcal{P}_{2,n_1}$	$\mathcal{P}_{2,n_1}$			$\mathcal{P}_{\mathrm{TRI},n_1}$			$\mathcal{P}_{\tilde{\mathcal{F}}_{\alpha},n_1}$		
α	$n_1 + 1$	(it)	(ms)	к	(it)	(ms)	к	(it)	(ms)	ĸ	(it)	(ms)	ĸ	
1.2	2 <sup>6</sup>	8.0	1.1	1.2	9.0	1.0	2.1	5.0	0.7	1.3	7.5	2.1	29.2	
	2 <sup>7</sup>	8.0	7.5	1.3	10.0	8.6	2.2	5.0	5.9	1.4	8.5	12.2	58.7	
	2 <sup>8</sup>	7.0	32.0	1.3	10.0	37.4	2.4	5.0	32.0	1.5	9.9	52.0	118.6	
	2 <sup>9</sup>	7.0	180.9	1.4	10.0	191.2	2.6	5.0	171.0	1.5	9.9	254.3	239.7	
	$2^{10}$	6.0	959.7	1.4	9.0	1066.0	2.8	5.0	928.7	1.6	11.0	1363.0	484.0	
	$2^{11}$	6.0	7026.0	1.5	9.0	7675.0	3.0	5.0	6914.0	1.7	12.0	10787.0	976.3	
1.5	2 <sup>6</sup>	16.0	1.5	2.5	8.0	1.0	2.1	7.0	1.0	2.4	8.7	2.7	13.6	
	27	20.0	14.4	3.1	9.0	8.1	2.3	8.0	7.5	3.0	8.0	12.1	26.3	
	2 <sup>8</sup>	24.0	67.9	4.0	9.0	35.3	2.7	11.0	40.2	4.0	8.4	47.7	51.8	
	2 <sup>9</sup>	26.0	366.7	5.2	10.0	197.5	3.0	13.0	227.3	5.4	9.9	248.1	103.0	
	$2^{10}$	27.0	1810.0	6.9	10.0	1105.0	3.5	15.0	1331.0	7.4	10.0	1636.0	205.9	
	$2^{11}$	25.4	11212.0	9.0	11.0	8179.0	4.0	18.0	9684.0	10.4	11.0	10563.0	424.5	
1.8	2 <sup>6</sup>	25.0	2.5	8.4	6.0	0.8	1.6	7.0	1.0	3.5	8.0	2.3	9.0	
	27	40.0	27.3	14.3	6.0	6.3	1.7	10.0	8.7	5.6	7.8	11.3	17.0	
	2 <sup>8</sup>	61.0	159.8	25.3	7.0	31.0	1.8	15.0	48.3	9.4	6.9	43.3	33.1	
	2 <sup>9</sup>	88.0	1083.0	44.7	7.0	170.1	2.0	22.0	325.4	16.6	7.0	222.4	65.4	
	$2^{10}$	120.0	6277.0	78.8	7.0	999.3	2.3	31.0	1983.0	30.0	8.9	1569.0	130.1	
	$2^{11}$	158.0	46716.0	138.2	7.0	7309.0	2.6	44.7	15756.0	54.6	10.0	10249.0	259.8	

**TABLE 2** Example 1: 1D,  $\alpha = \{1.2, 1.5, 1.8\}$ : Numerical experiments with GMRES and different preconditioners

*Note*: For each preconditioner we present the average number of iterations for one time step (it), the total timing in milliseconds (ms) to attain the approximate solution at time *T*, and the condition number  $\kappa$  of the preconditioned mass matrix,  $\mathcal{P}^{-1}\mathcal{M}_{\alpha,n_i}$ . The best results are highlighted in bold.

	$\mathbb{I}_{\!\boldsymbol{N}}$			$\mathcal{P}_{2,N}$			$\mathcal{P}_{MGM}$	N		$\mathcal{P}_{\mathcal{F}_{(lpha,eta)},N}$			
$n_1 = n_2$	(it)	(ms)	ĸ	(it)	(ms)	κ	(it)	(ms)	к	(it)	(ms)	к	
$2^{4}$	37.0	32.2	57.4	21.0	64.8	48.6	10.0	40.8	3.7	8.0	35.1	1.9	
2 <sup>5</sup>	73.0	331.4	167.4	17.6	551.1	31.7	11.0	383.1	5.4	8.0	296.8	2.7	
$2^{6}$	137.0	35440.0	429.4	17.0	10465.0	310.7	11.0	16146.0	8.2	9.0	6569.0	4.3	
27	251.0	1644134.0	966.8	17.0	213713.0	678.4	10.0	352471.0	12.2	9.0	135535.0	7.7	

*Note*: For each preconditioner we present the average number of iterations for one time step (it), the total timing in milliseconds (ms) to attain the approximate solution at time *T*, and the condition number  $\kappa$  of the preconditioned matrix,  $\mathcal{P}^{-1}\mathcal{M}_{(\alpha,\beta),N}$ . The best results are highlighted in bold.

leading to the analytical solution  $u(x, t) = 4e^{-t}x^2(2-x)^2$ . We assume  $h_x = h_t = 2/(n_1 + 1)$ , that is,  $v_{M,n_1} = h_x^{\alpha-1}$  and the number of time steps  $M = (n_1 + 1)T/(R_1 - L_1) = (n_1 + 1)/2$ . The set of fractional derivatives  $\alpha$ , for which a solution is computed for, is {1.2, 1.5, 1.8} and in addition we consider the following set of partial dimensions for  $n_1$ : { $2^6 - 1$ ,  $2^7 - 1$ ,  $2^8 - 1$ ,  $2^9 - 1$ }.

In Table 1 we present the results for the following preconditioners

- Identity  $(\mathbb{I}_{n_1})$ : GMRES without any preconditioner.
- Circulant ( $\mathcal{P}_{C,n_1}$ ): Described in Reference 17 and implemented using FFT.
- "Full" symbol ( $\mathcal{P}_{\text{FULL},n_1}$ ): Defined as  $\mathbb{S}_{n_1}$ diag ( $v_{M,n_1} + d_{+,i}g_\alpha(\theta_{j,n_1}) + d_{-,i}g_\alpha(-\theta_{j,n_1}), j = 1, 2, ..., n_1$ )  $\mathbb{S}_{n_1}$  and implemented using FFT.
- Symbol  $(\mathcal{P}_{\mathcal{F}_{\alpha},n_1})$ : Proposed in Section 3.1,  $D_{n_1}\mathbb{S}_{n_1}$  diag  $(p_{\alpha}(\theta_{j,n_1}), j = 1, 2, ..., n_1)\mathbb{S}_{n_1}$ , and implemented using FFT.

**TABLE 4** Example 3: 2D,  $\alpha = 1.8$ ,  $\beta = 1.2$ : Numerical experiments with GMRES and different preconditioners

	$\mathbb{I}_N$			$\mathcal{P}_{2,N}$			$\mathcal{P}_{\mathrm{MGM},}$	N		$\mathcal{P}_{\mathcal{F}_{(\alpha,\beta)},N}$		
$n_1 = n_2$	(it)	(ms)	к	(it)	(ms)	к	(it)	(ms)	ĸ	(it)	(ms)	к
2 <sup>4</sup>	49.0	37.1	57.8	26.5	79.7	42.8	18.0	39.0	8.2	10.0	37.0	1.9
2 <sup>5</sup>	92.0	394.0	162.9	32.0	713.8	104.0	26.0	450.7	16.7	12.0	329.0	2.7
2 <sup>6</sup>	173.0	44532.0	401.7	41.0	17197.0	231.6	33.0	35021.0	32.8	13.0	7493.0	4.4
27	316.0	2070478.0	876.4	51.0	438344.0	515.8	41.0	1107711.0	62.9	14.5	171500.0	7.9

*Note:* For each preconditioner we present: average number of iterations for one time step (it), total timing in milliseconds (ms) to attain the approximate solution at time *T*, and the condition number  $\kappa$  of the preconditioned matrix,  $\mathcal{P}^{-1}\mathcal{M}_{(\alpha,\beta),N}$ . The best results are highlighted in bold.



**FIGURE 2** Example 1: 1D,  $\alpha = \{1.2, 1.5, 1.8\}$ : Scaled spectra of the resulting matrices when the preconditioners  $\mathcal{P}_{\mathcal{F}_{\alpha},n_1}$  are applied to the coefficient matrices  $\mathcal{M}_{\alpha,n_1}$  for  $n_1 = 2^6 - 1$ 

In Figure 1 we present the scaled spectra of the resulting matrices, when the preconditioners  $\mathbb{I}_{n_1}$ ,  $\mathcal{P}_{C,n_1}$ , and  $\mathcal{P}_{\text{FULL},n_1}$  are applied to the coefficient matrices  $\mathcal{M}_{\alpha,n_1}$  when  $n_1 = 2^6 - 1$  and  $\alpha = 1.2$  (left),  $\alpha = 1.5$  (middle), and  $\alpha = 1.8$  (right). We conclude that the spectral behavior resulting from the circulant and "full" symbol preconditioner resemble each other, but the condition number is lower for the "full" symbol preconditioner, as seen in Table 1. In Figure 2 we show the scaled spectra of the resulting matrices when the preconditioners  $\mathcal{P}_{F_{\alpha},n_1}$  are applied to the coefficient matrices  $\mathcal{M}_{\alpha,n_1}$  with  $n_1 = 2^6 - 1$  and  $\alpha = \{1.2, 1.5, 1.8\}$ . We note that strong clustering of the eigenvalues of the preconditioned matrices with few large eigenvalues. The condition number is higher for the symbol preconditioner, compared to the "full" symbol preconditioner, however, as seen in Table 1 both the number of iterations and execution time are lower for the symbol preconditioner. This numerically confirms what we mentioned in Section 3.1 explaining the Equation (17), and the motivation of using a diagonal times a proper  $\tau$  as preconditioner. In detail, this two terms preconditioner properly acts on the different sources affecting the spectrum of the matrix: the diagonal part operates on the spatial space treating the influence that the coefficients of the equation have on the matrix, while the  $\tau$  matrix focuses on the spectral space and the ill-conditioning generated by the discretization of the fractional differential operator. Consequently, this better clustering observed in Figure 2, is the reason that the preconditioned GMRES method (see Reference 27) performs in general very well with this preconditioner. In Table 2 we present the results for the following preconditioners:



**FIGURE 3** Example 1: 1D,  $\alpha = \{1.2, 1.5, 1.8\}$ : Scaled spectra of the resulting matrices when the preconditioners  $\mathcal{P}_{1,n_1}$ ,  $\mathcal{P}_{2,n_1}$ , and  $\mathcal{P}_{\text{TRI},n_1}$  are applied to the matrices  $\mathcal{M}_{\alpha,n_1}$  and  $n_1 = 2^6 - 1$ . Left:  $\alpha = 1.2$ . Middle:  $\alpha = 1.5$ . Right:  $\alpha = 1.8$ 

- First derivative ( $\mathcal{P}_{1,n_1}$ ): Tridiagonal preconditioner based on the finite difference discretization of the first derivative, proposed in Reference 3 and implemented using the Thomas algorithm.
- Second derivative ( $\mathcal{P}_{2,n_1}$ ): Tridiagonal preconditioner based on the finite difference discretization of the second derivative, proposed in Reference 3 and implemented using the Thomas algorithm.
- Tridiagonal ( $\mathcal{P}_{\text{TRI},n_1}$ ): Tridiagonal preconditioner based on the three main diagonals of the coefficient matrix and implemented using the Thomas algorithm.
- Alternative symbol based  $(\mathcal{P}_{\tilde{\mathcal{F}}_{\alpha},n_1})$ : Constructed by  $\mathbb{S}_{n_1}D_{n_1}\operatorname{diag}(p_{\alpha}(\theta_{j,n_1}))\mathbb{S}_{n_1}$  and implemented using FFT.

As in Figure 1, in Figure 3 we present the scaled spectra of the preconditioned matrix. The spectral behavior of the three preconditioners (first and second derivative and the tridiagonal) for different values of  $\alpha$  correlate well with the results presented in Table 2. In the left panel of Figure 3 the best clustering is obtained using the tridiagonal preconditioner, followed by the first derivative, and then by the second derivative. Since  $\alpha = 1.2$ , a value close to one, this behavior is expected. When  $\alpha = 1.5$ , as presented in the middle panel of Figure 3, the results are similar for the three preconditioners, but the second derivative preconditioner performs in the best way as  $n_1$  increases. In the right panel of Figure 3 we see that the best clustering is observed for the second derivative preconditioner, and also show the best performances for all  $n_1$  and all reported quantities (iterations, timings, and condition numbers). The better performance of the preconditioners reported in Table 2 as opposed the ones in Table 1 is expected: this is due to the computational complexity of O(n) for the Thomas algorithm, as opposed to  $O(n \log n)$  for the DFT.

In Figure 4 we present the scaled spectrum of an alternative symbol based preconditioner,  $\mathcal{P}_{\bar{F}_a,n_1}$ , which performs slightly better than the proposed preconditioner  $\mathcal{P}_{\bar{F}_a,n_1}$  in Section 3.1 (compare Tables 1 and 2). This is mainly due to the avoided multiplication with the inverse of  $D_n$  for  $\mathcal{P}_{\bar{F}_a,n_1}$ , since the spectrum of the resulted preconditioned matrices using  $\mathcal{P}_{\bar{F}_a,n_1}$  and  $\mathcal{P}_{F_a,n_1}$  are comparable. Furthermore, in this case it seems that the most efficient choice of preconditioner is problem specific, depending on  $d_{\pm}$ .

#### 4.2 | Example 2

The considered two-dimensional example is originally from Reference 28 (Example 4) and is also discussed in Reference 20 (Example 1). In (2), define  $\alpha = 1.8$ ,  $\beta = 1.6$ , and

$$\begin{aligned} &d_{+}(x,y) = \Gamma(3-\alpha)(1+x)^{\alpha}(1+y)^{2}, \qquad d_{-}(x,y) = \Gamma(3-\alpha)(3-x)^{\alpha}(3-y)^{2}, \\ &e_{+}(x,y) = \Gamma(3-\beta)(1+x)^{2}(1+y)^{\beta}, \qquad e_{-}(x,y) = \Gamma(3-\beta)(3-x)^{2}(3-y)^{\beta}. \end{aligned}$$

The spatial domain is  $\Omega = [0, 2] \times [0, 2]$  and the time interval is  $[t_0, T] = [0, 1]$ . The initial condition is

$$u(x, y, 0) = u_0(x, y) = 16x^2y^2(2 - x)^2(2 - y)^2$$



WILEY <u>19 of 22</u>

**FIGURE 4** Example 1: 1D,  $\alpha = \{1.2, 1.5, 1.8\}$ : Scaled spectra of the resulting matrices when the preconditioners  $\mathcal{P}_{\tilde{\mathcal{F}}_{\alpha}, n_1}$  are applied to the matrices  $\mathcal{M}_{\alpha, n_1}$  for  $n_1 = 2^6 - 1$ 

and the source term is

$$f(x,y,t) = -16e^{-t} \left( x^2 (2-x)^2 y^2 (2-y)^2 + g_\alpha(x,y) + g_\alpha(2-x,2-y) + g_\beta(y,x) + g_\beta(2-y,2-x) \right),$$

where

$$g_{\gamma}(x,y) = \left(8x^{2-\gamma} - \frac{24x^{3-\gamma}}{3-\gamma} + \frac{24x^{4-\gamma}}{(4-\gamma)(3-\gamma)}\right)(1+x)^{\gamma}(1+y)^{2}y^{2}(2-y)^{2},$$

such that the solution to the FDE is given by  $u(x, y, t) = 16e^{-t}x^2(2-x)^2y^2(2-y)^2$ . Let  $h = h_x = h_y = 2/(n+1)$ , with  $n = n_1 = n_2 = M$ , and  $h_t = 1/(M+1)$ . Then,

$$\frac{1}{r} = \frac{2h^{\alpha}}{h_t} = \frac{2^{\alpha+1}M}{(n+1)^{\alpha}} = \frac{2^{\alpha+1}n}{(n+1)^{\alpha}}, \qquad \frac{s}{r} = \frac{h^{\alpha}}{h^{\beta}} = 2^{\alpha-\beta}(n+1)^{\beta-\alpha}.$$

In Table 3 (and also Table 4) we present the results for the following preconditioners:

- Second derivative (P<sub>2,N</sub>): Preconditioner based on the finite difference discretization of the second derivative, proposed in Reference 20 and implemented using one Galerkin projection multigrid V-cycle.
- Algebraic multigrid (*P*<sub>MGM,N</sub>): Preconditioner based on algebraic multigrid, proposed in Reference 20 and implemented using one algebraic multigrid V-cycle.
- Symbol ( $\mathcal{P}_{\mathcal{F}_{(a,b)}N}$ ): Proposed preconditioner and implemented using FFT.

We mention that in multi-dimensional setting, a negative results holds concerning the optimality of circulant algebra when it is used for preconditioning Toeplitz matrices generated by function with zeros of order greater than one (e.g., see References 29,30). Thus, we find a comparison with such kind of preconditioners to be unnecessary.

For details on the multigrid based preconditioners,  $\mathcal{P}_{2,N}$  (Galerkin projection multigrid) and  $\mathcal{P}_{MGM,N}$  (algebraic multigrid), see Reference 20. The proposed symbol-based preconditioner,  $\mathcal{P}_{\mathcal{F}_{(\alpha,\beta)},N}$ , performs better than the multigrid-based



**FIGURE 5** Example 2: 2D,  $\alpha = 1.8$ ,  $\beta = 1.6$ : Scaled spectra of the resulting matrices when the preconditioners are applied to the matrices  $\mathcal{M}_{(\alpha,\beta),n^2}$  and  $n_1 = 2^4$ . Left: Preconditioners  $\mathbb{I}_N$ ,  $\mathcal{P}_{2,N}$ , and  $\mathcal{P}_{MGM,N}$  Right: Preconditioner  $\mathcal{P}_{\mathcal{F}_{(\alpha,\beta)},N}$ 

preconditioners, as seen in Table 3. In Figure 5 we present the scaled spectra of the preconditioned matrices for  $N = n_1n_2 = 2^8$ . The clustering is better for the proposed symbol-based preconditioners than the other three, as seen comparing the left and right panels. We note in Table 3 that the number of iterations are essentially constant both for the algebraic multigrid and the symbol-based preconditioners.

By fine tuning of the parameters for the multigrid-based preconditioners, such as number of smoothing steps, W-cycles and so forth, these results might be improved. However, the simplicity of the proposed preconditioner, where no fine-tunings are required, is advantageous.

### 4.3 | Example 3

By modifying the coefficients  $\alpha = 1.8$  and  $\beta = 1.6$  in Example 2, to  $\alpha = 1.8$  and  $\beta = 1.2$  we obtain Example 3. In Table 4 we present the same type of computations as in Table 3. As discussed in Reference 20, the performance of the proposed multigrid-based preconditioners depend on the fractional derivatives  $\alpha$  and  $\beta$ . Since, in this example,  $\alpha$  and  $\beta$  differ more than in Example 2, and  $\beta$  is far away from two, we clearly see in Table 4 that the multigrid-based preconditioners perform worse than in Example 2. Especially note the worse behavior of the condition number for the algebraic multigrid-based preconditioner  $\mathcal{P}_{MGM,N}$ . The condition numbers are essentially the same for the symbol-based preconditioner  $\mathcal{P}_{F_{(\alpha,\beta)},N}$  in Examples 2 and 3.

In Figure 6 we present the same scaled spectra as in Figure 5, but regarding Example 3. Again, we note the advantageous clustering properties of the proposed symbol-based preconditioner in the right panel.

### 5 | CONCLUSIONS

The purpose of the article was the theoretical and numerical exploration of proper preconditioners based on the spectral symbols of the coefficient matrix for FDE problems. Beside the theoretical study, we have compared our results with past ones, especially those presented in References 3,20. As expected, and numerically shown in Example 1 which concerns the one-dimensional case, our the proposed preconditioners performs slightly worse, at least in sequential computations, than the tridiagonal preconditions, because of the computational complexity involved. However, in the more challenging



**FIGURE 6** Example 3: 2D,  $\alpha = 1.8$ ,  $\beta = 1.2$ : Scaled spectra of the resulting matrices when the preconditioners are applied to the coefficient matrices  $\mathcal{M}_{(\alpha,\beta),n_1^2}$ , and  $n_1 = 2^4$ . Left: Preconditioners  $\mathbb{I}_N$ ,  $\mathcal{P}_{2,N}$ , and  $\mathcal{P}_{MGM,N}$  Right: Preconditioner  $\mathcal{P}_{\mathcal{F}_{(\alpha,\beta)},N}$ 

two-dimensional case, as discussed in Examples 2 and 3, the proposed preconditioners do indeed perform better than the previously proposed multigrid-based preconditioners proposed and studied in Reference 20.

We note that future directions of research may include more complex problems, further analysis, and more extensive numerical experimentation. Also, problems where the fractional derivatives are close to three may be considered, since then we expect the symbol-based preconditioners to be even more advantageous, maybe even in the one-dimensional case.

#### ACKNOWLEDGMENTS

The authors thank the anonymous reviewer's for suggestions improving the quality of the manuscript. The second author was supported by the Grant "DRASI II" of Athens University of Economics and Business (Registration Number ER-3238).

#### **CONFLICT OF INTEREST**

This study does not have any conflicts to disclose.

#### DATA AVAILABILITY STATEMENT

Data sharing not applicable to this article as no datasets were generated or analyzed during the current study.

#### ORCID

Sven-Erik Ekström b https://orcid.org/0000-0002-7875-7543 Paris Vassalos b https://orcid.org/0000-0002-2131-7643

#### REFERENCES

- 1. Meerschaert MM, Tadjeran C. Finite difference approximations for fractional advection–Dispersion flow equations. J Comput Appl Math. 2004;172(1):65–77.
- 2. Saichev AI, Zaslavsky GM. Fractional kinetic equations: solutions and applications. Chaos. 1997;7(4):753-64.
- 3. Donatelli M, Mazza M, Serra-Capizzano S. Spectral analysis and structure preserving preconditioners for fractional diffusion equations. J Comput Phys. 2016;307:262–79.
- 4. Serra-Capizzano S. On the extreme eigenvalues of Hermitian (block) Toeplitz matrices. Linear Algebra Appl. 2021;270:109–29.
- Serra S. On the extreme spectral properties of Toeplitz matrices generated by L<sup>1</sup> functions with several minima (maxima). BIT Numer Math. 1996;34:135–42.

## 22 of 22 | WILEY-

- 6. Bottcer S, Grudsky S. On the condition numbers of large semi-definite Toeplitz matrices. Linear Algebra Appl. 1998;279:285-301.
- 7. Serra-Capizzano S, Tilli P. Extreme singular values and eigenvalues of non-Hermitian block Toeplitz matrices. J Comput Appl Math. 1999;108:113–30.
- 8. Chan RH, Ng MK. Conjugate gradient methods for Toeplitz systems. SIAM Rev. 1996;38:427-82.
- 9. Ng MK. Iterative methods for Toeplitz systems. New York, NY: Oxford University Press, Inc; 2004.
- 10. Huckle T, Serra-Capizzano S, Tablino-Possio C. Preconditioning strategies for non-Hermitian Toeplitz linear systems. Numer Linear Algebra. 2005;12:211–20.
- 11. Baerland T, Kuchta M, Mardal K. Multigrid methods for discrete fractional Sobolev spaces. SIAM J Sci Comput. 2019;41(2):A948-72.
- 12. Harizanov S, Lazarov R, Margenov S, Marinov P, Pasciak J. Analysis of numerical methods for spectral fractional elliptic equations based on the best uniform rational approximation. J Comput Phys. 2020;408:109285.
- 13. Khristenko U, Wohlmuth B. Solving time-fractional differential equation via rational approximation; 2021. arXiv:210205139.
- Meerschaert MM, Tadjeran C. Finite difference approximations for two-sided space-fractional partial differential equations. Appl Numer Math. 2006;56(1):80–90.
- Wang H, Wang K, Sircar T. A direct 𝒪 (Nlog<sup>2</sup>N) finite difference method for fractional diffusion equations. J Comput Phys. 2010;229(21):8095–104.
- 16. Pang HK, Sun HW. Multigrid method for fractional diffusion equations. J Comput Phys. 2012;231(2):693–703.
- 17. Lei SL, Sun HW. A circulant preconditioner for fractional diffusion equations. J Comput Phys. 2013;242:715–25.
- Pan J, Ng MK, Wang H. Fast iterative solvers for linear systems arising from time-dependent space-fractional diffusion equations. SIAM J Sci Comput. 2016;38(5):A2806–26.
- Lin XL, Ng MK, Sun HW. A splitting preconditioner for Toeplitz-like linear systems arising from fractional diffusion equations. SIAM J Matrix Anal Appl. 2017;38(4):1580–614.
- 20. Moghaderi H, Dehghan M, Donatelli M, Mazza M. Spectral analysis and multi-grid preconditioners for two-dimensional space-fractional diffusion equations. J Comput Phys. 2017;350:992–1011.
- 21. Noutsos D, Serra-Capizzano S, Vassalos P. Essential spectral equivalence via multiple step preconditioning and applications to ill conditioned Toeplitz matrices. Linear Algebra Appl. 2016;491:276–91.
- 22. Garoni C, Serra-Capizzano S. Generalized locally Toeplitz sequences: theory and applications. Vol 1. Berlin, Germany: Springer International Publishing; 2017.
- 23. Zamarashkin NL, Tyrtyshnikov EE. Distribution of the eigenvalues and singular numbers of Toeplitz matrices under weakened requirements on the generating function. Mat Sb. 1997;188(8):83–92.
- 24. Noutsos D, Vassalos P. Superlinear convergence for PCG using band plus algebra preconditioners for Toeplitz systems. Comput Math Appl. 2008;56(5):1255–70.
- 25. Noutsos D, Serra-Capizzano S, Vassalos P. The conditioning of FD matrix sequences coming from semi-elliptic differential equations. Linear Algebra Appl. 2008;428(2-3):600–24.
- 26. Vassalos P. Asymptotic results on the condition number of FD matrices approximating semi-elliptic PDEs. Electron J Linear Algebra. 2018;34:566–81.
- 27. Axelsson O, Lindskog G. On the rate of convergence of the preconditioned conjugate gradient method. Numer Math. 1986;48(5):499-523.
- 28. Pang HK, Sun HW. Fast numerical contour integral method for fractional diffusion equations. J Sci Comput. 2015;66(1):41–66.
- 29. Noutsos D, Serra-Capizzano S, Vassalos P. Matrix algebra preconditioners for multilevel Toeplitz systems do not insure optimal convergence rate. Theor Comput Sci. 2004;315(2):557–79.
- 30. Noutsos D, Serra-Capizzano S, Vassalos P. Spectral equivalence and matrix algebra pre-conditioners for multilevel Toeplitz systems: a negative result. Contemp Math. 2003;323:313–22.

**How to cite this article:** Barakitis N, Ekström S-E, Vassalos P. Preconditioners for fractional diffusion equations based on the spectral symbol. Numer Linear Algebra Appl. 2022;e2441. https://doi.org/10.1002/nla.2441