# Preconditioning Techniques for Diagonal-times-Toeplitz Matrices in Fractional Diffusion Equations

Jianyu Pan [*]     Rihuan Ke [†]     Michael K. Ng [‡]     Hai-Wei Sun [§]

## Abstract

The fractional diffusion equation is discretized by an implicit finite difference scheme with the shifted Grünwald formula, which is unconditionally stable. The coefficient matrix of the discretized linear system is equal to the sum of a scaled identity matrix and two diagonal-times-Toeplitz matrices. Standard circulant preconditioners may not work for such Toeplitz-like linear systems. The main aim of this paper is to propose and develop approximate inverse preconditioners for such Toeplitz-like matrices. The construction of an approximate inverse preconditioner is to approximate the inverses of weighted Toeplitz matrices by circulant matrices, and then combine them together row-by-row. Because of Toeplitz structure, both the discretized coefficient matrix and the preconditioner can be implemented very efficiently by using fast Fourier transforms. Theoretically, we show that the spectra of the resulting preconditioned matrices are clustered around one. Thus Krylov subspace methods with the proposed preconditioner converge very fast. Numerical examples are given to demonstrate the effectiveness of the proposed preconditioner and show that its performance is better than the other testing preconditioners.

**Key words:** Fractional diffusion equation, Toeplitz matrix, approximate inverse, circulant matrix, Fast Fourier transform, Krylov subspace methods
**Mathematics Subject Classification:** 65F10, 65L12, 65L20, 65T50, 26A33

# 1 Introduction

In this paper, we consider an initial boundary value problem of a fractional diffusion equation (FDE)

$$\frac{\partial u(x,t)}{\partial t} = d_+(x,t)\frac{\partial^\alpha u(x,t)}{\partial_+ x^\alpha} + d_-(x,t)\frac{\partial^\alpha u(x,t)}{\partial_- x^\alpha} + f(x,t),$$
$$x \in [x_L, x_R], \quad t \in (0, T_f], \tag{1}$$
$$u(x_L, t) = u(x_R, t) = 0, \ 0 \le t \le T_f,$$
$$u(x, 0) = u_0(x), \ x_L \le x \le x_R,$$

where $1 < \alpha < 2$, $f(x,t)$ is the source term, and the diffusion coefficients satisfying $d_+(x,t) \ge 0$ and $d_-(x,t) \ge 0$. Here the left-sided and the right-sided fractional derivatives are defined in the Grünwald-Letnikov form [27]:

$$\frac{\partial^\alpha u(x,t)}{\partial_+ x^\alpha} = \lim_{h \to 0^+} \frac{1}{h^\alpha} \sum_{k=0}^{(\lfloor x-x_L \rfloor)/h} g_k^{(\alpha)} \, u(x - kh, t),$$

$$\frac{\partial^\alpha u(x,t)}{\partial_- x^\alpha} = \lim_{h \to 0^+} \frac{1}{h^\alpha} \sum_{k=0}^{(\lfloor x_R-x \rfloor)/h} g_k^{(\alpha)} \, u(x + kh, t),$$

where $\lfloor x \rfloor$ denotes the floor of $x$ and the coefficients $g_k^{(\alpha)}$ are defined as follows

$$g_0^{(\alpha)} = 1 \quad \text{and} \quad g_k^{(\alpha)} = (-1)^k \frac{\alpha(\alpha-1)\cdots(\alpha-k+1)}{k!}, \quad k = 1, 2, \ldots. \tag{2}$$

In last few decades, more and more anomalous diffusion phenomena have been found in the real world, which lead to FDEs. The FDEs were shown to provide an adequate and accurate description for these anomalous diffusions which include modeling chaotic dynamics of classical conservative systems [40], groundwater contaminant transport [3, 4], turbulent flow [7, 29], and applications in biology [19], finance [28], image processing [1], and physics [30]. In [39], Wang and Yang studied and analyzed variable-coefficient conservative fractional elliptic differential equations. Usually, closed-form analytical solutions of FDEs are not available. Several numerical methods for solutions of FDEs are proposed and developed; see, for instances, [6, 13, 14, 16, 18, 20, 21, 22, 23, 31, 33, 35].

One of the main characteristics of the fractional differential operator is nonlocal. It was shown that a simple discretization scheme of the FDE, even though implicit, leads to be unconditionally unstable [21, 22]. Moreover, most numerical methods for FDEs tend to generate full coefficient matrices, which require computational cost of $\mathcal{O}(N^3)$ and storage of $\mathcal{O}(N^2)$, where $N$ is the number of grid points [37]. It is quite different from the second-order diffusion equations which usually yield sparse coefficient matrices with $\mathcal{O}(N)$ nonzero entries and can be solved very efficiently by fast iterative methods with $\mathcal{O}(N)$ complexity. In order to keep a stable numerical scheme, Meerschaet and Tadjeran [21, 22] proposed a shifted Grünwald discretization to approximate FDEs. Their method has been shown to be unconditionally stable. Furthermore, the full coefficient matrix by the

Meerschaet-Tadjeran's method holds Toeplitz-like structure [37], which can be written as the sum of the scaled identify matrix and two diagonal-times-Toeplitz matrices. Thus the storage requirement is significantly reduced from $\mathcal{O}(N^2)$ to $\mathcal{O}(N)$. Using the fast Fourier transform (FFT), the Toeplitz matrix-vector multiplication can be done in $\mathcal{O}(N \log N)$ operations [10, 24]. Therefore, the computational cost per iteration keeps $\mathcal{O}(N \log N)$ operations when the conjugate gradient normal residual (CGNR) method is applied to solve the discretized system [38]. Nevertheless, the resulting system in general is ill-conditioned and hence the iterative method converges slowly. To speed up the convergence rate, Pang and Sun [26] proposed a multigrid method, which is developed from [8, 34], to solve the discretized system of the FDE by the Meerschaet-Tadjeran method. With the damped-Jacobi method as the smoother, the multigrid algorithm can preserve the computational cost per iteration as $\mathcal{O}(N \log N)$ operations. Numerical results have shown that their multigrid method converges very fast. However, from the theoretical point of view, the linear convergence of their multigrid method, even for the case where both diffusion coefficients are equal and constant, has not been shown in the literature.

As the resulting discretized systems are Toeplitz-like, we may consider circulant preconditioning techniques for such systems. Circulant preconditioners for Toeplitz matrices have been theoretically and numerically studied with numerous applications for over twenty years; see [9, 10, 24]. Recently, Lei and Sun [17] applied the preconditioned CGNR method with a circulant preconditioner, which is extended from the Strang circulant preconditioner [11], to solve the discretized system of the FDE by the Meerschaet-Tadjeran method. The spectrum of the preconditioned matrix is theoretically proven to be clustered around one providing that both diffusion coefficients are constant, and hence the superlinear convergence rate is obtained. However, when the diffusion coefficients are not constant, the spectrum of the preconditioned matrix is no longer clustered around one. One possible approach is to approximate a Toeplitz matrix by a circulant matrix and then use diagonal-times-circulant matrix as the preconditioner for the discretized system. However, the main difficulty of this approach is that the resulting preconditioner is not circulant, and its inverse cannot be determined efficiently. Indeed, the cost of the inversion is about the same as that of the inversion of the original discretized matrix.

Recently, Ng and Pan [25] proposed approximate inverse circulant-plus-diagonal preconditioners for solving Toeplitz-plus-diagonal systems. Their idea is to use circulant matrices to approximate the inverses of Toeplitz matrices and then combine these circulant matrices together. As the resulting preconditioner is already of the inverted form, only matrix-vector multiplications are required in the preconditioning step. Therefore, the resulting preconditioner can be efficiently implemented. They also showed that the spectrum of the preconditioned Toeplitz-plus-diagonal matrix is clustered around one. Numerical examples including the application of image restoration have demonstrated that their approximate inverse preconditioner is very effective, and the Krylov subspace method converges very fast when it is applied to solve these preconditioned systems.

The main aim of this paper is to propose and develop approximate inverse preconditioners for the sum of the scaled identity matrix and two diagonal-times-Toeplitz matrices arising from the discretization of FDEs. The construction of an approximate inverse preconditioner is to approximate the inverses of scaled Toeplitz matrices by circulant matrices,

and then combine them together row-by-row. We remark that this idea of construction is similar to that in [25], but two linear systems are different. Because of Toeplitz structure, both the discretized coefficient matrix and the preconditioner can be implemented very efficiently by the FFT. The computational cost per iteration is of $\mathcal{O}(N \log N)$. Theoretically, we show that the spectra of the resulting preconditioned matrices are clustered around one, and thus Krylov subspace methods with the proposed preconditioner converge very quickly. Numerical examples are given to demonstrate the effectiveness of the proposed preconditioner and show that its performance is better than the other testing preconditioners.

The paper is organized as follows. In Section 2, we present the discretized system. In Section 3, we construct the proposed preconditioner. We also analyze the spectrum of the preconditioned matrices. In Section 4, numerical examples are given to demonstrate the performance of the proposed preconditioner. Finally, concluding remarks are given in Section 5.

## 2  Discretization of FDEs

Let $h = (x_R - x_L)/(N + 1)$ and $\Delta t = T_f/M$ be the sizes of spatial grid and time step, respectively, where $N$ and $M$ are positive integers. We define a spatial and temporal partition $x_i = x_L + ih$ for $i = 0, 1, 2, \ldots, N + 1$ and $t_m = m\Delta t$ for $m = 0, 1, 2, \ldots, M$, and denote

$$u_i^{(m)} = u(x_i, t_m), \quad d_{+,i}^{(m)} = d_+(x_i, t_m), \quad d_{-,i}^{(m)} = d_-(x_i, t_m), \quad \text{and} \quad f_i^{(m)} = f(x_i, t_m).$$

The first-order time derivative in (1) can be discretized by a standard first-order time difference quotient. For the discretization of the fractional spatial derivative, we employ the following shifted Grünwald approximations

$$\frac{\partial^\alpha u(x_i, t_m)}{\partial_+ x^\alpha} = \frac{1}{h^\alpha} \sum_{k=0}^{i+1} g_k^{(\alpha)} u_{i-k+1}^{(m)} + O(h),$$

$$\frac{\partial^\alpha u(x_i, t_m)}{\partial_- x^\alpha} = \frac{1}{h^\alpha} \sum_{k=0}^{N-i+2} g_k^{(\alpha)} u_{i+k-1}^{(m)} + O(h),$$

which is proposed by Meerschaert and Tadjeran [21, 22]. They proved that the corresponding implicit finite difference scheme

$$\frac{u_i^{(m)} - u_i^{(m-1)}}{\Delta t} = \frac{d_{+,i}^{(m)}}{h^\alpha} \sum_{k=0}^{i+1} g_k^{(\alpha)} u_{i-k+1}^{(m)} + \frac{d_{-,i}^{(m)}}{h^\alpha} \sum_{k=0}^{N-i+2} g_k^{(\alpha)} u_{i+k-1}^{(m)} + f_i^{(m)}, \tag{3}$$

$$i = 1, 2, \ldots, N, \ m = 1, 2, \ldots, M,$$

is unconditionally stable. By the boundary condition, we have $u_0^{(m)} = u_{N+1}^{(m)} = 0$ for $m = 1, 2, \ldots, M$. Denote

$$u^{(m)} = \left[u_1^{(m)}, u_2^{(m)}, \ldots, u_N^{(m)}\right]^\mathsf{T} \in \mathbb{R}^N, \quad f^{(m)} = \left[f_1^{(m)}, f_2^{(m)}, \ldots, f_N^{(m)}\right]^\mathsf{T} \in \mathbb{R}^N.$$

Then we can rewrite the numerical scheme (3) into the following matrix form

$$A^{(m)}u^{(m)} = \eta u^{(m-1)} + h^\alpha f^{(m)}, \quad m = 1, 2, \ldots, M, \tag{4}$$

where the coefficient matrix $A^{(m)}$ is of the form [37]

$$A^{(m)} = \eta I + D^{(m)}T + W^{(m)}T^\intercal, \tag{5}$$

with $\eta = h^\alpha/\Delta t$. Here $T^\intercal$ denotes the transpose of $T$, $I \in \mathbb{R}^{N \times N}$ is the identity matrix, $D^{(m)}$ and $W^{(m)}$ are diagonal matrices defined by

$$\begin{cases} D^{(m)} = \mathrm{diag}\left(d_{+,1}^{(m)}, d_{+,2}^{(m)}, \ldots, d_{+,N}^{(m)}\right) \\ W^{(m)} = \mathrm{diag}\left(d_{-,1}^{(m)}, d_{-,2}^{(m)}, \ldots, d_{-,N}^{(m)}\right) \end{cases}$$

and $T$ is the Toeplitz matrix

$$T = - \begin{bmatrix} g_1^{(\alpha)} & g_0^{(\alpha)} & 0 & \cdots & 0 \\ g_2^{(\alpha)} & \ddots & \ddots & \ddots & \vdots \\ g_3^{(\alpha)} & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & g_0^{(\alpha)} \\ g_N^{(\alpha)} & \cdots & g_3^{(\alpha)} & g_2^{(\alpha)} & g_1^{(\alpha)} \end{bmatrix}. \tag{6}$$

# 3 The Preconditioning Method

As $A^{(m)}$ is nonsymmetric, we can apply Krylov subspace methods, such as GMRES, to solve the linear systems (4). In order to improve the performance and reliability of the Krylov subspace methods, preconditioning is often employed. It is widely recognized that preconditioning is the most critical ingredient in the development of efficient solvers for challenging problems in scientific computations [5].

In the following, we consider the preconditioners for the matrix $A \in \mathbb{R}^{N \times N}$ of the following form

$$A = \eta I + DT + WT^\intercal, \tag{7}$$

where $\eta > 0$, $I$ is the identity matrix, $D$ and $W$ are diagonal matrices with nonnegative diagonal entries and $T$ is the Toeplitz matrix defined in (6). Here we assume that the diagonals of $D = \mathrm{diag}(d_1, d_2, \ldots, d_N)$ and $W = \mathrm{diag}(w_1, w_2, \ldots, w_N)$ are determined by the nonnegative functions $d(x)$ and $w(x)$ on $[x_L, x_R]$, respectively; i.e.,

$$d_i = d(x_i) \quad \text{and} \quad w_i = w(x_i), \quad i = 1, 2, \ldots, N.$$

Define

$$K_i \triangleq \eta I + d_i T + w_i T^\intercal = \eta I + d(x_i)T + w(x_i)T^\intercal, \quad i = 1, 2, \ldots, N. \tag{8}$$

Clearly, all $K_i$'s are Toeplitz matrices. According to the fact that

$$e_i^\intercal A = e_i^\intercal K_i,$$

5

our preconditioner is based on the following approximation

$$e_i^\mathsf{T} A^{-1} \approx e_i^\mathsf{T} K_i^{-1},$$

where $e_i$ denotes the $i$-th column of the identity matrix. This means that the $i$-th row of the inverse of $A$ is approximated by the $i$-th row of the inverse of $K_i$. Therefore, we propose the following preconditioner $B_1$ whose inverse is defined by

$$B_1^{-1} = \sum_{i=1}^{N} e_i e_i^\mathsf{T} K_i^{-1}. \tag{9}$$

To construct $B_1^{-1}$, we need to compute the inverse of $K_i$ $(i = 1, 2, \ldots, N)$, which is impractical. However, as $K_i$ is a Toeplitz matrix, we can approximate $K_i$ by a circulant matrix. Let $C$ be the Strang circulant approximation [11] of the Toeplitz matrix $T$, that is, the first column of the circulant matrix $C$ is given by

$$- [g_1^{(\alpha)}, \ g_2^{(\alpha)}, \ \cdots, \ g_{\lfloor (N+1)/2 \rfloor}^{(\alpha)}, \ 0, \ \cdots, \ 0, \ g_0^{(\alpha)}]^\mathsf{T} \tag{10}$$

We remark that other successful preconditioners can be considered and used; see, for instance, [24]. Let

$$C_i = \eta I + d_i C + w_i C^\mathsf{T} = \eta I + d(x_i) C + w(x_i) C^\mathsf{T}, \quad i = 1, 2, \ldots, N. \tag{11}$$

Then we obtain the preconditioner $B_2$ with

$$B_2^{-1} = \sum_{i=1}^{N} e_i e_i^\mathsf{T} C_i^{-1}, \tag{12}$$

which is based on the circulant matrices. It is well known that circulant matrices can be diagonalized in $O(N \log N)$ operations by making use of FFT. Hence the product $C_i^{-1} y$ for any vector $y$ can be computed by FFT in $O(N \log N)$ operations.

By the definition of $B_2$, we know that implementing a preconditioner based on $B_2$ would require $O(N)$ FFT's per iteration. This is still expensive. In order to reduce the computational workload, we propose to use the interpolation method to construct the practical preconditioner; see [25]. We first choose a small number $\ell$ $(\ell \ll N)$ of values $\{\tilde{x}_j\}_{j=1}^{\ell} \subset \{x_i\}_{i=1}^{N}$, which covers (most of) the range of values of $\{x_i\}_{i=1}^{N}$. The idea is given as follows. First, let $\lambda$ be a certain complex number with positive real part, i.e., $\mathrm{Re}(\lambda) > 0$, and define the function

$$q_\lambda(x) \triangleq \frac{1}{\eta + \lambda d(x) + \bar{\lambda} w(x)}, \quad x \in [x_L, x_R], \tag{13}$$

where $\bar{\lambda}$ denotes the complex conjugate of $\lambda$. Let

$$p_\lambda(x) = \phi_1(x) q_\lambda(\tilde{x}_1) + \phi_2(x) q_\lambda(\tilde{x}_2) + \cdots + \phi_\ell(x) q_\lambda(\tilde{x}_\ell) \tag{14}$$

be the piecewise linear interpolation for $q_\lambda(x)$ based on the $\ell$ points

$$\left\{ \left( \tilde{x}_j, q_\lambda(\tilde{x}_j) \right) \right\}_{j=1}^\ell \subset \left\{ \left( x_i, q_\lambda(x_i) \right) \right\}_{i=1}^N .$$

Next, we precompute the eigenvalues of the circulant matrix $C$ by FFT, that is

$$C = F \Lambda F^*,$$

where $F$ is the Fourier matrix and $\Lambda_j = \mathrm{diag}(\lambda_1, \lambda_2, \ldots, \lambda_N)$ is a diagonal matrix whose diagonals are the eigenvalues of $C$. Let $\tilde{C}_j \triangleq \eta I + d(\tilde{x}_j)C + w(\tilde{x}_j)C^\mathsf{T}$. Then we have

$$\tilde{C}_j = F \tilde{\Lambda}_j F^*, \quad j = 1, 2, \ldots, \ell,$$

where $F^*$ is the conjugate transpose of $F$, and $\tilde{\Lambda}_j = \eta I + d(\tilde{x}_j)\Lambda + w(\tilde{x}_j)\Lambda^*$.

Finally, we apply interpolation formula (14) to approximate $C_i^{-1}$ :

$$C_i^{-1} \approx F \left( \sum_{j=1}^\ell \phi_j(x_i) \tilde{\Lambda}_j^{-1} \right) F^* = \sum_{j=1}^\ell \phi_j(x_i) \tilde{C}_j^{-1}, \quad i = 1, 2, \ldots, N,$$

where $\phi_j(x_i)$ are the interpolation coefficients. Then we can get the practical preconditioner $B_3$ with

$$
\begin{aligned}
B_3^{-1} &= \sum_{i=1}^N e_i e_i^\mathsf{T} \left( \sum_{j=1}^\ell \phi_j(x_i) \tilde{C}_j^{-1} \right) \\
&= \sum_{i=1}^N e_i e_i^\mathsf{T} F \left( \sum_{j=1}^\ell \phi_j(x_i) \tilde{\Lambda}_j^{-1} \right) F^* \\
&= \sum_{i=1}^N \sum_{j=1}^\ell e_i e_i^\mathsf{T} F \left( \phi_j(x_i) \tilde{\Lambda}_j^{-1} \right) F^* \\
&= \sum_{j=1}^\ell \left( \sum_{i=1}^N e_i e_i^\mathsf{T} \phi_j(x_i) \right) F \tilde{\Lambda}_j^{-1} F^* \\
&= \sum_{j=1}^\ell \Phi_j F \tilde{\Lambda}_j^{-1} F^*,
\end{aligned}
\tag{15}
$$

where $\Phi_j = \mathrm{diag}\left( \phi_j(x_1), \phi_j(x_2), \ldots, \phi_j(x_N) \right)$ are diagonal matrices. Now applying $B_3^{-1}$ to any vector requires about $O(\ell N \log N)$ operations which is acceptable for a moderate number $\ell$.

It is expected that as the number of interpolation points increases, the number of iterations required for convergence decreases. However, the cost of forming and applying the preconditioner grows proportionally to the number of interpolation points. Hence, there is a trade-off to determine the number of interpolation points. In the next section, we will analyze the spectra of the preconditioned matrices.

# 4   Analysis of Preconditioners

One of the important aspect that affect the convergence property of the Krylov subspace methods is the eigenvalue distribution of the (preconditioned) coefficient matrix. In general, a clustered spectrum away from zero often results in faster convergence, especially for those matrices close to normal [5]. In this section, we consider the spectral properties of the preconditioned matrix $B_2^{-1}A$.

We first introduce the off-diagonal decay property.

**Definition 4.1.** [32] Let $A = [a_{i,j}]_{i,j \in \mathcal{I}}$ be a matrix, where the index set is $\mathcal{I} = \mathbb{Z}, \mathbb{N}$, or $\{1, 2, \ldots, N\}$. We say $A$ belongs to the class $\mathcal{L}_s$ if

$$|a_{i,j}| \leq \frac{c}{(1 + |i - j|)^s} \quad \text{for } s > 1, \tag{16}$$

and some constant $c > 0$.

The following result is due to [15, 32].

**Lemma 4.1.** *Let $A = [a_{i,j}]_{i,j \in \mathcal{I}}$ be a nonsingular matrix, where the index set is $\mathcal{I} = \mathbb{Z}, \mathbb{N}$, or $\{1, 2, \ldots, N\}$. If $A \in \mathcal{L}_s$ for some $s > 1$, then $A^{-1} \in \mathcal{L}_s$.*

Now we investigate the off-diagonal decay properties of the Toeplitz matrix $T$ defined by (6). For the entries of $T$, we have the following result [21, 22, 27, 37].

**Lemma 4.2.** *Let $g_k^{(\alpha)}$ be defined in (2) with $1 < \alpha < 2$. Then the following recursive relationship holds*

$$g_k^{(\alpha)} = \left(1 - \frac{\alpha + 1}{k}\right) g_{k-1}^{(\alpha)}, \quad k = 1, 2, \ldots.$$

*Moreover, we have*

$$g_0^{(\alpha)} = 1, \quad g_1^{(\alpha)} = -\alpha < 0, \quad 1 > g_2^{(\alpha)} > g_3^{(\alpha)} > \cdots > 0,$$

*and*

$$\sum_{k=0}^{\infty} g_k^{(\alpha)} = 0, \quad \sum_{k=0}^{m} g_k^{(\alpha)} < 0, \quad 1 \leq m < \infty.$$

By Lemma 4.2 and the definition of the matrix $T$, we conclude that $T$ is a strictly diagonally dominant $M$-matrix. In fact, as $d_+(x, t) \geq 0$ and $d_-(x, t) \geq 0$, it is easily to show that the matrix $A^{(m)}$, which is defined by (5), is a strictly diagonally dominant $M$-matrix [37] and, hence, all its eigenvalues have positive real parts.

It was shown in [37] that $g_k^{(\alpha)}$ decreases monotonically to zero as $k$ tends to infinity with the rate of $\alpha + 1$.

**Lemma 4.3.** [37] *Let $g_k^{(\alpha)}$ be defined by (2) with $1 < \alpha < 2$. Then*

$$g_k^{(\alpha)} = \frac{1}{\Gamma(-\alpha) k^{\alpha+1}} \left(1 + O\left(\frac{1}{k}\right)\right),$$

*where $\Gamma(x)$ is the Gamma function.*

Therefore, the matrix $T$ has the off-diagonal decay property (16), that is, $T \in \mathcal{L}_{\alpha+1}$. Thus, the matrices $A$ and $K_i$ also have the off-diagonal decay property, that is, $A \in \mathcal{L}_{\alpha+1}$ and $K_i \in \mathcal{L}_{\alpha+1}$

As $1 < \alpha < 2$, by Lemma 4.1, we have the following result.

**Lemma 4.4.** *Let $T$ be defined by* (6) *with $1 < \alpha < 2$. Assume $d(x), w(x) \in C[x_L, x_R]$. Then $T^{-1}, A^{-1}, K_i^{-1} \in \mathcal{L}_{\alpha+1}$, where $A$ and $K_i$ are defined by* (7) *and* (8), *respectively.*

That is to say, there exists a constant $c_0 > 0$ such that

$$|L(i,j)| \leq \frac{c_0}{(1 + |i - j|)^{\alpha+1}}, \tag{17}$$

where $L$ can be $T, A, K_i, T^{-1}, A^{-1}$ and $K_i^{-1}$.

Let $q > 0$ and $\alpha > 0$. Then we have

$$\sum_{x=q+1}^{\infty} \frac{1}{x^{\alpha+1}} = \frac{1}{(q+1)^{\alpha+1}} + \frac{1}{(q+2)^{\alpha+1}} + \cdots \leq \frac{1}{\alpha q^{\alpha}}. \tag{18}$$

This inequality will be used to estimate the norm of a matrix having the off-diagonal decay property (16).

In the following, we discuss the spectral properties of the preconditioned matrix $B_2^{-1}A$, or the approximation property of $B_2^{-1}$ to $A^{-1}$. As

$$B_2^{-1} - A^{-1} = (B_2^{-1} - B_1^{-1}) + (B_1^{-1} - A^{-1}),$$

we will investigate the properties of $B_2^{-1} - B_1^{-1}$ and $B_1^{-1} - A^{-1}$, respectively.

First, we consider the approximation property of $B_1^{-1}$ to $A^{-1}$.

**Lemma 4.5.** *Let $A$ and $K_j$ be defined by* (7) *and* (8), *respectively. Assume $d(x), w(x) \in C^1[x_L, x_R]$. Then for a given $\varepsilon > 0$, there exists a constant $c_1$ and an integer $N_1$ such that for $l \geq N_1$ we have*

$$\|e_i^{\mathsf{T}}(K_i^{-1} - A^{-1})\|_1 \leq c_1 \Delta(x_i, l) + \varepsilon, \tag{19}$$

*where $\Delta(x_i, l) \triangleq \max_{i-l < k < i+l} |x_k - x_i|$.*

*Proof.* It follows from the definitions of $A$ and $K_i$ that

$$\begin{aligned}
&\|e_i^{\mathsf{T}}(K_i^{-1} - A^{-1})\|_1 \\
&= \|e_i^{\mathsf{T}} A^{-1}(A - K_i)K_i^{-1}\|_1 \\
&= \|e_i^{\mathsf{T}} A^{-1}\left((D - d_i I)T + (W - w_i I)T^{\mathsf{T}}\right) K_i^{-1}\|_1 \\
&\leq \left\|e_i^{\mathsf{T}} A^{-1}(D - d_i I)\right\|_1 \|T\|_{\infty} \|K_i^{-1}\|_{\infty} + \left\|e_i^{\mathsf{T}} A^{-1}(W - w_i I)\right\|_1 \|T^{\mathsf{T}}\|_{\infty} \|K_i^{-1}\|_{\infty}. \quad (20)
\end{aligned}$$

As $T$ has the off-diagonal decay property (17), by (18), we have

$$\|T\|_\infty = \max_{1\le i\le N}\sum_{j=1}^N |T(i,j)| = \max_{1\le i\le N}\left(|T(i,i)| + \sum_{j\ne i}|T(i,j)|\right)$$

$$\le \max_{1\le i\le N}\left(c_0 + 2\sum_{k=2}^\infty \frac{c_0}{k^{\alpha+1}}\right)$$

$$\le c_0 + \frac{2c_0}{\alpha} = \frac{(2+\alpha)c_0}{\alpha}.$$

Analogously, we can show that

$$\|K_i^{-1}\|_\infty \le \frac{(2+\alpha)c_0}{\alpha}\quad\text{and}\quad \|A^{-1}\|_\infty \le \frac{(2+\alpha)c_0}{\alpha}.$$

It follows from (18) and Lemma 4.4 that

$$\left\|e_i^\mathsf{T} A^{-1}(D - d_i I)\right\|_1$$

$$= \sum_{j=1}^N |(d_j - d_i)A^{-1}(i,j)|$$

$$= \sum_{j=1}^{i-l} |d_j - d_i|\cdot|A^{-1}(i,j)| + \sum_{j=i-l+1}^{i+l-1} |d_j - d_i|\cdot|A^{-1}(i,j)| + \sum_{j=i+l}^N |d_j - d_i|\cdot|A^{-1}(i,j)|$$

$$\le \max_{1\le j\le N}|d_j - d_i|\sum_{j=1}^{i-l}\frac{c_0}{(1+i-j)^{\alpha+1}} + \max_{i-l<j<i+l}|d_j - d_i|\sum_{j=i-l+1}^{i+l-1}|A^{-1}(i,j)|$$

$$+ \max_{1\le j\le N}|d_j - d_i|\sum_{j=i+l}^N\frac{c_0}{(1+j-i)^{\alpha+1}}$$

$$\le \max_{1\le j\le N}|d_j - d_i|\sum_{k=l+1}^i\frac{c_0}{k^{\alpha+1}} + \|A^{-1}\|_\infty\cdot\max_{i-l<j<i+l}|d_j - d_i| + \max_{1\le j\le N}|d_j - d_i|\sum_{k=l+1}^i\frac{c_0}{k^{\alpha+1}}$$

$$\le \frac{2c_0}{\alpha l^\alpha}\max_{1\le j\le N}|d_j - d_i| + \frac{(2+\alpha)c_0}{\alpha}\max_{i-l<j<i+l}|d_j - d_i|.$$

As $d(x)\in C^1[x_L, x_R]$, we have

$$\max_{1\le j\le N}|d_j - d_i| \le 2\max_{1\le j\le N}|d_j| = 2\max_{1\le j\le N}|d(x_j)| \le 2\max_{x\in[x_L,x_R]}|d(x)|,$$

and

$$\max_{i-l<j<i+l}|d_j - d_i| = \max_{i-l<j<i+l}|d(x_j) - d(x_i)|$$

$$\le \max_{i-l<j<i+l}\max_{x\in[x_L,x_R]}|d'(x)|\cdot|x_j - x_i|$$

$$\le \max_{x\in[x_L,x_R]}|d'(x)|\cdot\Delta(x_i, l).$$

10

Hence,

$$\left\| e_i^\mathsf{T} A^{-1}(D - d_i I) \right\|_1 \leq \frac{4c_0}{\alpha l^\alpha} \max_{x \in [x_L, x_R]} |d(x)| + \frac{(2+\alpha)c_0}{\alpha} \max_{x \in [x_L, x_R]} |d'(x)| \cdot \Delta(x_i, l).$$

For any given $\varepsilon > 0$, let $\tilde{N}_1$ be the integer satisfying

$$\tilde{N}_1^\alpha > \frac{1}{\varepsilon} \cdot \frac{8c_0}{\alpha} \max_{x \in [x_L, x_R]} |d(x)| \cdot \left( \frac{(2+\alpha)c_0}{\alpha} \right)^2.$$

Then for each integer $l \geq \tilde{N}_1$ we have

$$\left\| e_i^\mathsf{T} A^{-1}(D - d_i I) \right\|_1 \|T\|_\infty \|K_i^{-1}\|_\infty \leq \frac{1}{2}\varepsilon + \frac{(2+\alpha)c_0}{\alpha} \max_{x \in [x_L, x_R]} |d'(x)| \cdot \Delta(x_i, l). \quad (21)$$

Analogously, we can prove that there exists an integer $\hat{N}_1$ such that for $l \geq \hat{N}_1$ we have

$$\left\| e_i^\mathsf{T} A^{-1}(W - w_i I) \right\|_1 \|T^\mathsf{T}\|_\infty \|K_i^{-1}\|_\infty \leq \frac{1}{2}\varepsilon + \frac{(2+\alpha)c_0}{\alpha} \max_{x \in [x_L, x_R]} |w'(x)| \cdot \Delta(x_i, l). \quad (22)$$

Let $N_1 = \max\{\tilde{N}_1, \hat{N}_1\}$ and $c_1 = \frac{(2+\alpha)c_0}{\alpha} \left( \max_{x \in [x_L, x_R]} |d'(x)| + \max_{x \in [x_L, x_R]} |w'(x)| \right)$. Combining (20), (21) and (22), we obtain the conclusion. $\square$

It follows from the definition of $B_1$ that

$$\|B_1^{-1} - A^{-1}\|_\infty = \max_{1 \leq i \leq N} \|e_i^\mathsf{T}(B_1^{-1} - A^{-1})\|_1 = \max_{1 \leq i \leq N} \|e_i^\mathsf{T}(K_i^{-1} - A^{-1})\|_1.$$

Note that the constant $c_1$ and the integer $N_1$ in Lemma 4.5 are independent on $i$. Hence, we have the following result.

**Lemma 4.6.** *Let $A$ and $B_1$ be defined by (7) and (9) respectively. Assume $d(x), w(x) \in C^1[x_L, x_R]$. Then for a given $\varepsilon > 0$, there exists a constant $c_1$ and an integer $N_1$ such that for $l \geq N_1$ we have*

$$\|B_1^{-1} - A^{-1}\|_\infty \leq c_1 \max_{1 \leq i \leq N} \Delta(x_i, l) + \varepsilon. \quad (23)$$

We remark that, as $x_i = x_L + ih$ for $i = 0, 1, 2, \ldots, N+1$, it holds that

$$\Delta(x_i, l) = \max_{i-l < k < i+l} |x_k - x_i| = (l-1)h, \quad (24)$$

which will tend to zero as $N$ tends to infinity.

11

## 4.1 The Preconditioner $B_2^{-1}$

Next, we consider the approximation of $B_2^{-1}$ to $B_1^{-1}$. By the definitions of $B_1$ and $B_2$, we have

$$
\begin{aligned}
B_2^{-1} - B_1^{-1} &= \sum_{i=1}^{N} e_i e_i^{\mathsf{T}} (C_i^{-1} - K_i^{-1}) \\
&= \sum_{i=1}^{N} e_i e_i^{\mathsf{T}} K_i^{-1} (K_i - C_i) C_i^{-1} \\
&= d_i \sum_{i=1}^{N} e_i e_i^{\mathsf{T}} K_i^{-1} (T - C) C_i^{-1} + w_i \sum_{i=1}^{N} e_i e_i^{\mathsf{T}} K_i^{-1} (T - C)^{\mathsf{T}} C_i^{-1}.
\end{aligned}
$$

Now we want to show is that $B_2^{-1} - B_1^{-1}$ can be written into a sum of two matrices where one matrix is of small norm and another is of low rank.

We first prove that $\|C_i^{-1}\|_\infty$ is bounded by $\eta^{-1}$.

**Lemma 4.7.** *Let $C_i$ be defined by* (11). *Then we have*

$$
\|C_i^{-1}\|_\infty < \eta^{-1}.
$$

*Proof.* For any given vector $y = [y_1, y_2, \ldots, y_N]^{\mathsf{T}} \in \mathbb{R}^N$, we define

$$
x = [x_1, x_2, \ldots, x_N]^{\mathsf{T}} \triangleq C_i^{-1} y.
$$

Then we have

$$
y = C_i x = \eta x + d_i C x + w_i C^{\mathsf{T}} x.
$$

Let $x_k$ be the entry of $x$ such that

$$
|x_k| = \|x\|_\infty.
$$

Thus, $|x_i| \leq |x_k|$, $i = 1, 2, \ldots, N$. It follows from Lemma 4.2 and (10) that the $k$-th entry of $y$ satisfies

$$
\begin{aligned}
|y_k| &= \left| \eta x_k - (d_i + w_i) g_1^{(\alpha)} x_k + [Cx]_k - d_i g_1^{(\alpha)} x_k + [C^{\mathsf{T}} x]_k - w_i g_1^{(\alpha)} x_k \right| \\
&\geq \left| \eta x_k - (d_i + w_i) g_1^{(\alpha)} x_k \right| - \left| [Cx]_k - d_i g_1^{(\alpha)} x_k \right| - \left| [C^{\mathsf{T}} x]_k - w_i g_1^{(\alpha)} x_k \right| \\
&\geq \eta |x_k| - (d_i + w_i) g_1^{(\alpha)} |x_k| - d_i \sum_{j=1, j \neq k}^{N} g_j^{(\alpha)} |x_k| - w_i \sum_{j=1, j \neq k}^{N} g_j^{(\alpha)} |x_k| \\
&= \eta |x_k| - d_i |x_k| \sum_{j=1}^{N} g_j^{(\alpha)} - w_i |x_k| \sum_{j=1}^{N} g_j^{(\alpha)} \\
&\geq \eta |x_k|.
\end{aligned}
$$

Here $[Cx]_k$ denotes the $k$-th entry of $Cx$. Therefore,

$$\|x\|_\infty = |x_k| \le \frac{1}{\eta}|y_k| \le \frac{1}{\eta}\|y\|_\infty,$$

which holds for any vector $y \in \mathbb{R}^N$. Hence,

$$\|C_i^{-1}\|_\infty = \max_{y \in \mathbb{R}^N, y \ne 0} \frac{\|C_i^{-1}y\|_\infty}{\|y\|_\infty} = \frac{\|x\|_\infty}{\|y\|_\infty} \le \frac{1}{\eta}.$$

$\square$

As $C$ is the Strang circulant approximation of $T$, we have [11]

$$T - C = E_T + S_T,$$

where $E_T$ is of small norm and $S_T$ is of low rank. More precisely, for a given $\varepsilon > 0$ we have $\|E_T\|_\infty < \varepsilon$ and $S_T$ is of the following form

$$S_T = \begin{bmatrix} 0 & 0 & \cdots & 0 & S_1 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & 0 \\ S_2 & 0 & \cdots & 0 & 0 \end{bmatrix},$$

where the dimension of $S_1$ and $S_2$ is dependent on $\varepsilon$ and the decay property of $T$, but is independent on $N$ (the dimension of $A$). Since $K_i^{-1}$ has the off-diagonal decay property (17), we can write $K_i^{-1}$ as a banded matrix plus another one with small norm, that is, $K_i^{-1} = \tilde{K}_i + \hat{K}_i$ where

$$\tilde{K}_i = \begin{bmatrix} * & \cdots & * & 0 & \cdots & 0 \\ \vdots & \ddots & & \ddots & \ddots & \vdots \\ * & & \ddots & & \ddots & 0 \\ 0 & \ddots & & \ddots & & * \\ \vdots & \ddots & \ddots & & \ddots & \vdots \\ 0 & \cdots & 0 & * & \cdots & * \end{bmatrix} \quad \text{and} \quad \hat{K}_i = \begin{bmatrix} 0 & \cdots & 0 & * & \cdots & * \\ \vdots & \ddots & & \ddots & \ddots & \vdots \\ 0 & & \ddots & & \ddots & * \\ * & \ddots & & \ddots & & 0 \\ \vdots & \ddots & \ddots & & \ddots & \vdots \\ * & \cdots & * & 0 & \cdots & 0 \end{bmatrix}.$$

Here we use "$*$" to denote the nonzero entries and $\|\hat{K}_i\|_\infty < \varepsilon$. It follows from (17) and (18) that the bandwidth of $\tilde{K}_i$ is the smallest integer larger than $(\frac{c_0}{\alpha\varepsilon})^{1/\alpha} + 1$. We remark that all matrices $\tilde{K}_i$'s have the same bandwidth for $i = 1, 2, \ldots, N$. It is clear that $\tilde{K}_i$ can be written into a block triangular form

$$\tilde{K}_i = \begin{bmatrix} + & + & & \\ + & \ddots & \ddots & \\ & \ddots & \ddots & + \\ & & + & + \end{bmatrix},$$

13

where "$+$" denotes the nonzero block submatrix whose size is equal to the bandwidth of $\tilde{K}_i$. Therefore,

$$
\begin{aligned}
K_i^{-1}(C-T)C_i^{-1} &= K_i^{-1}E_T C_i^{-1} + (\tilde{K}_i + \hat{K}_i)S_T C_i^{-1} \\
&= (K_i^{-1}E_T + \hat{K}_i S_T)C_i^{-1} + \tilde{K}_i S_T C_i^{-1}.
\end{aligned}
$$

It holds that

$$
\begin{aligned}
\|(K_i^{-1}E_T + \hat{K}_i S_T)C_i^{-1}\|_\infty &\leq \left( \|K_i^{-1}\|_\infty \|E_T\|_\infty + \|\hat{K}_i\|_\infty \|S_T\|_\infty \right) \|C_i^{-1}\|_\infty \\
&\leq \varepsilon \left( \|K_i^{-1}\|_\infty + \|S_T\|_\infty \right) \|C_i^{-1}\|_\infty,
\end{aligned}
$$

and hence

$$
\begin{aligned}
\left\| \sum_{i=1}^N e_i e_i^\intercal \left( (K_i^{-1}E_T + \hat{K}_i S_T)C_i^{-1} \right) \right\|_\infty &= \max_{1 \leq i \leq N} \left\| e_i^\intercal (K_i^{-1}E_T + \hat{K}_i S_T)C_i^{-1} \right\|_1 \\
&\leq \max_{1 \leq i \leq N} \left\| (K_i^{-1}E_T + \hat{K}_i S_T)C_i^{-1} \right\|_\infty \\
&\leq \varepsilon \cdot \max_{1 \leq i \leq N} \left( \|K_i^{-1}\|_\infty + \|S_T\|_\infty \right) \|C_i^{-1}\|_\infty.
\end{aligned}
$$

Note that both $\|K_i^{-1}\|_\infty$ and $\|S_T\|_\infty$ are bounded because of the off-diagonal decay property of $K_i^{-1}$ and $T$. It follows from Lemma 4.7 that $\|C_i^{-1}\|_\infty$ is bounded. Hence, there exists a constant $c_2 > 0$ such that

$$
\left\| \sum_{i=1}^N e_i e_i^\intercal \left( (K_i^{-1}E_T + \hat{K}_i S_T)C_i^{-1} \right) \right\|_\infty < c_2 \varepsilon.
$$

Now we look at the matrix product $\tilde{K}_i S_T C_i^{-1}$. Without loss of generality, we can assume that the dimension of the blocks of $\tilde{K}_i$ and $S_T$ have the same size. Otherwise, we can enlarge the smaller. By direct computations, we can show that $\tilde{K}_i S_T C_i^{-1}$ has the following block structure

$$
\tilde{K}_i S_T C_i^{-1} = \begin{bmatrix}
0 & 0 & \cdots & 0 & + \\
0 & 0 & \cdots & 0 & + \\
0 & 0 & \cdots & 0 & 0 \\
\vdots & \vdots & & \vdots & \vdots \\
0 & 0 & \cdots & 0 & 0 \\
+ & 0 & \cdots & 0 & 0 \\
+ & 0 & \cdots & 0 & 0
\end{bmatrix}
\quad
C_i^{-1} = \begin{bmatrix}
+ & + & \cdots & + & + \\
+ & + & \cdots & + & + \\
0 & 0 & \cdots & 0 & 0 \\
\vdots & \vdots & & \vdots & \vdots \\
0 & 0 & \cdots & 0 & 0 \\
+ & + & \cdots & + & + \\
+ & + & \cdots & + & +
\end{bmatrix}.
$$

Hence,

$$
\operatorname{rank}\left( \sum_{i=1}^N e_i e_i^\intercal \tilde{K}_i S_T C_i^{-1} \right) \leq 4\xi, \tag{25}
$$

14

which is bounded. Here $\xi$ denotes the dimension of the blocks in $\tilde{K}_i S_T C_i^{-1}$, which is independent on $N$ (the dimension of $A$). Therefore,

$$\sum_{i=1}^{N} e_i e_i^\intercal K_i^{-1}(T-C)C_i^{-1} = \sum_{i=1}^{N} e_i e_i^\intercal \left( (K_i^{-1}E_T + \hat{K}_i S_T)C_i^{-1} \right) + \sum_{i=1}^{N} e_i e_i^\intercal \hat{K}_i S_T C_i^{-1}$$

is a sum of a small norm matrix with a low rank matrix.

Analogously, the similar result holds for $\sum_{i=1}^{N} e_i e_i^\intercal K_i^{-1}(T-C)^\intercal C_i^{-1}$. As

$$B_2^{-1} - B_1^{-1} = d_i \sum_{i=1}^{N} e_i e_i^\intercal K_i^{-1}(T-C)C_i^{-1} + w_i \sum_{i=1}^{N} e_i e_i^\intercal K_i^{-1}(T-C)^\intercal C_i^{-1},$$

we have

**Lemma 4.8.** *Let $B_1$ and $B_2$ be defined by* (9) *and* (12) *respectively. Then we have*

$$B_2^{-1} - B_1^{-1} = E_B + S_B,$$

*where $E_B$ is of small norm and $S_B$ is of low rank, that is, $\|E_B\|_\infty < 2c_2\varepsilon$ and $\mathrm{rank}(S_B) \le 4\xi$.*

Therefore, we have the following result.

**Theorem 4.1.** *Let $A$ and $B_2$ be defined by* (7) *and* (12), *respectively. Then, there exists an integer $N_2$ such that for $N > N_2$, we have*

$$B_2^{-1} - A^{-1} = E + S,$$

*where $E$ is of a small norm and $S$ is of a low rank.*

*Proof.* we have

$$
\begin{aligned}
B_2^{-1} - A^{-1} &= (B_2^{-1} - B_1^{-1}) + (B_1^{-1} - A^{-1}) \\
&= E_B + (B_1^{-1} - A^{-1}) + S_B \\
&\triangleq E + S,
\end{aligned}
$$

where $E = E_B + (B_1^{-1} - A^{-1})$ and $S = S_B$. It follows from Lemma 4.6 and (24) that

$$\|B_1^{-1} - A^{-1}\|_\infty \le c_1 \max_{1 \le i \le N} \Delta(x_i, N_1) + \varepsilon = c_1(N_1-1)h + \varepsilon.$$

Let $N_2$ be an integer large enough such that

$$(N_1-1)h = (N_1-1)\frac{x_R - x_L}{N_2+1} < \varepsilon.$$

Then

$$\|E\|_\infty \le \|E_B\|_\infty + \|B_1^{-1} - A^{-1}\|_\infty < (c_1 + 2c_2 + 1)\varepsilon.$$

The conclusion follows. $\qquad\square$

## 4.2 The Preconditioner $B_3^{-1}$

We first establish the difference between $B_3^{-1} - B_2^{-1}$ in terms of interpolation polynomial and interpolation points. We note that

$$
\begin{aligned}
|(B_3^{-1} - B_2^{-1})_{i,j}| &= \left\| e_i e_i^\intercal (B_3^{-1} - B_2^{-1}) e_j e_j^\intercal \right\|_2 \\
&= \left\| e_i e_i^\intercal \left[ \sum_{k=1}^{N} e_k e_k^\intercal \left( \sum_{u=1}^{\ell} \phi_u(x_k) \tilde{C}_u^{-1} \right) - \sum_{k=1}^{N} e_k e_k^\intercal C_k^{-1} \right] e_j e_j^\intercal \right\|_2 \\
&= \left\| e_i e_i^\intercal \left[ \left( \sum_{u=1}^{\ell} \phi_u(x_i) \tilde{C}_u^{-1} \right) - C_i^{-1} \right] e_j e_j^\intercal \right\|_2 \\
&= \left\| e_i e_i^\intercal \left[ \left( \sum_{u=1}^{\ell} F(\phi_u(x_i) \tilde{\Lambda}_u^{-1}) F^* \right) - C_i^{-1} \right] e_j e_j^\intercal \right\|_2 \\
&= \left\| e_i e_i^\intercal \left[ F P_\lambda(x_i) F^* - F Q_\lambda(x_i) F^* \right] e_j e_j^\intercal \right\|_2 \\
&\leq \left\| F P_\lambda(x_i) F^* - F Q_\lambda(x_i) F^* \right\|_2 \\
&\leq \max_k \left\| P_\lambda(x_k) - Q_\lambda(x_k) \right\|_2,
\end{aligned}
$$

where $(B_3^{-1} - B_2^{-1})_{i,j}$ donates the $(i,j)$-th entry of the matrix $B_3^{-1} - B_2^{-1}$, $P_\lambda(x_k) = \mathrm{diag}(p_{\lambda_1}(x_k), p_{\lambda_2}(x_k), \cdots, p_{\lambda_N}(x_k))$ and $Q_\lambda(x_k) = \mathrm{diag}(q_{\lambda_1}(x_k), q_{\lambda_2}(x_k), \cdots, q_{\lambda_N}(x_k))$ are the diagonal matrices, and $\lambda_1, \lambda_2, \cdots, \lambda_N$ are the eigenvalues of $C$. It follows that

$$
|(B_3^{-1} - B_2^{-1})_{i,j}| \leq \max_{1 \leq k \leq N} \| p_\lambda(x_k) - q_\lambda(x_k) \|_2 = \max_{1 \leq k \leq N} \max_{1 \leq u \leq N} \{ |p_{\lambda_u}(x_k) - g_{\lambda_u}(x_k)| \}.
$$

**Theorem 4.2.** *Suppose $\ell$ is sufficiently smaller than $N$. Then $B_3^{-1}$, $B_2^{-1}$ and $B_3^{-1} - B_2^{-1}$ can be expressed as $X + Y + Z$, where $X$ has off-diagonal decay property, $Y$ is of a small norm matrix and $Z$ is of a low rank matrix.*

*Proof.* Let $C_x = \eta I + d(x)C + w(x)C^\intercal$. In Section 4.1, we have shown that $C_x^{-1} - K_x^{-1}$ is equal to a sum of a small norm matrix and a low rank matrix, and the low rank matrix is given in the form of (25). As $K_x^{-1}$ has the off-diagonal decay property (see Lemma 4.4), $C_x^{-1}$ can be written as the sum of three matrices $X_x + Y_x + Z_x$, where $X_x$ has the off-diagonal decay property, $Y_x$ is of a small norm matrix and $Z_x$ is of a low rank matrix. The above results can be applied to the points $x = x_i$ and $x = \tilde{x}_j$ used in Section 3. We note from (15) that

$$
B_3^{-1} = \sum_{j=1}^{\ell} \Phi_j \tilde{C}_j^{-1} = \sum_{j=1}^{\ell} \Phi_j \tilde{X}_j + \sum_{j=1}^{\ell} \Phi_j \tilde{Y}_j + \sum_{j=1}^{\ell} \Phi_j \tilde{Z}_j.
$$

We know that for linear interpolation scheme, $|\Phi_i(x)|$ is bounded above by 1 for all $i = 1, 2, \cdots, \ell$ and $x \in [x_L, x_R]$. Therefore, when $\ell$ is sufficiently smaller than $N$, $\sum_{j=1}^{\ell} \Phi_j \tilde{X}_j$ have off-diagonal decay property, $\| \sum_{j=1}^{\ell} \Phi_j \tilde{Y}_j \|_\infty \leq \sum_{j=1}^{\ell} \| \Phi_j \|_\infty \| \tilde{Y}_j \|_\infty$ is small and $\sum_{j=1}^{\ell} \Phi_j \tilde{Z}_j$ is a low rank matrix given in the form of (25).

On the other hand, in Lemma 4.8, we have shown that $B_2^{-1} - B_1^{-1}$ is equal to the sum of a small norm matrix and a matrix of low rank. In Lemma 4.4, we know that $B_1^{-1}$ has the off-diagonal property. Also it is clear that $B_3^{-1} - B_2^{-1}$ can be expressed as the sum of a matrix with off-diagonal decay property, a small norm matrix and a low rank matrix. $\square$

**Theorem 4.3.** *Let* $\theta = \max_{x \in [x_L, x_R]} \max_{1 \le u \le N} \{p_{\lambda_u}(x) - g_{\lambda_u}(x)\}$. *Suppose $\ell$ is sufficiently smaller than $N$. Then for each $\epsilon > 0$, there exists an integer $N_3 > 0$ (which is independent of $N$) such that $B_3^{-1} - B_2^{-1} = E + S$, where $E$ is a matrix with $\|E\|_\infty \le \theta(2N_3 + 1) + \epsilon$. and $S$ is a low rank matrix.*

*Proof.* According to Theorem 4.2, we know that $B_3^{-1} - B_2^{-1} = X + Y + Z$, where $X$ has off-diagonal decay property, $\|Y\|_\infty \le \frac{\epsilon}{2}$ and $Z$ is a low rank matrix given in the form of (25). Now we can be written as $B_3^{-1} - B_2^{-1} = \hat{X} + \hat{Y} + \hat{Z}$ with the property that

$$|(B_3^{-1} - B_2^{-1})_{i,j}| = |\hat{X}_{i,j}| + |\hat{Y}_{i,j}| + |\hat{Z}_{i,j}|, \tag{26}$$

where $\tilde{X}$ has off-diagonal decay property, $\|\hat{Y}\|_\infty \le \frac{\epsilon}{2}$ and $\hat{Z}$ is a low rank matrix given in the form of (25). According to the structure of $Z$ (see the matrix structure in (25)), when $Z_{i,j} \ne 0$, we set $\hat{X}_{i,j} = \hat{Y}_{i,j} = 0$ and $\hat{Z}_{i,j} = (B_3^{-1} - B_2^{-1})_{i,j}$; when $Z_{i,j} \ne 0$, we set $\hat{Y}_{i,j} = 0$ and $\hat{X}_{i,j} = X_{i,j} + Y_{i,j}$ for $X_{i,j} Y_{i,j} < 0$ or we set $\hat{X}_{i,j} = X_{i,j}$ and $\hat{Y}_{i,j} = Y_{i,j}$ for $X_{i,j} Y_{i,j} \ge 0$. It is clear that (26) satisfies.

Next we would like to show that $\hat{X} + \hat{Y}$ can be a small norm matrix. W first note that for $\hat{X}$, there exists $N_3 > 0$ such that

$$\sum_{|j-i|>N_3} |\hat{X}_{i,j}| \le \frac{\epsilon}{2}, \quad i = 1, 2, \cdots, N.$$

It follows that

$$
\begin{aligned}
\|\hat{X}\|_\infty &= \max_i \|e_i^T \hat{X}\|_1 \\
&\le \max_i \sum_{j-i \le N_3} |\hat{X}_{i,j}| + \frac{\epsilon}{2} \\
&= \max_i \sum_{j-i \le N_3} |(B_3^{-1} - B_2^{-1})_{i,j}| - |\hat{Y}_{i,j}| - |\hat{Z}_{i,j}| + \frac{\epsilon}{2} \\
&\le \max_i \sum_{j-i \le N_3} |(B_3^{-1} - B_2^{-1})_{i,j}| + \frac{\epsilon}{2} \\
&\le (2N_3 + 1) \max_{i,j} |(B_3^{-1} - B_2^{-1})_{i,j}| + \frac{\epsilon}{2} \\
&\le (2N_3 + 1) \max_k \max_u \{p_{\lambda_u}(x_k) - q_{\lambda_u}(x_k)\} + \frac{\epsilon}{2} \\
&\le (2N_3 + 1)\theta + \frac{\epsilon}{2}.
\end{aligned}
$$

Let $E = \hat{X} + \hat{Y}$, we have $\|E\|_\infty \le \|\hat{X}\|_\infty + \|\hat{Y}\|_\infty \le (2N_3 + 1)\theta + \epsilon$. The result follows. $\square$

By combining the results in Theorems 4.1 and 4.3, we conclude that $B_3^{-1} - A^{-1}$ can be written as the sum of a small norm matrix and a low rank matrix. In the next section, we present numerical examples to demonstrate the usefulness of the proposed preconditioner.

# 5 Numerical Examples

In this section, we carry out numerical experiments to study the performance of our new preconditioner $B_3$. We employ the preconditioned GMRES method to solve the linear system (4). In all numerical experiments, the stopping criterion is

$$\frac{\|r_k\|_2}{\|r_0\|_2} < 10^{-7},$$

where $r_k$ is the residual vector after $k$ iterations and $r_0$ is the initial residual vector.

**Example 1.** We first consider the FDE (1) tested in [26, 37] where $\alpha = 1.8$, $[x_L, x_R] = [0, 2]$ and $T_f = 1$. The diffusion coefficients are given by

$$d_+(x) = \Gamma(1.2)\, x^\alpha, \quad d_-(x) = \Gamma(1.2)\,(2-x)^\alpha,$$

and the source term is

$$f(x,t) = -32\,\mathrm{e}^{-t}\left(x^2 + \frac{1}{8}(2-x)^2(8+x^2) - \frac{5}{2}[x^3 + (2-x)^3] + \frac{25}{22}[x^4 + (2-x)^4]\right).$$

The initial condition is chosen as

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

Then the true solution to the corresponding FDE is given by [22, 37, 26]

$$u(x,t) = 4\,\mathrm{e}^{-t}x^2(2-x)^2.$$

The numerical results are listed in Table 1, where "GMRES" denotes the GMRES method without preconditioner, "$B_3(\ell)$-GMRES" denotes the GMRES method with the preconditioner $B_3$ with $\ell$ being the number of interpolation points, and "$C$-GMRES" denote the GMRES method with the following circulant preconditioner

$$\eta I + \tilde{d}C + \tilde{w}C^\mathsf{T},$$

where $\tilde{d}$ and $\tilde{w}$ are the mean values of the diagonals of the diagonal matrices $D$ and $W$, respectively. As for comparisons, we also carry out the Gaussian elimination method, which is denoted by "GE".

In Table 1, "$N$" denotes the number of spatial grid points, "$M$" denotes the number of time steps, "Iter" denotes the average number of iterations required to solve (4) at each time step, and "CPU" denotes the total CPU time in seconds for solving the whole discretized system.

Since the true solution is available, we also report the error between the true solution and the approximation under the infinity norm at the last time step. It is denoted by "Error" in Table 1. As this value is almost same for all GMRES methods, we just list the error for GE and $B_3(4)$-GMRES.

Table 1: Numerical results for Example 1.

| $N$ | $M$ | GE | | GMRES | | $C$-GMRES | | $B_3(2)$-GMRES | | $B_3(4)$-GMRES | | |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| | | CPU | Error | Iter | CPU | Iter | CPU | Iter | CPU | Iter | CPU | Error |
| $2^7$ | $2^6$ | 0.13 | 8.3059e-3 | 65.98 | 1.41 | 8.02 | 0.12 | 6.00 | 0.10 | 5.00 | 0.09 | 8.3059e-3 |
| $2^8$ | $2^7$ | 0.99 | 4.0727e-3 | 115.02 | 7.54 | 7.05 | 0.23 | 5.00 | 0.19 | 4.00 | 0.17 | 4.0727e-3 |
| $2^9$ | $2^8$ | 9.25 | 2.0159e-3 | 159.02 | 31.22 | 7.00 | 0.59 | 4.00 | 0.39 | 4.00 | 0.45 | 2.0159e-3 |
| $2^{10}$ | $2^9$ | 90.26 | 1.0028e-3 | 196.55 | 117.47 | 6.00 | 1.44 | 3.00 | 0.88 | 3.00 | 1.02 | 1.0027e-3 |
| $2^{11}$ | $2^{10}$ | 1057.99 | 5.0009e-4 | 225.01 | 449.25 | 5.00 | 3.93 | 3.00 | 2.77 | 3.00 | 3.28 | 5.0008e-4 |

We see that the preconditioned GMRES methods exhibit excellent performance both in terms of iteration steps and CPU time, and the iteration number does not increase as the number of the spatial grid points increases. For this example, $\ell = 2$ is good enough, which means that we only need to choose two interpolation points. The performance of the proposed preconditioner is better than that of the circulant preconditioner by taking the average of the coefficient values in the fractional diffusion equations.

**Example 2.** This example is a modification of Example 1. We replace the right-sided diffusion coefficient $d_-(x)$ with

$$d_-(x) = \Gamma(1.2) \, (2 - x)^{(1+\alpha)}.$$

Other data are the same as that in Example 1.

The numerical results are listed in Table 2. We can see that the performance is improved significantly as the number of interpolation points increasing. Again, the performance of the proposed preconditioner is better than that of the circulant preconditioner.

Table 2: Numerical results for Example 2.

| $N$ | $M$ | GE | GMRES | | $C$-GMRES | | $B_3(2)$-GMRES | | $B_3(4)$-GMRES | | $B_3(6)$-GMRES | |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| | | CPU | Iter | CPU | Iter | CPU | Iter | CPU | Iter | CPU | Iter | CPU |
| $2^7$ | $2^6$ | 0.15 | 81.59 | 1.98 | 115.31 | 3.80 | 51.42 | 1.14 | 12.58 | 0.22 | 8.86 | 0.17 |
| $2^8$ | $2^7$ | 0.91 | 114.16 | 7.63 | 176.84 | 17.49 | 65.23 | 3.64 | 12.54 | 0.50 | 8.48 | 0.37 |
| $2^9$ | $2^8$ | 9.23 | 146.46 | 27.65 | 239.08 | 70.08 | 80.52 | 12.00 | 15.30 | 1.59 | 9.18 | 1.05 |
| $2^{10}$ | $2^9$ | 88.45 | 172.40 | 96.10 | 295.32 | 267.78 | 97.40 | 42.60 | 18.77 | 5.62 | 10.83 | 3.59 |
| $2^{11}$ | $2^{10}$ | 1062.40 | 186.62 | 335.89 | 347.31 | 1195.35 | 119.48 | 178.55 | 21.29 | 21.05 | 12.12 | 13.26 |

**Example 3.** This example comes from [26], which is a modification of an example in [37]. In this example, we consider the FDE with an anomalously diffused Gaussian pulse. The initial condition is given by

$$u(x, 0) = x^2(2 - x^2) \exp\left( -\frac{(x - x_c)^2}{2\sigma^2} \right),$$

with the mean $x_c = 1.2$ and the standard deviation $\sigma = 0.08$. The diffusion coefficients are

$$d_+(x) = \delta(1 + x^2 + t^2), \quad d_-(x) = \delta(1 + (2-x)^2 + t^2),$$

which are dependent on $x$ and $t$. Here, $\delta$ is a parameter. We will test the problem for different values of $\delta$. Other data are as follows:

$$\alpha = 1.5, \quad [x_L, x_R] = [0, 2], \quad T_f = 1, \quad f(x, t) = 0.$$

The numerical results are reported in Table 3. In the table, we also list the iteration number of the first time step. Here we use the initial condition as the initial guess in the first time step, and find that it takes more iteration numbers to converge. The performance of the proposed preconditioner is still better than the circulant preconditioner.

Table 3: Numerical results for Example 3. The number inside the bracket refers to the number of iterations required at the first time step.

| $\delta$ | $N$ | $M$ | GE | GMRES | | $C$-GMRES | | $B_3(2)$-GMRES | | $B_3(4)$-GMRES | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | CPU | Iter | CPU | Iter | CPU | Iter | CPU | Iter | CPU |
| 1 | $2^7$ | $2^6$ | 0.17 | 59.66 (78) | 1.24 | 11.42 (14) | 0.19 | 8.45 (12) | 0.15 | 6.08 (7) | 0.12 |
| | $2^8$ | $2^7$ | 1.00 | 79.52 (111) | 4.11 | 10.70 (13) | 0.36 | 7.27 (11) | 0.28 | 5.90 (7) | 0.27 |
| | $2^9$ | $2^8$ | 9.65 | 93.46 (145) | 12.84 | 10.09 (13) | 0.87 | 6.14 (11) | 0.61 | 5.17 (7) | 0.62 |
| | $2^{10}$ | $2^9$ | 90.82 | 99.79 (183) | 37.70 | 9.24 (12) | 2.22 | 5.59 (10) | 1.57 | 4.95 (6) | 1.70 |
| | $2^{11}$ | $2^{10}$ | 1072.44 | 98.26 (219) | 112.42 | 8.21 (11) | 6.31 | 4.85 (8) | 4.39 | 4.68 (6) | 5.21 |
| 10 | $2^7$ | $2^6$ | 0.13 | 67.58 (99) | 1.45 | 15.41 (18) | 0.23 | 11.22 (14) | 0.19 | 7.52 (9) | 0.14 |
| | $2^8$ | $2^7$ | 1.00 | 121.16 (173) | 8.27 | 15.48 (18) | 0.53 | 10.31 (14) | 0.39 | 6.61 (9) | 0.29 |
| | $2^9$ | $2^8$ | 9.41 | 210.39 (306) | 51.43 | 15.04 (18) | 1.30 | 9.29 (14) | 0.89 | 6.25 (9) | 0.72 |
| | $2^{10}$ | $2^9$ | 96.25 | 331.91 (512) | 306.73 | 14.02 (17) | 3.36 | 8.17 (14) | 2.18 | 5.32 (9) | 1.80 |
| | $2^{11}$ | $2^{10}$ | 1072.33 | 442.54 (792) | 1827.28 | 12.77 (16) | 10.01 | 6.91 (13) | 5.90 | 5.19 (8) | 5.63 |
| 100 | $2^7$ | $2^6$ | 0.14 | 69.70 (105) | 1.53 | 17.72 (20) | 0.27 | 13.05 (16) | 0.22 | 8.75 (11) | 0.16 |
| | $2^8$ | $2^7$ | 1.02 | 127.02 (189) | 8.98 | 18.60 (21) | 0.65 | 12.62 (16) | 0.48 | 8.41 (11) | 0.36 |
| | $2^9$ | $2^8$ | 9.36 | 231.53 (343) | 61.27 | 18.82 (21) | 1.67 | 12.24 (16) | 1.16 | 8.07 (11) | 0.89 |
| | $2^{10}$ | $2^9$ | 94.26 | 417.55 (625) | 472.63 | 18.61 (21) | 4.57 | 11.62 (16) | 3.03 | 7.63 (11) | 2.42 |
| | $2^{11}$ | $2^{10}$ | 1072.72 | 742.34 (1123) | 5330.16 | 17.91 (21) | 14.28 | 10.78 (16) | 8.87 | 7.41 (10) | 7.55 |

In this example, We use the initial condition as the initial guess in the first time step, and find that it takes more iteration numbers to converge. In Table **??**, we list the iteration number of the first time step. The performance of the proposed preconditioner is still better than the circulant preconditioner.

**Example 4.** We consider the two dimensional fractional diffusion equation

$$\frac{\partial u(x,y,t)}{\partial t} = d_+(x,y,t)\frac{\partial^\alpha u(x,y,t)}{\partial_+ x^\alpha} + d_-(x,y,t)\frac{\partial^\alpha u(x,y,t)}{\partial_- x^\alpha}$$

$$+ e_+(x,y,t)\frac{\partial^\beta u(x,y,t)}{\partial_+ y^\beta} + e_-(x,y,t)\frac{\partial^\beta u(x,y,t)}{\partial_- y^\beta} + f(x,y,t),$$

$$(x,y) \in \Omega \triangleq [x_L, x_R] \times [y_L, y_R], \quad t \in (0, T_f],$$

$$u(x,y,t) = 0 \text{ for } (x,y) \in \partial\Omega \text{ and } 0 \le t \le T_f,$$

$$u(x,y,0) = u_0(x,y) \text{ for } (x,y) \in \overline{\Omega}.$$

In the numerical tests, we choose $\alpha = \beta = 1.2$, $[x_L, x_R] = [y_L, y_R] = [0,2]$, $T_f = 1$, $f(x,y,t) = 0$ and

$$d_+(x,y,t) = d_-(x,y,t) = d_{x,t}d_{y,t} = \mathrm{e}^{4t}x^{4\alpha}y^{4\beta}, \quad e_+(x,y,t) = e_-(x,y,t) = e_{x,t}e_{y,t} = \mathrm{e}^{4t}(2-x)^{4\alpha}(2-y)^{4\beta}.$$

The initial condition is given by

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

Assume the number of spatial discretization points in $x$-direction and $y$-direction are $N_x$ and $N_y$ respectively. The coefficient matrix in this example is given by

$$A = \eta I + \begin{bmatrix} D_{x,1} & & & \\ & D_{x,2} & & \\ & & \ddots & \\ & & & D_{x,N_y} \end{bmatrix}(I \otimes T_\alpha) + \begin{bmatrix} W_{x,1} & & & \\ & W_{x,2} & & \\ & & \ddots & \\ & & & W_{x,N_y} \end{bmatrix}(I \otimes T_\alpha^\intercal)$$

$$+ \begin{bmatrix} D_{y,1} & & & \\ & D_{y,2} & & \\ & & \ddots & \\ & & & D_{y,N_x} \end{bmatrix}(T_\beta \otimes I) + \begin{bmatrix} W_{y,1} & & & \\ & W_{y,2} & & \\ & & \ddots & \\ & & & W_{y,N_x} \end{bmatrix}(T_\beta^\intercal \otimes I),$$

where $T_\alpha$ and $T_\beta$ are the discretization matrices for the fractional orders $\alpha$ and $\beta$, respectively (similar to $T$ defined in (6)). The diagonals of diagonal matrices $D_{x,j}$, $W_{x,j}$, $D_{y,j}$ and $W_{y,j}$ are defined by

$$(D_{x,j})_{i,i} = d_+(x_i, y_j), \quad (W_{x,j})_{i,i} = d_-(x_i, y_j), \quad 1 \le i \le N_x, \ 1 \le j \le N_y,$$

$$(D_{i,y})_{j,j} = e_+(x_i, y_j), \quad (W_{i,y})_{j,j} = e_-(x_i, y_j), \quad 1 \le i \le N_x, \ 1 \le j \le N_y.$$

We consider the sample points to be $\tilde{x}_1, \tilde{x}_2, \cdots, \tilde{x}_{\ell_x}$ and $\tilde{y}_1, \tilde{y}_2, \cdots, \tilde{y}_{\ell_y}$. The corresponding matrices are

$$\tilde{C}_{u,v} = \eta I + d_+(\tilde{x}_u, \tilde{y}_v)(I \otimes T_\alpha) + d_-(\tilde{x}_u, \tilde{y}_v)(I \otimes T_\alpha^\intercal) + e_+(\tilde{x}_u, \tilde{y}_v)(T_\beta \otimes I) + e_-(\tilde{x}_u, \tilde{y}_v)(T_\beta^\intercal \otimes I),$$

for $1 \leq i \leq \ell_x$ and $1 \leq j \leq \ell_y$. Then we obtain the preconditioner

$$
\begin{aligned}
B_3^{-1} &= \sum_{i=1}^{N_x}\sum_{j=1}^{N_y}(e_i \otimes e_j)(e_i \otimes e_j)^\intercal \left(\sum_{u=1}^{\ell_x}\sum_{v=1}^{\ell_y}\phi_{u,v}(x_i,y_j)\tilde{C}_{u,v}^{-1}\right) \\
&= \sum_{i=1}^{N_x}\sum_{j=1}^{N_y}(e_i \otimes e_j)(e_i \otimes e_j)^\intercal F \left(\sum_{u=1}^{\ell_x}\sum_{v=1}^{\ell_y}\phi_{u,v}(x_i,y_j)\tilde{\Lambda}_{u,v}^{-1}\right) F^* \\
&= \sum_{i=1}^{N_x}\sum_{j=1}^{N_y}(e_i \otimes e_j)(e_i \otimes e_j)^\intercal \sum_{u=1}^{\ell_x}\sum_{v=1}^{\ell_y} F\left(\phi_{u,v}(x_i,y_j)\tilde{\Lambda}_{u,v}^{-1}\right) F^* \\
&= \sum_{u=1}^{\ell_x}\sum_{v=1}^{\ell_y}\left(\sum_{i=1}^{N_x}\sum_{j=1}^{N_y}(e_i \otimes e_j)(e_i \otimes e_j)^\intercal \phi_{u,v}(x_i,y_j)\right) F\tilde{\Lambda}_{u,v}^{-1} F^* \\
&= \sum_{u=1}^{\ell_x}\sum_{v=1}^{\ell_y}\Phi_{u,v} F\tilde{\Lambda}_{u,v}^{-1} F^*.
\end{aligned}
$$

Here $F$ is the the two dimensional discrete Fourier transform matrix of size $N_x N_y$, $\Phi_{u,v} = \mathrm{diag}\left(\phi_{u,v}(x_1,y_1), \phi_{u,v}(x_2,y_1), \ldots, \phi_{u,v}(x_N,y_N)\right)$ are diagonal matrices. Now applying $B_3^{-1}$ to any vector requires about $O(\ell_x \ell_y N_x N_y \log(N_x N_y))$ operations which is acceptable for a moderate number $\ell_x$ and $\ell_y$.

The results are reported in Table 4, where $M$ denotes the number of time steps. Here we set $\ell = \ell_x = \ell_y$ for the proposed preconditioner. In the experiment, we test $\ell = 2, 3, 4$. As the Gaussian Elimination is more expensive for this example, we do not report its results. In the table, "$-$" means that the methods do not converge on at least one of the time steps within 1000 iterations in the iterative solver. We can see that GMRES(C) does not work well for this example. But our preconditioners still exhibit excellent performance.

Table 4: Numerical results for Example 4

| $N_x$ $N_y$ $M$ | GMRES | | $C$-GMRES | | $B_3(2)$-GMRES | | $B_3(3)$-GMRES | | $B_3(4)$-GMRES | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Iter | CPU | Iter | CPU | Iter | CPU | Iter | CPU | Iter | CPU |
| $2^4$ $2^4$ $2^4$ | 149.38 | 1.83 | 160.88 | 2.40 | 138.88 | 2.52 | 68.31 | 1.45 | 58.19 | 1.74 |
| $2^5$ $2^5$ $2^5$ | 284.59 | 17.16 | 345.34 | 28.61 | 248.91 | 18.56 | 89.47 | 5.83 | 69.09 | 5.99 |
| $2^6$ $2^6$ $2^6$ | 393.42 | 615.73 | $-$ | $-$ | 330.44 | 455.23 | 108.16 | 70.63 | 73.17 | 46.26 |
| $2^7$ $2^7$ $2^7$ | 442.41 | 6205.37 | $-$ | $-$ | 374.66 | 4820.09 | 115.77 | 773.45 | 62.59 | 384.41 |

# 6   Concluding Remarks

In this paper, we have considered discretized linear systems arising from fractional diffusion equations. The matrix structure of their coefficient martrices is the sum of a scaled identity matrix and two diagonal-times-Toeplitz matrices. Preconditioning techniques for

such Toeplitz-like matrices have not been studied and developed. The main contribution of this paper is to develop approximate inverse preconditioners for these Toeplitz-like matrices so that Krylov subspace methods for solving these preconditioned systems can converge very quickly. We have presented numerical examples and have shown that the proposed preconditioning technique is very effective and efficient.

# References

[1] J. Bai and X. Feng, *Fractional-order anisotropic diffusion for image denoising*, IEEE Tran. Image Proc., 16 (2007), pp. 2492–2502.

[2] R. Barrett, M. Berry, T. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. van der Vorst, *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*, SIAM, Philadelphia, 1994.

[3] D. Benson, S. W. Wheatcraft, and M. M. Meerschaert, *Application of a fractional advection-dispersion equation*, Water Resour. Res., 36 (2000), pp. 1403–1413.

[4] D. Benson, S. W. Wheatcraft, and M. M. Meerschaert, *The fractional-order governing equation of Lévy motion*, Water Resour. Res., 36 (2000), pp. 1413–1423.

[5] M. Benzi, *Preconditioning techniques for large linear systems: A survey*, J. Comput. Phys., 182 (2002), pp. 418–477.

[6] B. Beumer, M. Kovács, and M. M. Meerschaert, *Numerical solutions for fractional reaction-diffusion equations*, Comput. Math. Appl., 55 (2008), pp. 2212–2226.

[7] B. A. Carreras, V. E. Lynch, and G. M. Zaslavsky, *Anomalous diffusion and exit time distribution of particle tracers in plasma turbulence models*, Phys. Plasma, 8 (2001), pp. 5096–5103.

[8] R. Chan, Q. Chang, and H. Sun, *Multigrid method for ill-conditioned symmetric Toeplitz systems*, SIAM J. Sci. Comput., 19 (1998), pp. 516–529.

[9] R. Chan and X. Jin, *An Introduction to Iterative Toeplitz Solvers*, SIAM, Philadelphia, 2007.

[10] R. Chan and M. Ng, *Conjugate gradient methods for Toeplitz systems*, SIAM Rev., 38 (1996), pp. 427–482.

[11] R. Chan and G. Strang, *Toeplitz equations by conjugate gradients with circulant pre-conditioner*, SIAM J. Sci. Statist. Comput., 10 (1989), pp. 104–119.

[12] T. Chan, *An optimal circulant preconditioner for Toeplitz systems*, SIAM J. Sci. Statist. Comput., 9 (1988), pp. 766–771.

[13] W. Deng, *Finite element method for the space and time fractional Fokker-Planck equation*, SIAM J. Numer. Anal., 47 (2008), pp. 204–226.

[14] V. J. Ervin, N. Heuer, and J. P. Roop, *Numerical approximation of a time dependent, nonlinear, space-fractional diffusion equation*, SIAM J. Numer. Anal., 45 (2007), pp. 572–591.

[15] S. Jaffard, *Propriétés des matrices "bien localisées" près de leur diagonale et quelques applications*, Ann. Inst. H. Poincaré Anal. Non Linéaire, 7 (1990), pp. 461–476.

[16] T. A. M. Langlands and B. I. Henry, *The accuracy and stability of an implicit solution method for the fractional diffusion equation*, J. Comput. Phys., 205 (2005), pp. 719–736.

[17] S. Lei and H. Sun, *A circulant preconditioner for fractional diffusion equations*, J. Comput. Phys., 242 (2013), pp. 715–725.

[18] F. Liu, V. Anh, and I. Turner, *Numerical solution of the space fractional Fokker-Planck equation*, J. Comput. Appl. Math., 166 (2004), pp. 209–219.

[19] R. L. Magin, *Fractional Calculus in Bioengineering*, Begell House Publishers, 2006.

[20] M. M. Meerschaert, H. P. Scheffler, and C. Tadjeran, *Finite difference methods for two-dimensional fractional dispersion equation*, J. Comput. Phys., 211 (2006), pp. 249–261.

[21] M. M. Meerschaert and C. Tadjeran, *Finite difference approximations for fractional advection-dispersion flow equations*, J. Comput. Appl. Math., 172 (2004), pp. 65–77.

[22] M. M. Meerschaert and C. Tadjeran, *Finite difference approximations for two-sided space-fractional partial differential equations*, Appl. Numer. Math., 56 (2006), pp. 80–90.

[23] D. A. Murio, *Implicit finite difference approximation for time fractional diffusion equations*, Comput. Math. Appl., 56 (2008), pp. 1138–1145.

[24] M. Ng, *Iterative Methods for Toeplitz Systems*, Numerical Mathematics and Scientific Computation, Oxford University Press, New York, 2004.

[25] M. Ng and J. Pan, *Approximate inverse circulant-plus-diagonal preconditioners for Toeplitz-plus-diagonal matrices*, SIAM J. Sci. Comput., 32 (2010), pp. 1442–1464.

[26] H. Pang and H. Sun, *Multigrid method for fractional diffusion equations*, J. Comput. Phys., 231 (2012), pp. 693–703.

[27] I. Podlubny, *Fractional Differential Equations*, Academic Press, New York, 1999.

[28] M. Raberto, E. Scalas, and F. Mainardi, *Waiting-times and returns in high-frequency financial data: an empirical study*, Physica A, 314 (2002), pp. 749–755.

[29] M. F. Shlesinger, B. J. West, and J. Klafter, *Lévy dynamics of enhanced diffusion: application to turbulence*, Phys. Rev. Lett., 58 (1987), pp. 1100–1103.

[30] I. M. Sokolov, J. Klafter, and A. Blumen, *Fractional kinetics*, Phys. Today Nov., (2002), pp. 28–53.

[31] E. Sousa, *Finite difference approximates for a fractional advection diffusion problem*, J. Comput. Phys., 228 (2009), pp. 4038–4054.

[32] T. Strohmer, *Four short stories about Toeplitz matrix calculations*, Linear Algebra Appl., 343-344 (2002), pp. 321–344.

[33] L. Su, W. Wang, and Z. Yang, *Finite difference approximations for the fractional advection-diffusion equation*, Phys. Lett. A, 373 (2009), pp. 4405–4408.

[34] H. Sun, R. Chan, and Q. Chang, *A note on the convergence of the two-grid method for Toeplitz systems*, Comput. Math. Appl., 34 (1997), pp. 11–18.

[35] C. Tadjeran, M. M. Meerschaert, and H. P. Scheffler, *A second-order accurate numerical approximation for the fractional diffusion equation*, J. Comput. Phys., 213 (2006), pp. 205–213.

[36] H. Wang and K. Wang, *An $\mathcal{O}(N \log^2 N)$ alternating-direction finite difference method for two-dimensional fractional diffusion equations*, J. Comput. Phys., 230 (2011), pp. 7830–7839.

[37] H. Wang, K. Wang, and T. Sircar, *A direct $O(N \log^2 N)$ finite difference method for fractional diffusion equations*, J. Comput. Phys., 229 (2010), pp. 8095–8104.

[38] K. Wang and H. Wang, *A fast characteristic finite difference method for fractional advection-diffusion equations*, Adv. Water Resour., 34 (2011), pp. 810–816.

[39] K. Wang and D. Yang, Wellposedness of variable-coefficient conservative fractional elliptic differential equations, SIAM J. Numer. Anal., 51 (2013), pp. 1088-1107.

[40] G. M. Zaslavsky, D. Stevens, and H. Weitzner, *Self-similar transport in incomplete chaos*, Phys. Rev. E, 48 (1993), pp. 1683–1694.