# A STABLE FAST TIME-STEPPING METHOD FOR FRACTIONAL INTEGRAL AND DERIVATIVE OPERATORS [*]

FANHAI ZENG[†], IAN TURNER[†,‡], AND KEVIN BURRAGE[†,§]

**Abstract.** A unified fast time-stepping method for both fractional integral and derivative operators is proposed. The fractional operator is decomposed into a local part with memory length $\Delta T$ and a history part, where the local part is approximated by the direct convolution method and the history part is approximated by a fast memory-saving method. The fast method has $O(n_0 + \sum_\ell^L q_\alpha(N_\ell))$ active memory and $O(n_0 n_T + (n_T - n_0) \sum_\ell^L q_\alpha(N_\ell))$ operations, where $L = \log(n_T - n_0)$, $n_0 = \Delta T/\tau$, $n_T = T/\tau$, $\tau$ is the stepsize, $T$ is the final time, and $q_\alpha(N_\ell)$ is the number of quadrature points used in the truncated Laguerre–Gauss (LG) quadrature. The error bound of the present fast method is analyzed. It is shown that the error from the truncated LG quadrature is independent of the stepsize, and can be made arbitrarily small by choosing suitable parameters that are given explicitly. Numerical examples are presented to verify the effectiveness of the current fast method.

**Key words.** Fast convolution, the (truncated) Laguerre–Gauss quadrature, short memory principle, fractional differential equations, fractional Lorenz system.

**AMS subject classifications.** 26A33, 65M06, 65M12, 65M15, 35R11

## 1. Introduction.

The convolution of the form

$$(1.1) \qquad (k * u)(t) = \int_0^t k(t - s)u(s)\,\mathrm{d}s$$

arises in many physical models, such as integral equations, integrodifferential equations, fractional differential equations, and integer-order differential equations such as wave propagation with nonreflecting boundary conditions, see for example, [3, 8, 28, 29, 35, 22].

The aim of this paper is to present a stable and fast memory-saving time-stepping algorithm for the convolution (1.1) with a kernel $k(t) = t^{\alpha-1}/\Gamma(\alpha)$. This kind of kernel has found wide applications in science and engineering [8, 28, 29]. When $\alpha \geq 0$, Eq. (1.1) gives a fractional integral of order $\alpha$. If $\alpha < 0$, then Eq. (1.1) can be interpreted as the Hadamard finite part integral, which is equivalent to the Riemann–Liouville (RL) fractional derivative of order $-\alpha$, see [30, p. 112].

The direct discretization of (1.1) takes the following form

$$(1.2) \qquad \sum_{k=0}^n \omega_{n,k} u(t_k), \quad n = 1, 2, \ldots, n_T,$$

where $\omega_{n,k}$ are the convolution quadrature weights. The direct computation of (1.2) requires $O(n_T)$ active memory and $O(n_T^2)$ operations, which is expensive for long time computations. The computational difficulty in both memory requirement and computational cost will increase greatly when the direct approximation (1.2) is applied

[†]School of Mathematical Sciences, Queensland University of Technology, Brisbane, QLD 4001, Australia (f2.zeng@qut.edu.au).

[‡]Australian Research Council Centre of Excellence for Mathematical and Statistical Frontiers, Queensland University of Technology, Brisbane, QLD 4001, Australia (i.turner@qut.edu.au)

[§]Visiting Professor, Department of Computer Science, University of Oxford, OXI 3QD, UK (kevin.burrage@qut.edu.au)

to resolve high-dimensional time evolution equations and/or a large system of time-fractional partial differential equations (PDEs) involving (1.1), see e.g., [2, 10, 22, 41, 43]. However, we believe this computational difficulty for fractional operators has not been fully addressed in literature. The short memory principle (see [6, 29]) seems promising to resolve this difficulty, but it has not been widely applied in fractional calculus due to its inaccuracy.

Up to now, some progress has been made in reducing storage requirements and computational cost for resolving fractional models. The basic idea is to seek a suitable sum-of-exponentials to approximate the kernel function $k_\alpha(t) = t^{\alpha-1}/\Gamma(\alpha)$, i.e.,

$$(1.3) \qquad k_\alpha(t) = t^{\alpha-1}/\Gamma(\alpha) = \sum_{j=1}^{Q} w_j e^{\lambda_j t} + O(\epsilon), \quad t \in [\delta, T]$$

where $\delta, T > 0$ and $\epsilon > 0$ is a given precision. The key is to determine $w_j$ and $\lambda_j$ in (1.3) in order that the desired accuracy up to $O(\epsilon)$ can be achieved.

In order to derive (1.3), Lubich and Schädle [22] expressed $k_\alpha(t)$ in terms of its inverse Laplace transform $k_\alpha(t) = \frac{1}{2\pi i} \int_{\mathcal{C}} \mathcal{L}[k_\alpha] e^{t\lambda} d\lambda$, where $\mathcal{L}[k_\alpha](\lambda)$ denotes the Laplace transform of $k_\alpha(t)$ and $\mathcal{C}$ is a suitable contour. Then a suitable quadrature was applied to approximate $\int_{\mathcal{C}} \mathcal{L}[k_\alpha] e^{t\lambda} d\lambda$, which leads to (1.3). This approach can be applied to construct fast methods for a wide class of nonlocal models. The method in [22] was then extended to calculate the discrete convolution (1.2) in [3, 31], where the coefficients $\omega_{n,j}$ are generated from the generating functions. A graded mesh version of [22] was developed in [20] and an application of [22] in the simulation of fractional-order viscoelasticity in complicated arterial geometries was proposed in [41]. The storage and computational cost of the fast methods in [3, 20, 22, 31] are $O(\log n_T)$ and $O(n_T \log n_T)$, respectively, which are much less than the direct methods with $O(n_T)$ memory and $O(n_T^2)$ operations. Recently, Baffet and Hesthaven [1, 2] have proposed to approximate $\mathcal{L}[k_\alpha]$ using the multipole approximation, which yields (1.3).

Another approach to derive (1.3) is based on the following integral expression

$$(1.4) \qquad k_\alpha(t) = \frac{\sin(\alpha\pi)}{\pi} \int_0^\infty \lambda^{-\alpha} e^{-t\lambda} d\lambda, \quad \alpha < 1.$$

Eq. (1.4) can be derived by inserting $\mathcal{L}[k_\alpha] = \lambda^{-\alpha}$ into Henrici's formula (see [5]), i.e., $k_\alpha(t) = \frac{1}{2\pi i} \int_0^\infty \left[ \mathcal{L}[k_\alpha](\lambda e^{-i\pi}) - \mathcal{L}[k_\alpha](\lambda e^{i\pi}) \right] e^{-t\lambda} d\lambda$. Li [18] transformed the above integral into its equivalent form, then multi-domain Legendre–Gauss quadrature was applied to approximate the transformed integral to obtain (1.3). Jiang et al [14] combined Jacobi–Gauss quadrature and multi-domain Legendre–Gauss quadrature to discretize (1.4) for $-1 < \alpha < 1$, then a global after-processing optimization technique was applied to further reduce the number of quadrature points, see also [40]. The exponential sum approximation for $t^{-\beta} (\beta > 0)$ has been studied in the literature, which can be used to design fast algorithm to approximate the fractional operators. The interested readers can refer to [4, 26]. In the references [1, 2, 3, 20, 22, 31], $w_j$ and $\lambda_j$ are complex, however, they are real in [18, 14, 4, 26]. Apart from the above mentioned fast methods for fractional operators, McLean [25] proposed to use a degenerate kernel for evaluating $k_\alpha * u(t)$.

In this work, we derive (1.3) by approximating (1.4) using a truncated Laguerre–Gauss (LG) quadrature. We list the main contributions of this work as follows.
- We follow and generalize the framework in [22] to resolve the short memory principle with a lag-term, see Section 3. By choosing suitable parameters, the current

method can be simplified as that in [1, 2, 18, 14], where the time domain is not divided into exponentially increasing subintervals. This approach simplifies the implementation of the algorithm. The present fast method unifies the calculation of the discrete convolutions to the approximation of both fractional integral and derivative operators with arbitrary accuracy (see Tables 4.2 and 5.5), for example, the trapezoidal rule for the fractional integral operator (see [8]) and the L1 method for the fractional derivative operator (see [7, 34]); see Figure 2.1.

- Given any basis $B$ ($B > 1$ is an integer), any stepsize $\tau$, any memory length $\Delta T \geq \tau$, and any precisions $\epsilon, \epsilon_0 > 0$, the truncation number $q_\alpha(N_\ell)$ (see (4.7) and (4.10)) of the truncated LG quadrature is determined in order that the overall error of the present fast method from the truncated LG quadrature is $O(\epsilon + \epsilon_0)$, see (4.14). The truncated LG quadrature and/or a relatively smaller basis $B$ saves memory and computational cost, see numerical results in Tables 4.1–4.2 and Figure 4.2.

We would like to emphasize that the *truncated* LG quadrature reduces the memory and computational cost significantly, see Table 4.1. The memory and computational cost in [22, 20] can be halved due to the symmetry of the trapezoidal rule, but operations with complex numbers are involved. In addition, the Gauss–Jacobi and Gauss–Legendre quadrature used in [18, 14] may not be truncated. Furthermore, the discretization error caused from the LG quadrature is independent of the stepsize and the regularity of the solution to the considered fractional differential equation (FDE), and does not appear to be sensitive to the fractional order $\alpha \in (-2, 1)$ as exhibited in Figure 4.2, which is competitive with the mutlipole approximation [2] in both accuracy and memory requirement, see Figure 4.3.

In real applications, analytical solutions to FDEs are unknown and often non-smooth, and may have strong singularity at $t = 0$, see, for example, [16, 23, 27, 33]. In order to resolve the singularity of the solution to the considered FDE, a graded mesh approach was adopted by some researchers [16, 27, 33]. In [33], an optimal graded mesh was obtained to achieve the global convergence of order $2 - \alpha$, $\alpha \in (0, 1)$, which is very effective when the fractional order is relatively large, but is less effective when the fractional order tends to zero. In this paper, we follow Lubich's approach [21] to deal with the singularity by introducing correction terms.

As we mentioned above, a rational approximation was made in [2] to approximate the Laplace transform of the fractional kernel, while the method in [18] used the Legendre–Gauss quadrature to approximate the transformed integral of (1.4). These approaches were initially designed for the fractional integral operator. However, we have found that they can be applied to the RL fractional derivative operator directly as done in the present work.

This paper is organized as follows. We follow the approach in [22] to present our fast method in Section 2. The short memory principle with lag terms is resolved in Section 3, it unifies the calculation of the discrete convolutions to the approximation of the fractional integral and derivative operators. The error analysis of the fast method is presented in Section 4, where all the parameters needed in numerical simulations are explicitly given. Numerical simulations are presented to verify the effectiveness of the fast method in Section 5 before the conclusion in the last section.

**2. A stable fast convolution.** In this section, we follow the approach given in [22] to develop our fast convolution.

The goal is to discretize the right-hand side of (1.4) using a highly accurate

3

numerical method. It is natural to use LG quadrature to approximate (1.4), i.e.,

$$(2.1) \qquad k_\alpha(t) = \frac{\sin(\alpha\pi)}{\pi} \int_0^\infty \lambda^{-\alpha} e^{-T\lambda} e^{-(t-T)\lambda} \, d\lambda \approx \frac{\sin(\alpha\pi)}{\pi} \sum_{j=0}^N \omega_j e^{-(t-T)\lambda_j},$$

where $\{\omega_j\}$ and $\{\lambda_j\}$ are the LG quadrature weights and points that correspond to the weight function $\lambda^{-\alpha} e^{-T\lambda}$. The quadrature (2.1) is exponentially convergent for any $t \geq T$ if $N$ is sufficiently large, which is discussed in Section 4.

Denote $t_n = n\tau$ ($n = 0, 1, ..., n_T$) as the grid point, where $\tau$ is the stepsize. We first restrict ourselves to $\alpha \in [0, 1)$. Using (2.1) and following the idea in [22], we present our stable fast convolution for approximating $k_\alpha * u(t)$ as follows:

- Step 1) Decompose the convolution $k_\alpha * u(t)$ as

$$(2.2)$$
$$k_\alpha * u(t) = \int_{t-\tau}^t k_\alpha(t-s)u(s)\, ds + \int_0^{t-\tau} k_\alpha(t-s)u(s)\, ds \equiv L^\alpha(u, t) + H^\alpha(u, t),$$

where we call $L^\alpha(u, t)$ and $H^\alpha(u, t)$ the local and history parts, respectively. Let $I_\tau^{(1)} u(t)$ be the linear interpolation of $u(t)$. Then the local and history parts can be approximated by

$$L^\alpha(u, t_n) \approx L^\alpha(I_\tau^{(1)}u, t_n) = L_\tau^{(\alpha,n)}u, \quad H^\alpha(u, t_n) \approx H^\alpha(I_\tau^{(1)}u, t_n) = H_\tau^{(\alpha,n)}u.$$

- Step 2) For every $t = t_n$, let $L$ be the smallest integer satisfying $t_n < 2B^L\tau$, where $B > 1$ is a positive integer. For $\ell = 1, 2, ..., L-1$, determine the integer $q_\ell$ such that

$$(2.3) \qquad s_\ell = q_\ell B^\ell \tau \quad \text{satisfies} \quad t_n - s_\ell \in [B^\ell \tau, (2B^\ell - 1)\tau].$$

Set $s_0 = t_n - \tau$ and $s_L = 0$. Then $t_n - \tau = s_0 > s_1 > \cdots > s_{L-1} > s_L = 0$.
- Step 3) Using (2.1), we approximate the history part $H_\tau^{(\alpha,n)}u = H^\alpha(I_\tau u, t_n)$ by

$$(2.4) \quad H_\tau^{(\alpha,n)}u \approx \frac{\sin(\alpha\pi)}{\pi} \sum_{\ell=1}^L \sum_{j=0}^N \omega_j^{(\ell)} e^{-(t_n - s_{\ell-1} - T_{\ell-1})\lambda_j^{(\ell)}} y(s_{\ell-1}, s_\ell, \lambda_j^{(\ell)}) = {}_F H_\tau^{(\alpha,n)}u$$

with $y(s_{\ell-1}, s_\ell, \lambda)$ given by

$$(2.5) \qquad y(s_{\ell-1}, s_\ell, \lambda) = \int_{s_\ell}^{s_{\ell-1}} e^{(s - s_{\ell-1})\lambda} I_\tau^{(1)} u(s)\, ds,$$

where $\{\omega_j^{(\ell)}\}$ and $\{\lambda_j^{(\ell)}\}$ are the LG quadrature weights and points that corresponds to the weight function $\lambda^{-\alpha} e^{-T_{\ell-1}\lambda}$ (see (7.27) in [32]), and $T_{\ell-1} = B^{\ell-1}\tau$ satisfying $t_n - s - T_{\ell-1} \geq 0$ for all $t_n - s \in [B^{\ell-1}\tau, (2B^\ell - 1)\tau]$, $s \in [s_\ell, s_{\ell-1}]$. Here $y(s) = y(s, s_\ell, \lambda_j^{(\ell)})$ used in (2.4) that is defined by (2.5) satisfies the following ODE

$$(2.6) \qquad y'(s) = -\lambda_j^{(\ell)} y(s) + I_\tau^{(1)} u(s), \qquad y(s_\ell) = 0,$$

which can be exactly solved by the following recursive relation

$$(2.7) \qquad y(t_{m+1}) = e^{-\lambda_j^{(\ell)}\tau} y(t_m) + e^{-\lambda_j^{(\ell)}\tau} \int_{t_m}^{t_{m+1}} e^{\lambda_j^{(\ell)}(s - t_m)} I_\tau^{(1)} u(s)\, ds.$$

4

- Step 4) Calculate the local part $L_\tau^{(\alpha,n)}u = \int_{t_n-\tau}^{t_n} k_\alpha(t_n-s)I_\tau^{(1)}u(s)\,ds$ with

$$(2.8) \qquad L_\tau^{(\alpha,n)}u = L(I_\tau^{(1)}u, t_n) = \frac{\tau^\alpha}{\Gamma(2+\alpha)}(u_n - u_{n-1}).$$

Combining Steps 1)–4), we obtain our fast convolution for approximating $k_\alpha * u(t)$. The above fast convolution has the same storage and computational cost as that in [22], the main differences are listed below:

i) The LG quadrature is applied instead of the trapezoidal rule to approximate the history part $H^\alpha(I_\tau u, t) = \int_0^{t-\tau} k_\alpha(t-s)I_\tau u(s)\,ds$, that is, the history part $H^\alpha(I_\tau u, t)$ in [22] was approximated by

$$(2.9) \quad {}_F\hat{H}_\tau^{(\alpha,n)}u = \mathrm{Im}\left\{\sum_{\ell=1}^{L}\sum_{j=-N}^{N-1} \hat{\omega}_j^{(\ell)}\mathcal{L}[k_\alpha](\hat{\lambda}_j^{(\ell)})e^{(t_n-s_{\ell-1})\hat{\lambda}_j^{(\ell)}}\hat{y}(s_{\ell-1}, s_\ell, \hat{\lambda}_j^{(\ell)})\right\},$$

where $\{\hat{\omega}_j^{(\ell)}\}$ and $\{\hat{\lambda}_j^{(\ell)}\}$ are the weights and quadrature points for the Talbot contour $\Gamma_\ell$, and $\hat{y}(s) = \hat{y}(s, s_\ell, \hat{\lambda}_j^{(\ell)}) = \int_{s_\ell}^{s} e^{-(s-s_{\ell-1})\hat{\lambda}_j^{(\ell)}}I_\tau^{(1)}u(s)\,ds$ satisfies the following ODE

$$(2.10) \qquad \hat{y}'(s) = \hat{\lambda}_j^{(\ell)}\hat{y}(s) + I_\tau^{(1)}u(s), \qquad y(s_\ell) = 0.$$

ii) We solve a stable ODE (2.6) instead of a possibly unstable ODE (2.10) that may affect the stability and accuracy of (2.9). Indeed, for the Talbot contour used in [22] (see also the parabolic contour or hyperbolic contour discussed in [38]), there exist $\hat{\lambda}_j^{(\ell)}$'s, whose real parts are positive. Numerical tests show that (2.9) still works well since one may not solve (2.10) for a long time, which reduces the iteration error from solving (2.10) even though the real part of some $\hat{\lambda}_j^{(\ell)}$ is positive.

The above fast convolution Step 1) – Step 4) holds for $\alpha < 1$, that is, the RL fractional derivative operator of order $-\alpha$ is thus discretized if $\alpha < 0$.

EXAMPLE 2.1. *Let $u(s) = 1 + s$ in (2.4) and define the relative error*

$$(2.11) \qquad e_n = \frac{\left|H^\alpha(u, t_n) - {}_F H_\tau^{(\alpha,n)}u\right|}{\left|H^\alpha(u, t_n)\right|}, \quad \hat{e}_n = \frac{\left|H^\alpha(u, t_n) - {}_F\hat{H}_\tau^{(\alpha,n)}u\right|}{\left|H^\alpha(u, t_n)\right|}$$

*where ${}_F H_\tau^{(\alpha,n)}u$ and ${}_F\hat{H}_\tau^{(\alpha,n)}u$ are defined by (2.4) and (2.9), respectively.*

We choose $u(s) = 1 + s$ in order that the errors $e_n$ and $\hat{e}_n$ mainly come from the quadrature used in the discretization of the kernel $k_\alpha(t)$. We show the errors $e_n$ and $\hat{e}_n$ for different $\alpha$ in Figure 2.1. We can see that the LG quadrature shows better accuracy than the trapezoidal rule based on the Talbot contour in this example.

**3. Short memory principle with lag terms.** In this section, we generalize the fast convolution in the previous section to resolve the short memory principle (see [6, 29]). The error analysis of the method is given in the next section.

Let $\Delta T \geq 0$ be a memory length. Divide the convolution $k_\alpha * u(t)$ into the local part $L_{\Delta T}^\alpha(u, t)$ and the history part $H_{\Delta T}^\alpha(u, t)$ as shown below

$$(3.1) \qquad k_\alpha * u(t) \equiv L_{\Delta T}^\alpha(u, t) + H_{\Delta T}^\alpha(u, t),$$

5

(a) $\alpha = -0.5$.　　　　　　　　　　　(a) $\alpha = 0.5$.

FIG. 2.1. *Comparison between the trapezoidal rule based on the Talbot contour (black curve) and the LG quadrature (red curve), $\tau = 1$, $B = 5$, $N = 100$. The optimal contour $z(\theta, t) = t\,(-0.4814 + 0.6443(\theta \cot(\theta) + i0.5653\theta))$ obtained in [37] is applied here, i.e., $\hat{\lambda}_j^{(\ell)}$ in (2.9) is given by $\hat{\lambda}_j^{(\ell)} = z(\theta_j, N/(2T_\ell - \tau))$ with the corresponding weight $\hat{\omega}_j^{(\ell)} = \partial_\theta z(\theta_j, N/(2T_\ell - \tau))$, where $\theta_j = (2j+1)\pi/(2N), j = -N, ..., N-1$ and $T_\ell = B^\ell \tau$, $B = 5, N = 32$.*

$$(3.2) \qquad L_{\Delta T}^\alpha(u, t) = \int_{\max\{0, t-\Delta T\}}^{t} k_\alpha(t-s) u(s)\, ds,$$

$$(3.3) \qquad H_{\Delta T}^\alpha(u, t) = \int_{0}^{\max\{0, t-\Delta T\}} k_\alpha(t-s) u(s)\, ds.$$

If we drop the history part in (3.1), then the remaining part is the famous short memory principle rule (see [29]). However, the short memory principle has not been widely applied, since $L_{\Delta T}^\alpha(u, t)$ is not a good approximation of $k_\alpha * u(t)$. Our goal is to develop a good numerical approximation of $H_{\Delta T}^\alpha(u, t)$, such that the storage and computational cost are reduced significantly compared with the direct approximation. The local part $L_{\Delta T}^\alpha(u, t)$ is just the fractional operator defined on $[\max\{0, t-\Delta T\}, t]$, which can be discretized directly by the known methods, see [11, 21, 24, 34]. Next, we introduce how to discretize (3.1) efficiently and accurately.

**3.1. Interpolations.** In this work, the discretization of (3.2) and (3.3) is based on the interpolation of $u$. Specifically, $L_{\Delta T}^\alpha(u, t)$ and $H_{\Delta T}^\alpha(u, t)$ are approximated by $L_{\Delta T}^\alpha(I_\tau^L u, t)$ and $H_{\Delta T}^\alpha(I_\tau^H u, t)$, respectively, where $I_\tau^L$ and $I_\tau^H$ are two suitable piecewise interpolation operators, which will be discussed in the following.

Linear interpolation is simple and has been applied widely in the discretization of the fractional operators, see [8, 15, 17, 34]. Several quadratic interpolations have been used to discretize the Caputo fractional operators, see, for example, [40, 19, 24]. We adopt the quadratic interpolation in [24] to illustrate the implementation of the present fast algorithm. Define the local interpolation operator $\Pi_\tau^j$ as

$$(3.4) \qquad \Pi_\tau^j u(t) = \sum_{k=1}^{3} u_{j+k-1} F_k^{(j)}(t), \quad t \in [t_j, t_{j+1}], j \geq 0,$$

where $F_1^{(j)} = \frac{(t-t_{j+1})(t-t_{j+2})}{(t_j-t_{j+1})(t_j-t_{j+2})}$, $F_2^{(j)} = \frac{(t-t_j)(t-t_{j+2})}{(t_{j+1}-t_j)(t_{j+1}-t_{j+2})}$, and $F_3^{(j)} = \frac{(t-t_j)(t-t_{j+1})}{(t_{j+2}-t_j)(t_{j+2}-t_{j+1})}$. Let $\Pi_\tau^{-1} u(t) = \Pi_\tau^0 u(t)$. Then, for each $n \geq 1$, the quadratic interpolation $I_\tau^{(2,n)}$ is

6

defined by

$$(3.5) \qquad I_\tau^{(2,n)} u(t) = \begin{cases} \Pi_\tau^j u(t), & t \in [t_j, t_{j+1}], 0 \le j \le n-2, \\ \Pi_\tau^{n-2} u(t), & t \in [t_{n-1}, t_n]. \end{cases}$$

For each $n \ge 1$, we define $I_\tau^L$ and $I_\tau^H$ as follows

$$(3.6) \qquad I_\tau^L u(t) = I_\tau^{(2,n)} u(t), \qquad t \in [\max\{t_n - \Delta T, 0\}, t_n],$$

$$(3.7) \qquad I_\tau^H u(t) = I_\tau^{(2,n)} u(t), \qquad t \in [0, \max\{t_n - \Delta T, 0\}].$$

Let $j_n = \max\{n - n_0, 0\}$. Then $L_{\Delta T, \tau}^{(\alpha,n)} u = L_{\Delta T}^\alpha(I_\tau^L u, t_n)$ and $H_{\Delta T, \tau}^{(\alpha,n)} u = H_{\Delta T}^\alpha(I_\tau^H u, t_n)$ are given by

$$(3.8) \qquad L_{\Delta T, \tau}^{(\alpha,n)} u = \sum_{j=j_n}^{n-2} \left( b_{n-1-j}^{(1)} u_j + b_{n-1-j}^{(2)} u_{j+1} + b_{n-1-j}^{(3)} u_{j+2} \right) + \sum_{j=0}^{2} d_j u_{n+j-2},$$

$$(3.9) \qquad H_{\Delta T, \tau}^{(\alpha,n)} u = \sum_{j=0}^{j_n-1} \left( b_{n-1-j}^{(1)} u_j + b_{n-1-j}^{(2)} u_{j+1} + b_{n-1-j}^{(3)} u_{j+2} \right),$$

where $b_{n-j-1}^{(k)} = \int_{t_j}^{t_{j+1}} k_\alpha(t_n - s) F_k^{(j)}(s)\, ds$ and $d_j = \int_{t_{n-1}}^{t_n} k_\alpha(t_n - s) F_j^{(n-2)}(s)\, ds$, i.e.,

$$(3.10) \qquad d_0 = \frac{-\alpha \tau^\alpha}{2\Gamma(\alpha+3)}, \quad d_1 = \frac{\alpha(3+\alpha)\tau^\alpha}{\Gamma(\alpha+3)}, \quad d_2 = \frac{(4+\alpha)\tau^\alpha}{2\Gamma(\alpha+3)},$$

$$(3.11) \qquad b_j^{(1)} = \frac{\tau^\alpha}{2\Gamma(\alpha)} \left[ \ell_j^{(\alpha+2)} - (2j-1)\ell_j^{(\alpha+1)} + j(j-1)\ell_j^{(\alpha)} \right],$$

$$(3.12) \qquad b_j^{(2)} = -\frac{\tau^\alpha}{\Gamma(\alpha)} \left[ \ell_j^{(\alpha+2)} - 2j\ell_j^{(\alpha+1)} + (j+1)(j-1)\ell_j^{(\alpha)} \right],$$

$$(3.13) \qquad b_j^{(3)} = \frac{\tau^\alpha}{2\Gamma(\alpha)} \left[ \ell_j^{(\alpha+2)} - (2j+1)\ell_j^{(\alpha+1)} + j(j+1)\ell_j^{(\alpha)} \right],$$

$$(3.14) \qquad \ell_j^{(\alpha)} = \frac{1}{\alpha} \left[ (j+1)^\alpha - j^\alpha \right].$$

Some researchers used other interpolations in the discretization of the fractional integral and derivative operators, we refer readers to [19, 7, 8, 17, 10, 42].

**3.2. Corrections.** From (3.8)–(3.9) and for any $\Delta T \ge 0$, we always have

$$(3.15) \qquad D_{\Delta T, \tau}^{(\alpha,n)} u = L_{\Delta T, \tau}^{(\alpha,n)} u + H_{\Delta T, \tau}^{(\alpha,n)} u = \sum_{j=0}^{n} w_{n,j} u_j,$$

where $w_{n,j}$ can be derived from (3.10)–(3.13), which do not give explicitly. Clearly, $D_{\Delta T, \tau}^{(\alpha,n)} u$ is just the second-order trapezoidal rule (or the $(2+\alpha)$-order L1 method) for the fractional integral of order $\alpha > 0$ (or the RL fractional derivative of order $0 < -\alpha < 1$) if the linear interpolation is used and $u(t)$ is sufficiently smooth, see [7, 8]. For quadratic interpolation (3.5), $D_{\Delta T, \tau}^{(\alpha,n)} u$ achieves the $(3+\alpha)$-order accuracy for the RL derivative of order $0 < -\alpha < 1$ when $u(t)$ is smooth. However, the approximation $D_{\Delta T, \tau}^{(\alpha,n)} u$ defined by (3.15) is not a good approximation of $k_\alpha * u(t_n)$ when $u(t)$ has strong singularities.

7

In this work, we follow Lubich's idea (see [21]) to use correction terms to capture the singularity of the solution $u(t)$ to the considered FODE. The correction method is based on the assumption that the solution $u(t)$ has the following form

$$(3.16) \qquad u(t) - u(0) = \sum_{j=1}^{m} c_j t^{\sigma_j} + t^{\sigma_{m+1}} \tilde{u}(t), \quad 0 < \sigma_j < \sigma_{j+1},$$

where $\tilde{u}(t)$ is uniformly bounded for $t \in [0, T]$. Readers can refer to [8, 23, 13, 29] for more detailed results of the regularity of FODEs.

Combining (3.15) and (3.16) gives the following correction method

$$(3.17) \qquad D_{\Delta T, \tau}^{(\alpha, n, m)} u = D_{\Delta T, \tau}^{(\alpha, n)} u + \tau^\alpha \sum_{j=1}^{m} W_{n,j}(u_j - u_0),$$

where $W_{n,j}$ are the starting weights that are chosen such that

$$(3.18) \qquad D_{\Delta T, \tau}^{(\alpha, n)} u + \tau^\alpha \sum_{j=1}^{m} W_{n,j}(u_j - u_0) = k_\alpha * u(t_n) = \frac{\tau^\alpha n^{\sigma_k + \alpha}}{\Gamma(\sigma_k + 1 + \alpha)}, 1 \le k \le m$$

for $u(t) = t^{\sigma_k}, 0 < \sigma_k < \sigma_{k+1}$. For each $n > 0$, one can resolve $W_{n,j}(1 \le j \le m)$ from the above linear system and $W_{n,j}$ are independent of $\tau$. The error analysis of the direct method (3.17) is discussed in Section 4. Readers can refer to [9, 21, 44] for more discussions of the correction method.

**3.3. The fast implementation.** In this subsection, we generalize the fast algorithm in Section 2 to approximate $D_{\Delta T, \tau}^{(\alpha, n)}$ defined by (3.15), which is given as follows.
- Step A) Decompose $D_{\Delta T, \tau}^{(\alpha, n)} u$ into two parts as $D_{\Delta T, \tau}^{(\alpha, n)} u = L_{\Delta T, \tau}^{(\alpha, n)} u + H_{\Delta T, \tau}^{(\alpha, n)} u$.
- Step B) Assume that $\Delta T = n_0 \tau = t_{n_0}$. For every $t = t_n, n \ge n_0$, let $L$ be the smallest integer satisfying $t_{n-n_0+1} < 2B^L \tau$. For $\ell = 1, 2, ..., L - 1$, determine the integer $q_\ell$ such that

$$(3.19) \qquad s_\ell = q_\ell B^\ell \tau \quad \text{satisfies} \quad t_{n-n_0+1} - s_\ell \in [B^\ell \tau, (2B^\ell - 1)\tau].$$

Set $s_0 = t_{n-n_0+1} - \tau$ and $s_L = 0$.
- Step C) Let $\hat{t}_n = t_n - \Delta T + \tau = t_{n-n_0+1} \ge \tau$. Then the history part $H_{\Delta T, \tau}^{(\alpha, n)} u = H_{\Delta T}^\alpha(I_\tau^H u, t_n)$ is approximated by

$$H_{\Delta T, \tau}^{(\alpha, n)} u = \frac{\sin(\alpha\pi)}{\pi} \sum_{\ell=1}^{L} \int_0^\infty \lambda^{-\alpha} e^{-(T_{\ell-1} + \Delta T - \tau)\lambda} e^{-(\hat{t}_n - s_{\ell-1} - T_{\ell-1})\lambda} y(s_{\ell-1}, s_\ell, \lambda) \, d\lambda$$

$$(3.20) \qquad \approx \frac{\sin(\alpha\pi)}{\pi} \sum_{\ell=1}^{L} \sum_{j=0}^{q_\alpha(N_\ell)} \omega_j^{(\ell)} e^{-(\hat{t}_n - s_{\ell-1} - T_{\ell-1})\lambda_j^{(\ell)}} y(s_{\ell-1}, s_\ell, \lambda_j^{(\ell)}) := {}_F H_{\Delta T, \tau}^{(\alpha, n)} u,$$

where $\{\omega_j^{(\ell)}\}$ and $\{\lambda_j^{(\ell)}\}$ are the LG quadrature weights and points that correspond to the weight function $\lambda^{-\alpha} e^{-(T_{\ell-1} + \Delta T - \tau)\lambda}$, $q_\alpha(N_\ell)$ is defined by (4.7), and $y(s_{\ell-1}, s_\ell, \lambda_j^{(\ell)})$ can be obtained exactly by solving the following linear ODE

$$(3.21) \qquad y'(s) = -\lambda_j^{(\ell)} y(s) + I_\tau^H u(s), \quad y(s_\ell) = 0,$$

see also (2.6) and (2.7).

8

- Step D) Calculate the local part $L_{\Delta T,\tau}^{(\alpha,n)} u = L_{\Delta T}^\alpha (I_\tau^L u, t_n)$.

The fast algorithm for the discretization of $k_\alpha * u(t)$ is now given by

$$(3.22) \qquad {}_F D_{\Delta T,\tau}^{(\alpha,n,m)} u = L_{\Delta T,\tau}^{(\alpha,n)} u + {}_F H_{\Delta T,\tau}^{(\alpha,n)} u + \tau^\alpha \sum_{j=1}^m W_{n,j}(u_j - u_0),$$

where $L_{\Delta T,\tau}^{(\alpha,n)}$ is given by (3.8), ${}_F H_{\Delta T,\tau}^{(\alpha,n)} u$ is given by (3.20), and the starting weights $W_{n,j}$ are determined by the linear system (3.18).

Next, we analyze the complexity of the present fast method (3.22). For the local part $L_{\Delta T,\tau}^{(\alpha,n)} u$, the memory requirement is $O(n_0)$ with the computational cost of $O(n_0(n_T - n_0))$ for all $n_0 < n \le n_T$. For the history part ${}_F H_{\Delta T,\tau}^{(\alpha,n)} u$, we have $O(\sum_\ell^L q_\alpha(N_\ell))$ active memory and $O((n_T - n_0)\sum_\ell^L q_\alpha(N_\ell))$ operations, where $L = \log_B(n_T - n_0)$. An additional cost is required to obtain the starting weights $W_{n,j}$ in (3.22), which can be performed by use of fast Fourier transform with $O(n_T \log(n_T))$ arithmetic operations. Hence, the overall active memory and computational cost are $O(n_0 + \sum_\ell^L q_\alpha(N_\ell))$ and $O(n_0 n_T + (n_T - n_0)\sum_\ell^L q_\alpha(N_\ell))$, respectively.

**4. Error analysis.** In this section, we analyse the overall discretization error of the fast method (3.22) in Section 3. Firstly, we present the exponential convergence rate of the LG quadrature used in (2.4) and (3.20). Then we show how to choose $N_\ell$ and $q_\alpha(N_\ell)$ used in (3.20), such that the desired accuracy is maintained with the use of the minimum number of the quadrature points.

For simplicity, we denote

$$(4.1) \qquad I^\alpha[T, f] = \int_0^\infty \lambda^\alpha e^{-T\lambda} f(\lambda) \, d\lambda \quad \text{and} \quad I^\alpha[f] = \int_0^\infty \lambda^\alpha e^{-\lambda} f(\lambda) \, d\lambda.$$

The LG quadrature for $I^\alpha[T, f]$ and $I^\alpha[f]$ are given by (see [32])

$$(4.2) \qquad Q_N^\alpha[T, f] = T^{-\alpha-1} \sum_{j=0}^N w_j^{(\alpha)} f(\lambda_j/T), \quad Q_N^\alpha[f] = \sum_{j=0}^N w_j^{(\alpha)} f(\lambda_j),$$

where $\lambda_j$ are the roots of the Laguerre polynomial $L_{N+1}^{(\alpha)}(\lambda)$, and $w_j^{(\alpha)}$ are the corresponding weights given by

$$(4.3) \qquad w_j^{(\alpha)} = \frac{\Gamma(N + \alpha + 1)\lambda_j}{(N + \alpha + 1)(N + 1)!} \left( L_N^{(\alpha)}(\lambda_j) \right)^{-2}.$$

We show the convergence of the quadrature $Q_N^\alpha[T, e^{-t\lambda}]$ and the property of the quadrature weight $w_j^{(\alpha)}$ (see (4.3)) in the following two theorems. The proofs are given in Appendix A.

THEOREM 4.1. *Let $t \ge 0, T > 0, \alpha > -1$, and $N$ be sufficiently large. Then*

$$(4.4) \qquad \left| I^\alpha[T, e^{-t\lambda}] - Q_N^\alpha[T, e^{-t\lambda}] \right| \le C_{\alpha,N} T^{-\alpha-1} \left( \frac{t/T}{1 + t/T} \right)^{2N},$$

*where $C_{\alpha,N}$ is bounded $-1 < \alpha \le 0$ and $C_{\alpha,N} \le C_\alpha N^\alpha$ for $\alpha > 0$.*

THEOREM 4.2. *Let $\alpha > -1$ and $w_j^{(\alpha)}$ be defined by (4.3). If $N$ and $j$ are sufficiently large, then there exists a positive constant $C$ independent of $N$ such that*

$$(4.5) \qquad w_j^{(\alpha)} \le C(N + 1)^\alpha e^{-0.25\pi^2(j+1)^2/(N+1)}.$$

9

Given a precision $\epsilon_0 > 0$, the truncated LG quadrature is given by

$$(4.6) \qquad Q_{N,\epsilon_0}^{\alpha}[T,f] = T^{-\alpha-1} \sum_{j=0}^{q_{-\alpha}(N)} w_j^{(\alpha)} f(\lambda_j/T),$$

where $q_{-\alpha}(N)$ is a positive number given by

$$(4.7) \qquad q_{-\alpha}(N) = \min\left\{N, \left\lceil 2\pi^{-1}\sqrt{(N+1)\log((N+1)^{\alpha}\epsilon_0^{-1})}\right\rceil - 1\right\}.$$

From (4.5), we can find the smallest integer $j$ satisfying $(N+1)^{\alpha}e^{-0.25\pi^2(j+1)^2/(N+1)} \leq \epsilon_0$, which yields (4.7). We choose $\epsilon_0 = 10^{-16}$ in this paper.

Figure 4.1 (a) shows the exponential decay of $w_j^{(-\alpha)}$ when $j \geq q_{\alpha}(128)$ for different fractional orders $\alpha = -1.8, -1.2, -0.8, -0.2, 0.2, 0.8$. Figure 4.1 (b) displays similar behaviors as shown in Figure 4.1 (a). We can see that Eq. (4.7) works well, and $q_{\alpha}(N)$ is not very sensitive to the fractional order $\alpha \in (-2,1)$. For example, $q_{\alpha}(128) = (48, 47, 46, 44, 43, 41)$ (or $q_{\alpha}(256) = (69, 67, 65, 62, 61, 58)$) for $\alpha = (-1.8, -1.2, -0.8, -0.2, 0.2, 0.8)$, and $q_{\alpha}(N) \ll N$ when $N$ is sufficiently large. For $\alpha \leq -2$, similar results are obtained, which is not shown here.



(a) $N = 128$.         (b) $N = 256$.

FIG. 4.1. *The exponential decay of the quadrature weights $w_j^{(-\alpha)}$ defined by (4.3).*

For notational simplicity, we denote

$$\widehat{T}_{\ell} = T_{\ell-1} + \Delta T - \tau, \quad \widehat{H}_{\ell}^n(s) = \hat{t}_n - T_{\ell-1} - s.$$

Next, we investigate how to estimate $N_{\ell}$ in (3.20), such that the LG quadrature $Q_{N_{\ell}}^{-\alpha}[\widehat{T}_{\ell}, e^{-\widehat{H}_{\ell}^n(s)\lambda}]$ to the integral $I^{-\alpha}[\widehat{T}_{\ell}, e^{-\widehat{H}_{\ell}^n(s)\lambda}]$ preserves the accuracy up to $O(\epsilon)$ for all $s \in [s_{\ell}, s_{\ell-1}]$.

From Theorem 4.1, we know that the error of $Q_{N_{\ell}}^{-\alpha}[\widehat{T}_{\ell}, e^{-\widehat{H}_{\ell}^n(s)\lambda}]$ mainly depends on the following term

$$\left(\widehat{H}_{\ell}^n(s)/\widehat{T}_{\ell}\right)^{2N} = \left(\frac{\hat{t}_n - s - T_{\ell-1}}{T_{\ell-1} + \Delta T - \tau}\right)^{2N}, \quad s \in [s_{\ell}, s_{\ell-1}].$$

Using (3.19) and $T_{\ell-1} = B^{\ell-1}\tau$ gives

$$(4.8) \qquad 0 \leq \frac{\hat{t}_n - s - T_{\ell-1}}{T_{\ell-1} + \Delta T - \tau} \leq \frac{2B - 1 - B^{1-\ell}}{1 + B^{1-\ell}(\Delta T/\tau - 1)} = \mathcal{T}_{\ell} \leq 2B - 1, \quad \forall \ell \geq 1.$$

Using the above inequality and Eq. (4.4) yields

(4.9)
$$\left| I^{-\alpha}[\widehat{T}_\ell, e^{-\widehat{H}_\ell^n(s)\lambda}] - Q_{N_\ell}^{-\alpha}[\widehat{T}_\ell, e^{-\widehat{H}_\ell^n(s)\lambda}] \right|$$
$$\leq C_{\alpha,N_\ell}\widehat{T}_\ell^{\alpha-1}\left(\frac{\widehat{H}_\ell^n(s)/\widehat{T}_\ell}{1+\widehat{H}_\ell^n(s)/\widehat{T}_\ell}\right)^{2N_\ell} \leq C_{\alpha,N_\ell}\widehat{T}_\ell^{\alpha-1}\left(\frac{\mathcal{T}_\ell}{1+\mathcal{T}_\ell}\right)^{2N_\ell}.$$

Since the relative error of (4.9) is independent of $\widehat{T}_\ell$, we can let $(\mathcal{T}_\ell/(\mathcal{T}_\ell+1))^{2N_\ell} \leq \epsilon$, which yields the minimum $N_\ell$ given by

(4.10)
$$N_\ell = \left\lceil \frac{\log \epsilon}{2\log(\frac{\mathcal{T}_\ell}{\mathcal{T}_\ell+1})} \right\rceil, \quad \mathcal{T}_\ell = \frac{2B-1-B^{1-\ell}}{1+B^{1-\ell}(\Delta T/\tau-1)}.$$

From (4.6) and (4.9), we derive that the pointwise error of the truncated quadrature $Q_{N_\ell,\epsilon_0}^{-\alpha}[\widehat{T}_\ell, e^{-\widehat{H}_\ell^n(s)\lambda}]$ for all $s \in [s_\ell, s_{\ell-1}]$ is given by

(4.11)  $E_{\epsilon,\epsilon_0}^{-\alpha,\ell}[e^{-\widehat{H}_\ell^n(s)\lambda}] = I^{-\alpha}[\widehat{T}_\ell, e^{-\widehat{H}_\ell^n(s)\lambda}] - Q_{N_\ell,\epsilon_0}^{-\alpha}[\widehat{T}_\ell, e^{-\widehat{H}_\ell^n(s)\lambda}] = O(\epsilon) + O(\epsilon_0),$

where $Q_{N_\ell,\epsilon_0}^{-\alpha}$ is defined by (4.6) and $N_\ell$ is given by (4.10).

Next, we present the error bound of the fast method in Section 3. Denote by

(4.12)
$$R^{(n)} = \int_0^{t_n} k_\alpha(t_n-s)u(s)\,\mathrm{d}s - D_{\Delta T,\tau}^{(\alpha,n,m)}u.$$

Note that the above discretization error $R^{(n)}$ depends on the smoothness of $u(t)$ and the discretization method $D_{\Delta T,\tau}^{(\alpha,n,m)}u$. If $u(t)$ is sufficiently smooth, no correction terms are needed to achieve $(2+\alpha)$-order (or $(3+\alpha)$-order) accuracy if linear (or quadratic) interpolation is applied for $\alpha < 0$. If $u(t)$ satisfies (3.16), then the global $(2+\alpha)$-order (or $(3+\alpha)$-order) accuracy can be achieved for $\sigma_{m+1} \geq 2$ (or $\sigma_{m+1} \geq 3$). In numerical simulations, the condition $\sigma_{m+1} \geq 2$ (or $\sigma_{m+1} \geq 3$) does not need to be satisfied, $(2+\alpha)$-order (or $(3+\alpha)$-order) accurate numerical solutions are observed far from $t = 0$. In the numerical simulations, only a small number of correction terms are sufficient to achieve very accurate numerical solutions; see numerical simulations in the following section and see also related results in [44].

From (3.20) and (4.11), we have

(4.13)
$$\left| H_{\Delta T}^\alpha(I_\tau^H u, t_n) - {}_F H_{\Delta T,\tau}^{(\alpha,n)}u \right| = \left| \frac{\sin(\alpha\pi)}{\pi}\sum_{\ell=1}^L E_{\epsilon,\epsilon_0}^{-\alpha,\ell}[e^{-(\hat{t}_n-T_{\ell-1}-s_{\ell-1})\lambda}y(s_{\ell-1},s_\ell,\lambda)] \right|$$

$$= \left| \frac{\sin(\alpha\pi)}{\pi}\sum_{\ell=1}^L \int_{s_\ell}^{s_{\ell-1}} E_{\epsilon,\epsilon_0}^{-\alpha,\ell}[e^{-\widehat{H}_\ell^n(s)\lambda}]I_\tau^H u(s)\,\mathrm{d}s \right|$$

$$\leq C(\epsilon+\epsilon_0)\int_0^{\hat{t}_n}\left| I_\tau^H u(s) \right|\,\mathrm{d}s \leq C\max\{0, t_{n+1}-\Delta T\}\|u\|_\infty(\epsilon+\epsilon_0).$$

Combining (4.12) and (4.13) yields

(4.14)
$$\left| k_\alpha * u(t_n) - {}_F D_{\Delta T,\tau}^{(\alpha,n,m)}u \right| = \left| H_{\Delta T}^\alpha(I_\tau^H u, t_n) - {}_F H_{\Delta T,\tau}^{(\alpha,n)}u + R^{(n)} \right|$$
$$\leq C\max\{0, t_{n+1}-\Delta T\}\|u\|_\infty(\epsilon+\epsilon_0) + \left| R^{(n)} \right|,$$

11

where the error in (4.14) originates from two parts: the LG quadrature for discretizing $I^{-\alpha}[\widehat{T}_\ell, e^{-\widehat{H}_\ell^n(s)\lambda}]$ (see (4.11)) and the discretization error defined by (4.12).

Next, we numerically study the error caused by the LG quadrature. Let $m = 0$ and $u(t) = 1 + t$. Then $R^{(n)}$ in (4.14) is zero. Denote the relative pointwise error

$$e_n = (k_\alpha * u(t_n))^{-1} \left| k_\alpha * u(t_n) - {}_F D_{\Delta T, \tau}^{(\alpha, n, 0)} u \right|, \quad 1 \le n \le n_T = T/\tau,$$

where $k_\alpha * u(t_n) = t_n^\alpha/\Gamma(1 + \alpha) + t_n^{\alpha+1}/\Gamma(2 + \alpha)$ for $u(t) = 1 + t$.

Given a precision $\epsilon = 10^{-10}$, the maximum relative error $\|e\|_\infty = \max_{0 \le n \le n_T} |e_n|$, the total number of the quadrature points $\sum N_\ell$, and the total number of the truncated quadrature points $\sum q_\alpha(N_\ell)$ are shown in Table 4.1 for different basis $B$, $\alpha = -0.5, \tau = 0.1, \Delta T = 1$, and $T = 10^4$. We can see that the truncated LG quadrature saves memory. A relatively smaller basis $B$ needs less quadrature points and thus saves memory, which can be explained from (4.8) and (4.9). Eq. (4.9) implies a faster convergence as $B$ decreases, due to $T_\ell \approx 2B - 1$, $\ell$ is sufficiently large. Hence, a relatively smaller $B$ means that smaller $N_\ell$ are needed to achieve high accuracy that leads to the use of less LG quadrature points.

We change the precision $\epsilon$ and show the corresponding relative maximum errors $\|e\|_\infty$ in Table 4.2. We can see that $\|e\|_\infty$ increases as $\epsilon$ increases and the total number of the quadrature points are reduced. It also shows that much better results are obtained than the theoretical prediction, see also the related results in Table 5.5.

TABLE 4.1
*The maximum relative error $\|e\|_\infty$ under the precision $\epsilon = 10^{-10}$, $\alpha = -0.5, \tau = 0.1, \Delta T = 1$, and $T = 10^4$.*

| $B$ | $\sum N_\ell$ | $\sum q_\alpha(N_\ell)$ | $\|e\|_\infty$ | $B$ | $\sum N_\ell$ | $\sum q_\alpha(N_\ell)$ | $\|e\|_\infty$ |
|---|---|---|---|---|---|---|---|
| 2 | 583 | 389 | 6.8651e-13 | 30 | 3333 | 545 | 7.2134e-13 |
| 3 | 605 | 323 | 6.6718e-13 | 40 | 3583 | 514 | 7.3190e-13 |
| 4 | 687 | 318 | 7.4246e-13 | 50 | 4521 | 576 | 7.3391e-13 |
| 5 | 783 | 320 | 7.4754e-13 | 60 | 5463 | 637 | 7.4632e-13 |
| 8 | 1151 | 370 | 7.4377e-13 | 70 | 6406 | 689 | 7.3451e-13 |
| 10 | 1246 | 359 | 7.1388e-13 | 80 | 7345 | 737 | 7.4678e-13 |
| 15 | 1598 | 376 | 6.9954e-13 | 90 | 8288 | 785 | 7.3892e-13 |
| 20 | 2173 | 438 | 7.0761e-13 | 100 | 9231 | 829 | 7.4941e-13 |



(a)                          (b)

FIG. 4.2. *(a) The maximum relative error $\|e\|_\infty$ against the basis $B$; (b) The total number $\sum q_\alpha(N_\ell)$ of the quadrature points against $B$; $\tau = 0.01, T = 10^4, \epsilon = 10^{-10}$.*

Figure 4.2 (a) shows the relative maximum error $\|e\|_\infty$ against the basis $B$ for different fractional orders, $\tau = 0.01, T = 10^4, \epsilon = 10^{-10}$. We can see that better results are obtained than the predicted precision $\epsilon = 10^{-10}$, even for $\alpha = -1.8$ (the fractional derivative of order 1.8). Figure 4.2 (b) shows the total number of the truncated LG quadrature points against the basis $B$. For a fixed $B$, the number of truncated quadrature points increases as $\alpha$ decreases, which is in line with the theoretical prediction (4.7). For a fixed fractional order $\alpha$, the number of the truncated quadrature points increases as $B$ increases, which agrees with (4.10), due to $\mathcal{T}_\ell \approx 2B-1$ when $\ell$ is sufficiently large.

It is reasonable to choose a relatively smaller precision $\epsilon$ and smaller basis $B$ in numerical simulations, since a smaller $\epsilon$ ensures a relatively large $N_\ell$ that guarantees the exponential convergence of the LG quadrature, and a small $B$ ensures a smaller radius of convergence of the quadrature error (4.9).

TABLE 4.2
*The maximum relative error $\|e\|_\infty$ under different precision $\epsilon$, $B = 5$, $\alpha = -0.5, \Delta T = 1$.*

| $\epsilon$ | $\tau = 0.01, T = 10^4$ | | | $\tau = 0.1, T = 10^5$ | | |
|---|---|---|---|---|---|---|
| | $\sum N_\ell$ | $\sum q_\alpha(N_\ell)$ | $\|e\|_\infty$ | $\sum N_\ell$ | $\sum q_\alpha(N_\ell)$ | $\|e\|_\infty$ |
| $10^{-12}$ | 1018 | 384 | 1.1937e-12 | 1203 | 442 | 6.8208e-12 |
| $10^{-10}$ | 848 | 349 | 4.7044e-12 | 1001 | 402 | 6.8191e-12 |
| $10^{-8}$ | 680 | 313 | 4.7044e-12 | 799 | 361 | 6.5369e-12 |
| $10^{-6}$ | 512 | 271 | 3.1858e-09 | 604 | 312 | 1.2450e-10 |
| $10^{-5}$ | 427 | 243 | 2.1577e-06 | 503 | 281 | 1.8889e-08 |
| $10^{-4}$ | 340 | 215 | 2.1577e-06 | 404 | 252 | 2.3317e-07 |

Finally in this section, we compare the truncated LG quadrature with the multipole approximation proposed in [2]. As the key idea of the existing fast methods aforementioned is to seek a sum-of-exponentials of the form $\sum \omega_j e^{-\lambda_j t}$ to approximate the kernel function $k_\alpha(t)$, we compare the accuracy of the sum-of-exponentials from the LG quadrature and the multipole approximation. The relative pointwise errors $e_n = |k_\alpha(t_n) - \sum \omega_j e^{-\lambda_j t_n}|/|k_\alpha(t_n)|$ are shown in Figure 4.3, where we set the precision $\epsilon = 10^{-10}$, $B = 5$, and $\Delta T = \tau$ when the LG quadrature is applied, and the precision in [2] is set to be $10^{-14}$. For $\alpha = 0.5$, the two methods achieve similar accuracy with, respectively, $P = 832$ and $Q = 777$ quadrature points for the multipole approximation and the truncated LG quadrature, see Figure 4.3 (a). Figure 4.3 (b) shows the pointwise errors for $\alpha = -0.5$, the truncated LG quadrature shows a slightly better approximation. For other fractional orders $\alpha \in (-1, 1)$, the truncated LG quadrature is competitive with the multipole approximation both in accuracy and the computational cost. These results are not shown here.

**5. Numerical examples and applications.** In this section, two examples are presented to verify the effectiveness of the present fast method when it is applied to solve FDEs. All the algorithms are implemented using MATLAB 2016a, which were run in a 3.40 GHz PC having 16GB RAM and Windows 7 operating system.

EXAMPLE 5.1. *Consider the following scaler FDE*

$$(5.1) \qquad {}_C D_{0,t}^\alpha u(t) = -Au(t) + F(u,t), \quad u(0) = u_0, \quad t \in (0, T],$$

*where $0 < \alpha \le 1$ and $A \ge 0$, and ${}_C D_{0,t}^\alpha$ is the Caputo fractional operator, which satisfies ${}_C D_{0,t}^\alpha u(t) = k_\alpha * u(t) - u(0)t^{-\alpha}/\Gamma(1 - \alpha)$, see [29].*

(a) $\alpha = 0.5, \tau = 0.0001$.　　　　(b) $\alpha = -0.5, \tau = 0.0001$.

Fig. 4.3. *Comparison of the truncated LG quadrature and the multipole approximation [2].*

We present our fast numerical method for (5.1) as follows: For a given length $\Delta T = n_0\tau$, find $U_n$ for $n > n_0$ such that

(5.2) $$ {}_F D_{\Delta T,\tau}^{(-\alpha,n,m)} U - u_0 t_n^{-\alpha}/\Gamma(1-\alpha) = -AU_n + F(U_n, t_n), $$

where ${}_F D_{\Delta T,\tau}^{(\alpha,n,m)}$ is defined by (3.22). The application of the direct convolution method means here that ${}_F D_{\Delta T,\tau}^{(-\alpha,n,m)} U$ in (5.2) is replaced by $D_{\Delta T,\tau}^{(-\alpha,n,m)} U$. The Newton method is used to solve the nonlinear system (5.2). The direct method is applied to obtain $U_k (\le k \le m)$, and the stating values $U_k (1 \le k \le m)$ are obtained using the direct method with one correction term and a smaller stepsize $\tau^{-2}$.

The following two cases are considered in this example.

Case I: For the linear case of $F = 0$, the exact solution of (5.1) is

$$ u(t) = E_\alpha(-At^\alpha), $$

where $E_\alpha(t)$ is the Mittag–Leffler function defined by $E_\alpha(t) = \sum_{k=0}^{\infty} \frac{t^k}{\Gamma(k\alpha+1)}$.

Case II: Let $F = u(1 - u^2)$ and the initial condition is taken as $u_0 = 1$.

The maximum error is defined by

$$ \|e\|_\infty = \max_{0 \le n \le T/\tau} |e_n|, \quad e_n = u(t_n) - U_n, T = 40. $$

In this example, we always choose the memory length $\Delta T = 0.5$, the basis $B = 5$, the truncation number $q_\alpha(N_\ell)$ in (3.20) is calculated by (4.7) and $N_\ell$ is determined by (4.10) under the precision $\epsilon = 10^{-10}$, and $A = 1$. We will reset these parameters if needed. The fast method based on quadratic interpolation (3.5) is applied in this example if there is no further illustration.

The purpose of Case I is to check the effectiveness of the present fast method for non-smooth solutions. We demonstrate that adding correction terms improves the accuracy of the numerical solutions significantly. We first let $\alpha = 0.8$, the maximum error and the error at $t = 40$ are shown in Tables 5.1 and 5.2, respectively. We can see that the expected convergence rate $3 - \alpha$ is almost achieved by adding two or three correction terms, and the numerical solution at the final time is much more accurate than that near the origin. This phenomenon can be observed from the existing time-stepping methods for time-fractional differential equations.

14

As the fractional order decreases, the singularity of the analytical solution becomes stronger. Tables 5.3 and 5.4, respectively, show the maximum errors and the errors at $t = 40$ for $\alpha = 0.1$. We observe that the maximum error, which occurs near the origin, almost does not improve, though the step size decreases. We find that as the number of correction terms increases, the accuracy of the numerical solutions increases significantly. We need almost 30 correction terms to derive the expected global convergence rate $3 - \alpha$ for $\alpha = 0.1$, which is impossible to achieve by double precision. The fact is a small number of correction terms is enough to yield satisfactory numerical solutions, which makes the present method more practical. We refer readers to [44] for more explanations and numerical results associated with this phenomenon.

TABLE 5.1

*The maximum error $\|e\|_\infty$ of the method (5.2), Case I, $\sigma_k = k\alpha$, $\alpha = 0.8$, $T = 40$, $B = 5$.*

| $\tau$ | $m = 0$ | Order | $m = 1$ | Order | $m = 2$ | Order | $m = 3$ | Order |
|---|---|---|---|---|---|---|---|---|
| $2^{-5}$ | 2.7966e-3 | | 1.5674e-3 | | 3.0233e-5 | | 3.4818e-5 | |
| $2^{-6}$ | 1.7545e-3 | 0.6726 | 5.4427e-4 | 1.5260 | 5.9024e-6 | 2.3568 | 7.7701e-6 | 2.1638 |
| $2^{-7}$ | 1.0604e-3 | 0.7264 | 1.8498e-4 | 1.5569 | 1.6683e-6 | 1.8229 | 1.7199e-6 | 2.1756 |
| $2^{-8}$ | 6.2689e-4 | 0.7584 | 6.2086e-5 | 1.5751 | 4.2485e-7 | 1.9733 | 3.7854e-7 | 2.1838 |
| $2^{-9}$ | 3.6575e-4 | 0.7773 | 2.0686e-5 | 1.5856 | 1.0205e-7 | 2.0576 | 8.3007e-8 | 2.1892 |

TABLE 5.2

*The absolute error $|e_n|$ of the method (5.2) at $t = 40$, Case I, $\sigma_k = k\alpha$, $\alpha = 0.8$, $B = 5$.*

| $\tau$ | $m = 0$ | Order | $m = 1$ | Order | $m = 2$ | Order | $m = 3$ | Order |
|---|---|---|---|---|---|---|---|---|
| $2^{-5}$ | 5.8890e-7 | | 1.9977e-7 | | 1.7531e-7 | | 4.8407e-6 | |
| $2^{-6}$ | 3.0861e-7 | 0.9322 | 5.8512e-8 | 1.7715 | 3.8771e-8 | 2.1769 | 1.1237e-6 | 2.1069 |
| $2^{-7}$ | 1.5741e-7 | 0.9713 | 1.6977e-8 | 1.7852 | 8.5112e-9 | 2.1875 | 2.5398e-7 | 2.1455 |
| $2^{-8}$ | 7.9372e-8 | 0.9878 | 4.8985e-9 | 1.7931 | 1.8745e-9 | 2.1829 | 5.6424e-8 | 2.1703 |
| $2^{-9}$ | 3.9821e-8 | 0.9951 | 1.4075e-9 | 1.7992 | 4.3619e-10 | 2.1035 | 1.2218e-8 | 2.2073 |

TABLE 5.3

*The maximum error $\|e\|_\infty$ of the method (5.2), Case I, $\sigma_k = k\alpha$, $\alpha = 0.1$, $T = 40$, $B = 5$.*

| $\tau$ | $m = 0$ | Order | $m = 1$ | Order | $m = 3$ | Order | $m = 5$ | Order |
|---|---|---|---|---|---|---|---|---|
| $2^{-5}$ | 2.5872e-3 | | 1.5852e-3 | | 2.1191e-5 | | 2.1191e-5 | |
| $2^{-6}$ | 2.5323e-3 | 0.0310 | 1.5017e-3 | 0.0781 | 9.6408e-6 | 1.1362 | 9.6408e-6 | 1.1362 |
| $2^{-7}$ | 2.4750e-3 | 0.0330 | 1.4177e-3 | 0.0831 | 4.3479e-6 | 1.1488 | 4.3479e-6 | 1.1488 |
| $2^{-8}$ | 2.4148e-3 | 0.0355 | 1.3338e-3 | 0.0880 | 1.9444e-6 | 1.1610 | 1.9444e-6 | 1.1610 |
| $2^{-9}$ | 2.3516e-3 | 0.0383 | 1.2507e-3 | 0.0928 | 1.2279e-6 | 0.6631 | 8.6285e-7 | 1.1721 |

TABLE 5.4

*The absolute error $|e_n|$ of the method (5.2) at $t = 40$, Case I, $\sigma_k = k\alpha$, $\alpha = 0.1$, $B = 5$.*

| $\tau$ | $m = 0$ | Order | $m = 1$ | Order | $m = 3$ | Order | $m = 5$ | Order |
|---|---|---|---|---|---|---|---|---|
| $2^{-5}$ | 6.4561e-6 | | 6.6140e-7 | | 1.0065e-8 | | 4.1985e-9 | |
| $2^{-6}$ | 3.2045e-6 | 1.0106 | 3.1760e-7 | 1.0583 | 3.8428e-9 | 1.3891 | 1.0230e-9 | 2.0371 |
| $2^{-7}$ | 1.5919e-6 | 1.0094 | 1.5234e-7 | 1.0599 | 1.5471e-9 | 1.3126 | 2.5282e-10 | 2.0166 |
| $2^{-8}$ | 7.9116e-7 | 1.0087 | 7.2987e-8 | 1.0616 | 6.4397e-10 | 1.2645 | 6.4081e-11 | 1.9801 |
| $2^{-9}$ | 3.9331e-7 | 1.0083 | 3.4930e-8 | 1.0632 | 2.7313e-10 | 1.2374 | 1.6914e-11 | 1.9217 |

Table 5.5 shows the difference $\eta = \max_{0 \leq n \leq T/\tau} |U_D^n - U_F^n|$, where $U_F^n$ are numerical solutions derived from the fast method under the precision $\epsilon$, and $U_D^n$ are numerical solutions from the direct method. We can see that the difference $\eta$ is independent of the stepsize $\tau$, which confirms (4.13). Table 5.5 also shows that the accuracy of the present fast calculation outperforms the predicted accuracy $\epsilon$, see (4.13).

TABLE 5.5

*The difference of the numerical solutions between the fast method and the direct method, Case I, $m = 0$, $\alpha = 0.1$, $T = 40$, $B = 5$.*

| $\epsilon$ | $\tau = 2^{-5}$ | $\tau = 2^{-6}$ | $\tau = 2^{-7}$ | $\tau = 2^{-8}$ | $\tau = 2^{-9}$ |
|---|---|---|---|---|---|
| $10^{-12}$ | 2.8255e-13 | 2.7850e-13 | 2.7167e-13 | 1.1102e-15 | 1.4988e-15 |
| $10^{-10}$ | 2.8239e-13 | 2.7839e-13 | 2.7162e-13 | 5.6066e-15 | 2.1094e-15 |
| $10^{-8}$ | 2.8116e-13 | 3.8219e-13 | 6.6397e-13 | 2.2093e-14 | 8.2682e-13 |
| $10^{-6}$ | 7.5677e-11 | 7.2366e-11 | 8.0116e-10 | 2.6198e-11 | 8.0352e-12 |
| $10^{-5}$ | 3.2916e-11 | 1.9709e-10 | 1.2477e-10 | 2.6461e-10 | 2.2769e-10 |
| $10^{-4}$ | 3.0868e-09 | 9.4674e-09 | 2.4896e-09 | 1.5458e-08 | 4.3178e-08 |

In Table 5.6, we compare the present fast method (5.2) based on the linear interpolation with the graded mesh method in [33], in which the nonuniform grid points are given by $t_j = (j\tau)^r$, $0 \leq j \leq 1/\tau$. We can see that the fast method with correction terms is competitive with the graded mesh method for $\alpha = 0.5$. In [33], the authors obtained the optimal grid mesh $t_j = (j\tau)^{(2-\alpha)/\alpha}$ to achieve the global $(2 - \alpha)$-order accuracy, which works well when $\alpha$ is relatively large and the expected convergence rate is achieved. For the correction method, too many correction terms may harm the accuracy, but a few number of correction terms can achieve highly accurate numerical solutions, which is not investigated here, readers can refer to [9, 21, 44] for more discussion. Note that for values of $\alpha$ close to zero, the correction method is even more effective than the optimal graded mesh approach in [33].

TABLE 5.6

*Comparison of the present method with the graded mesh method in [33] based on linear interpolation, Case I, $\sigma_k = k\alpha$, $\alpha = 0.5$, $B = 5$. The errors $\|e\|_\infty$ are displayed.*

The graded mesh method [33] with grid points $t_j = (j\tau)^r$.

| $\tau$ | $r = 1$ | Order | $r = 3/2$ | Order | $r = 3$ | Order | $r = 6$ | Order |
|---|---|---|---|---|---|---|---|---|
| $2^{-5}$ | 3.6491e-2 | | 1.6839e-2 | | 2.6568e-3 | | 2.5431e-3 | |
| $2^{-6}$ | 2.7048e-2 | 0.4320 | 1.0288e-2 | 0.7108 | 1.0077e-3 | 1.3986 | 9.2945e-4 | 1.4521 |
| $2^{-7}$ | 1.9772e-2 | 0.4521 | 6.2162e-3 | 0.7268 | 3.7277e-4 | 1.4347 | 3.3644e-4 | 1.4660 |
| $2^{-8}$ | 1.4312e-2 | 0.4663 | 3.7315e-3 | 0.7363 | 1.3582e-4 | 1.4566 | 1.2096e-4 | 1.4759 |
| $2^{-9}$ | 1.0288e-2 | 0.4762 | 2.2313e-3 | 0.7419 | 4.9041e-5 | 1.4696 | 4.3280e-5 | 1.4827 |

The fast method based on linear interpolation with correction terms

| $\tau$ | $m = 1$ | Order | $m = 2$ | Order | $m = 3$ | Order | $m = 4$ | Order |
|---|---|---|---|---|---|---|---|---|
| $2^{-5}$ | 5.6366e-3 | | 2.1893e-4 | | 1.4976e-4 | | 7.2351e-5 | |
| $2^{-6}$ | 3.1659e-3 | 0.8322 | 1.0737e-4 | 1.0279 | 6.9115e-5 | 1.1156 | 4.4890e-5 | 0.6886 |
| $2^{-7}$ | 1.7231e-3 | 0.8777 | 4.9196e-5 | 1.1260 | 2.9276e-5 | 1.2393 | 2.2604e-5 | 0.9898 |
| $2^{-8}$ | 9.1600e-4 | 0.9116 | 2.1427e-5 | 1.1991 | 1.1719e-5 | 1.3208 | 1.0105e-5 | 1.1614 |
| $2^{-9}$ | 4.7862e-4 | 0.9364 | 8.9763e-6 | 1.2552 | 4.5167e-6 | 1.3756 | 4.1916e-6 | 1.2696 |

Next, we implement the present fast method for a longer time computation and compare it with the direct convolution method. We show in Figure 5.1(a) numerical solutions for $t \in [0, T], T = 10000$, where we choose $\alpha = 0.1, 0.5, 0.9$, and the time stepsize $\tau = 0.01$. In Figure 5.1(b), we plot the computational time of the fast method and the direct method with two correction terms applied. It shows that the

16

computational cost of the fast convolution almost achieves linear complexity, which is much less than that of the direct convolution with $O(n_T^2)$ operations when $n_T$ is sufficiently large, where $n_0 = \Delta T/\tau$ and $n_T = T/\tau$.



(a) Numerical solutions.

(b) Computational time.

FIG. 5.1. *Numerical solutions and the computational time of the two different convolutions for Case II, $\tau = 0.01, B = 5, m = 2, \sigma_k = k\alpha$.*

Next, we consider a system of FODEs each having a different fractional index.

EXAMPLE 5.2. *Consider the following system of fractional differential equations each with different fractional index*

$$(5.3) \quad \begin{cases} _CD_{0,t}^{\alpha_1}u(t) = w + (v - c_1)u, \\ _CD_{0,t}^{\alpha_2}v(t) = 1 - c_2v - u^2, \\ _CD_{0,t}^{\alpha_3}w(t) = -u - c_3w, \end{cases}$$

*where $0 < \alpha_k \le 1$ $(k = 1, 2, 3)$, $c_1, c_2$ and $c_3$ are positive constants, $c_2 > 1/2$.*

Let $U_n, V_n$ and $W_n$ be the approximate solutions of $u(t_n), v(t_n)$ and $w(t_n)$, respectively. Then the fully discrete scheme for the system (5.3) is given by:

$$(5.4) \quad \begin{cases} _FD_{\Delta T,\tau}^{(-\alpha_1,n,m)}U - u_0t_n^{-\alpha_1}/\Gamma(1-\alpha_1) = W_n + (V_n - c_1)U_n, \\ _FD_{\Delta T,\tau}^{(-\alpha_2,n,m)}V - v_0t_n^{-\alpha_2}/\Gamma(1-\alpha_2) = 1 - c_2V_n - U_n^2, \\ _FD_{\Delta T,\tau}^{(-\alpha_3,n,m)}W - w_0t_n^{-\alpha_3}/\Gamma(1-\alpha_3) = -U_n - c_3W_n, \end{cases}$$

where $n \ge 2$ and $_FD_{\Delta T,\tau}^{(\alpha,n,m)}U$ is defined by (3.22). Because we use quadratic interpolation, $U_k, V_k, W_k$ for $k = 1, 2$ need to be known, which can be derived using the known methods with smaller stepsize. Here we use the L1 method with one correction term and smaller stepsize $\tau^2$ to derive these values. We take $\Delta T = \tau$ and $B = 5$ in numerical simulations of this example.

If $\alpha_1 = \alpha_2 = \alpha_3 = \alpha$, then (5.3) is just the fractional Lorenz system [36]. It has been proved that the fractional Lorenz system (5.3) is dissipative [36] and has an absorbing set defined by a ball $B(0, \sqrt{a/b} + \hat{\epsilon})$, where $a = 1/2$ and $b = \min\{c_1, c_2 - 1/2, c_3\}$. We take $c_1 = 1/4, c_2 = 1, c_3 = 1/4$ and the initial conditions as those in [36], i.e., $U_0 = u(0) = 2, V_0 = v(0) = 0.9, W_0 = w(0) = 0.2$, the time stepsize is taken as $\tau = 0.01$. We compute numerical solutions for $t \in [0, 1000]$, which is much larger than that in [36]. Numerical solutions for $\alpha = 0.9$ are shown in Figure 5.2. We can easily find that $U_n^2 + V_n^2 + W_n^2 < 2$, which means $(U_n, V_n, W_n) \in B(0, \sqrt{2})$. For other fractional orders $\alpha_k = \alpha$, we have similar results, see also [36]. Next, we

17

choose different $\alpha_1, \alpha_2$, and $\alpha_3$, and exhibit the numerical solutions in Figure 5.3 for $(\alpha_1, \alpha_2, \alpha_3) = (0.9, 0.8, 0.7)$ and $(\alpha_1, \alpha_2, \alpha_3) = (0.7, 0.8, 0.9)$. We can see that the numerical solutions $(U_n, V_n, W_n)$ are also in a ball. For other choices of the fractional orders, we have similar results, which are not provided here.



(a)

(b)

(c)

(d)

FIG. 5.2. *Numerical solutions for Example 5.2, $\tau = 0.01, m = 2, \sigma_k = k\alpha, \alpha = 0.9$.*



(a) $(\alpha_1, \alpha_2, \alpha_3) = (0.9, 0.8, 0.7)$.　　　(b) $(\alpha_1, \alpha_2, \alpha_3) = (0.7, 0.8, 0.9)$.

FIG. 5.3. *Numerical solutions for Example 5.2, $t \in [0, 1000], \tau = 0.01, m = 0$.*

**6. Conclusion and discussion.** We propose a unified fast memory-saving time-stepping method for both fractional integral and derivative operators, which is applied to solve FDEs. We generalize the fast convolution in [22] for fractional operators, while we use the (truncated) LG quadrature to discretize the kernel in the fractional operators instead of the trapezoidal rule in [22]. We also introduced correction terms

in the fast method such that the non-smooth solutions of the considered FDEs can be resolved accurately. The present fast method has $O(n_0 + \sum_\ell^L q_\alpha(N_\ell))$ active memory and $O(n_0 n_T + (n_T - n_0) \sum_\ell^L q_\alpha(N_\ell))$ operations. If a suitable memory length $\Delta T$ and a very large basis $B$ are applied, then the present fast method is similar to that in [14, 18], that is, the kernel $k_\alpha(t)$ is approximated via the truncated LG quadrature for all $t \in [\Delta T, T]$.

We present the error analysis of the present fast method. For a given precision, a criteria on how to choose the parameters used in the fast method is given explicitly, which works very well in numerical simulations.

We considered the fast method for the fractional operator of order $\alpha < 1$. In fact, our method can be extended to the fractional integral of order greater than one. For example, for $\alpha \in [1, 2)$, Eq. (1.4) still holds in the sense of the finite part integral, that is, Eq. (1.4) is equivalent to $k_\alpha(t) = \frac{t}{\Gamma(\alpha)\Gamma(2-\alpha)} \int_0^\infty \lambda^{1-\alpha} e^{-t\lambda} \, d\lambda$. In such a case, a coupled ODEs are needed to be resolved instead of one for $\alpha < 1$.

In the future, we will apply the fast method to solve a large system of time-fractional PDEs, perhaps in three spatial dimensions. We will also explore how to efficiently calculate the discrete convolution $\sum_{k=0}^n \omega_{n-k} u(t_k)$, where the quadrature weights $\{\omega_n\}$ are not from the interpolation as done in the present work, but from the generating functions, see, [21].

**Appendix A. Proofs.** Proof of Theorem 4.1.

*Proof.* We follow the proof of Theorem 2.2 in [39]. We first expand $g(\lambda) = e^{-t\lambda}$ in terms of the Laguerre polynomials, i.e., $g(\lambda) = \sum_{n=0}^\infty a_n L_n^{(\alpha)}(\lambda)$, where

$$a_n = \frac{\Gamma(n+1)}{\Gamma(n+1+\alpha)} \int_0^\infty \lambda^\alpha e^{-\lambda} g(\lambda) \, d\lambda = \frac{t^n}{(t+1)^{n+1+\alpha}}.$$

The following property will be used, see [39],

$$\left| Q_N^\alpha [L_n^{(\alpha)}] \right| \leq \begin{cases} 2\Gamma(1+\alpha), & -1 < \alpha \leq 0, \\ 2^{1+\alpha} \Gamma(n+1+\alpha)/\Gamma(n+1), & \alpha > 0. \end{cases}$$

With the above two equations and $I^\alpha[e^{-t\lambda}] - Q_N^\alpha[e^{-t\lambda}] = \sum_{n=2N}^\infty a_n Q_N^\alpha[L_n^{(\alpha)}]$ gives

(A.1) $\quad \left| I^\alpha[e^{-t\lambda}] - Q_N^\alpha[e^{-t\lambda}] \right| \leq \begin{cases} 2^{1+\alpha} \Gamma(1+\alpha) \sum\limits_{n=2N}^\infty a_n, & -1 < \alpha \leq 0, \\ c_\alpha 2^{1+\alpha} \sum\limits_{n=2N}^\infty n^\alpha a_n, & \alpha > 0. \end{cases}$

Let $q = t/(1+t)$. Then we have $\sum_{n=2N}^\infty a_n = (1+t)^{-\alpha} q^{2N}$ for $-1 < \alpha \leq 0$. For $\alpha > 0$, one hase

$$\sum_{n=2N}^\infty n^\alpha a_n = (1+t)^{-\alpha} q^{2N} (2N)^\alpha \sum_{n=0}^\infty q^n \left(1 + \frac{n}{2N}\right)^\alpha \leq C_{\alpha,t} 2^\alpha (1+t)^{1-\alpha} q^{2N} N^\alpha,$$

where $C_{\alpha,t} = \sum_{n=0}^\infty (n+2)^\alpha q^{2Nn}$ is used. With the above equation and (A.1) yields (4.4) for $T = 1$. Using the following relation

$$\left| I^\alpha[T, e^{-t\lambda}] - Q_N^\alpha[T, e^{-t\lambda}] \right| = T^{-\alpha-1} \left| I^\alpha[e^{-(t/T)\lambda}] - Q_N^\alpha[e^{-(t/T)\lambda}] \right|$$

19

leads to (4.4) for any $T > 0$. The proof is complete. ◻

Proof of Theorem 4.2.

*Proof.* The following results can be found in [12],

$$(A.2) \qquad \lambda_j < \frac{2j + \alpha + 3}{2N + \alpha + 3}\Big(2j + \alpha + 3 + \sqrt{(2j + \alpha + 3)^2 + 0.25 - \alpha^2}\Big),$$

$$(A.3) \qquad \lambda_j > \frac{2(J_{\alpha,j}/2)^2}{2N + \alpha + 3}, \quad J_{\alpha,j} = \pi(j + 3/4 + \alpha/2) + O(j^{-1}),$$

where (A.3) holds when $j$ is sufficiently large. For a sufficiently large $N$, the Laguerre polynomial satisfies (see (7.14) in [32])

$$(A.4) \qquad \big|L_N^{(\alpha)}(x)\big| \approx \pi^{-1/2}(Nx)^{-1/4}e^{x/2}, \quad \forall x \geq 0.$$

With (A.2)–(A.4) and (4.3), we have the following estimate

$$(A.5) \qquad w_j^{(\alpha)} \leq C_\alpha(N + 1)^\alpha \left(\frac{j + 1}{N + 1}\right)^3 e^{-\lambda_j} \leq C(N + 1)^\alpha e^{-\lambda_j},$$

for sufficiently large $j$, where $\lambda_j = \theta_j(j + 1)^2/(N + 1)$ and $\theta_j$ is bounded and approximately between $\pi^2/4$ and 4. The proof is completed. ◻

## REFERENCES

[1] D. BAFFET AND J. S. HESTHAVEN, *High-order accurate adaptive kernel compression time-stepping schemes for fractional differential equations*, J. Sci. Comput., 72 (2017), pp. 1169–1195.

[2] D. BAFFET AND J. S. HESTHAVEN, *A kernel compression scheme for fractional differential equations*, SIAM J. Numer. Anal., 55 (2017), pp. 496–520.

[3] L. BANJAI, M. LÓPEZ-FERNÁNDEZ, AND A. SCHÄDLE, *Fast and oblivious algorithms for dissipative and two-dimensional wave equations*, SIAM J. Numer. Anal., 55 (2017), pp. 621–639.

[4] G. BEYLKIN AND L. MONZN, *Approximation by exponential sums revisited*, Appl. Comput. Harmon. Anal., 28 (2010), pp. 131–149.

[5] L. D'AMORE, A. MURLI, AND M. RIZZARDI, *An extension of the Henrici formula for Laplace transform inversion*, Inverse Problems, 16 (2000), pp. 1441–1456.

[6] W. DENG, *Short memory principle and a predictor-corrector approach for fractional differential equations*, J. Comput. Appl. Math., 206 (2007), pp. 174–188.

[7] K. DIETHELM, *Generalized compound quadrature formulae for finite-part integrals*, IMA J. Numer. Anal., 17 (1997), pp. 479–493.

[8] ——, *The analysis of fractional differential equations*, Springer-Verlag, Berlin, 2010.

[9] K. DIETHELM, J. M. FORD, N. J. FORD, AND M. WEILBEER, *Pitfalls in fast numerical solvers for fractional differential equations*, J. Comput. Appl. Math., 186 (2006), pp. 482–503.

[10] N. J. FORD AND A. C. SIMPSON, *The numerical solution of fractional differential equations: Speed versus accuracy*, Numerical Algorithms, 26 (2001), pp. 333–346.

[11] L. GALEONE AND R. GARRAPPA, *Fractional Adams-Moulton methods*, Math. Comput. Simulation, 79 (2008), pp. 1358–1367.

[12] L. GATTESCHI, *Asymptotics and bounds for the zeros of laguerre polynomials: a survey*, J. Comput. Appl. Math., 144 (2002), pp. 7–27.

[13] H. JIANG, F. LIU, I. TURNER, AND K. BURRAGE, *Analytical solutions for the multi-term time-fractional diffusion-wave/diffusion equations in a finite domain*, Comput. Math. Appl., 64 (2012), pp. 3377–3388.

[14] S. JIANG, J. ZHANG, Q. ZHANG, AND Z. ZHANG, *Fast evaluation of the Caputo fractional derivative and its applications to fractional diffusion equations*, Commun. Comput. Phys., 21 (2017), pp. 650–678.

[15] B. JIN, R. LAZAROV, AND Z. ZHOU, *An analysis of the L1 scheme for the subdiffusion equation with nonsmooth data*, IMA J. Numer. Anal., 36 (2016), pp. 197–221.

[16] C. LI, Q. YI, AND A. CHEN, *Finite difference methods with non-uniform meshes for nonlinear fractional differential equations*, J. Comput. Phys., 316 (2016), pp. 614–631.

[17] C. Li and F. Zeng, *Numerical methods for fractional calculus*, CRC Press, Boca Raton, FL, 2015.

[18] J.-R. Li, *A fast time stepping method for evaluating fractional integrals*, SIAM J. Sci. Comput., 31 (2010), pp. 4696–4714.

[19] Z. Li, Z. Liang, and Y. Yan, *High-order numerical methods for solving time fractional partial differential equations*, J. Sci. Comput., 71 (2017), pp. 785–803.

[20] M. López-Fernández, C. Lubich, and A. Schädle, *Adaptive, fast, and oblivious convolution in evolution equations with memory*, SIAM J. Sci. Comput., 30 (2008), pp. 1015–1037.

[21] C. Lubich, *Discretized fractional calculus*, SIAM J. Math. Anal., 17 (1986), pp. 704–719.

[22] C. Lubich and A. Schädle, *Fast convolution for nonreflecting boundary conditions*, SIAM J. Sci. Comput., 24 (2002), pp. 161–182.

[23] Y. Luchko, *Initial-boundary problems for the generalized multi-term time-fractional diffusion equation*, J. Math. Anal. Appl., 374 (2011), pp. 538–548.

[24] C. Lv and C. Xu, *Error analysis of a high order method for time-fractional diffusion equations*, SIAM J. Sci. Comput., 38 (2016), pp. A2699–A2724.

[25] W. McLean, *Fast summation by interval clustering for an evolution equation with memory*, SIAM J. Sci. Comput., 34 (2012), pp. A3039–A3056.

[26] ———, *Exponential sum approximations for $t^{-\beta}$*, (2016), p. arXiv:1606.00123.

[27] W. McLean and K. Mustapha, *A second-order accurate numerical method for a fractional wave equation*, Numer. Math., 105 (2007), pp. 481–510.

[28] R. Metzler and J. Klafter, *The random walk's guide to anomalous diffusion: a fractional dynamics approach*, Phys. Rep., 339 (2000), p. 77.

[29] I. Podlubny, *Fractional differential equations*, Academic Press, Inc., San Diego, CA, 1999.

[30] S. G. Samko, A. A. Kilbas, and O. I. Marichev, *Fractional integrals and derivatives: Theory and applications*, Gordon and Breach Science Publishers, Yverdon, 1993.

[31] A. Schädle, M. López-Fernández, and C. Lubich, *Fast and oblivious convolution quadrature*, SIAM J. Sci. Comput., 28 (2006), pp. 421–438.

[32] J. Shen, T. Tang, and L.-L. Wang, *Spectral methods*, vol. 41 of Springer Series in Computational Mathematics, Springer, Heidelberg, 2011. Algorithms, analysis and applications.

[33] M. Stynes, E. O'Riordan, and J. L. Gracia, *Error analysis of a finite difference method on graded meshes for a time-fractional diffusion equation*, SIAM J. Numer. Anal., 55 (2017), pp. 1057–1079.

[34] Z.-z. Sun and X. Wu, *A fully discrete difference scheme for a diffusion-wave system*, Appl. Numer. Math., 56 (2006), pp. 193–209.

[35] C.-L. Wang, Z.-Q. Wang, and H.-L. Jia, *An hp-version spectral collocation method for nonlinear Volterra integro-differential equation with weakly singular kernels*, J. Sci. Comput., (2017), pp. 1–32.

[36] D. Wang and A. Xiao, *Dissipativity and contractivity for fractional-order systems*, Nonlinear Dynamics, 80 (2015), pp. 287–294.

[37] J. A. C. Weideman, *Optimizing talbots contours for the inversion of the laplace transform*, SIAM J. Numer. Anal., 44 (2006), pp. 2342–2362.

[38] J. A. C. Weideman and L. N. Trefethen, *Parabolic and hyperbolic contours for computing the bromwich integral*, Math. Comp., 76 (2007), pp. 1341–1356.

[39] S. Xiang, *Asymptotics on Laguerre or Hermite polynomial expansions and their applications in Gauss quadrature*, J. Math. Anal. Appl., 393 (2012), pp. 434–444.

[40] Y. Yan, Z.-Z. Sun, and J. Zhang, *Fast evaluation of the caputo fractional derivative and its applications to fractional diffusion equations: A second-order scheme*, Commun. Comput. Phys., 22 (2017), pp. 1028–1048.

[41] Y. Yu, P. Perdikaris, and G. E. Karniadakis, *Fractional modeling of viscoelasticity in 3d cerebral arteries and aneurysms*, J. Comput. Phys., 323 (2016), pp. 219–242.

[42] M. Zayernouri and A. Matzavinos, *Fractional Adams-Bashforth/Moulton methods: An application to the fractional Keller-Segel chemotaxis system*, J. Comput. Phys., 317 (2016), pp. 1–14.

[43] F. Zeng, Z. Zhang, and G. E. Karniadakis, *Fast difference schemes for solving high-dimensional time-fractional subdiffusion equations*, J. Comput. Phys., 307 (2016), pp. 15–33.

[44] F. Zeng, Z. Zhang, and G. E. Karniadakis, *Second-order numerical methods for multi-term fractional differential equations: Smooth and non-smooth solutions*, Comput. Methods Appl. Mech. Engrg., (2017), p. in press.