

OPTIMAL CONVERGENCE FOR TIME-DEPENDENT STOKES EQUATION: A NEW APPROACH

D. FISHELOV* AND J.-P. CROISILLE*

Abstract. In our book "Navier-Stokes Equations in Planar Domains", Imperial College Press, 2013, we have suggested a fourth-order compact scheme for the Navier-Stokes equations in streamfunction formulation $\partial_t(\Delta\psi) + (\nabla^\perp\psi) \cdot \nabla(\Delta\psi) = \nu\Delta^2\psi$. Here we present a new approach for the analysis of a high-order compact scheme for the Navier-Stokes equations in cases where the convective term vanishes, or in cases where the viscous term dominates the convective term. In these cases the Navier-Stokes equations is replaced by the time-dependent Stokes equation $\partial_t(\Delta\psi) = \nu\Delta^2\psi$. The same type of fourth-order compact schemes that were proposed for the Navier-Stokes equations, may be adopted to the time-dependent Stokes problem. For these methods the truncation error is only of first-order at near-boundary points, but is of fourth order at interior points. We prove that the rate of convergence is actually four, thus the error tends to zero as $O(h^4)$, where h is the size of the mesh.

Keywords: Biharmonic problem, high-order compact scheme, optimal convergence, time dependent Stokes problem, Navier-Stokes equations.

1. INTRODUCTION

The 2D incompressible Navier-Stokes (NS) equations in streamfunction formulation, $\partial_t(\Delta\psi) + (\nabla^\perp\psi) \cdot \nabla(\Delta\psi) = \nu\Delta^2\psi$, play an important role in various areas of physics. The streamfunction formulation of NS equations is due to Lagrange (1768); see [25]. In [5] we have suggested a fourth-order compact scheme for this equation, including a suboptimal error analysis.

In this paper we consider the time dependent Stokes equation in streamfunction formulation

$$(1.1) \quad \partial_t(\Delta\psi) = \nu\Delta^2\psi + f.$$

This equation coincides with the NS equation $\partial_t\Delta\psi(\mathbf{x}, t) + C(\psi) = \nu\Delta^2\psi(\mathbf{x}, t) + f(\mathbf{x}, t)$ when the convective term vanishes. Examples of the so-called Taylor flows and generalizations are given in [14, Chap. 4.6]. The Stokes flow is also an approximate version of the Navier-Stokes system in cases of low Reynolds numbers (and therefore the convective term is small compared to the viscous term). Typically it happens in microfluidics and in creeping flows [33]. It is therefore natural to use (1.1) as an intermediate model for the NS equations [3, 26]. In this paper we consider the fourth-order convergence of a compact high-order scheme suggested for the time-dependent Stokes equation.

The finite difference scheme studied here is based on the Discrete Biharmonic Operator (DBO), which is a compact ("3 point") operator of fourth-order accuracy. Compact high-order schemes for the biharmonic equation were suggested by Stephenson [35] for the two-dimensional biharmonic problem. The DBO studied here may be viewed as a one-dimensional analog of Stephenson's scheme. In our approach, the DBO is obtained as a fourth-order derivative of an interpolating polynomial. This polynomial requires not only values of the function at neighboring points, but also a high-order approximation for the first-order derivative. It turns out that if the first-order derivative is approximated by an Hermitian scheme, then an overall fourth-order accuracy for the fourth-order derivative at interior points is achieved. The DBO operator has a special property that it is compact and its truncation error is fourth order accurate at interior points, while it is only first order accurate at near boundary points. This numerical phenomenon is known in various other contexts. In [19, 20] a hyperbolic system of first order and a parabolic problem were analyzed in the case where extra boundary conditions were given in order to "close" the numerical scheme. It was shown that if the accuracy of the extra boundary conditions is one less compared to the accuracy of the inner scheme, then the overall accuracy of the scheme is determined by the accuracy at inner points. Similarly, in [1] it was proved for a parabolic equation that if the scheme is of order $O(h^\alpha)$ at inner points and of order $O(h^{\alpha-s})$ at near boundary points and if $s = 0, 1$, then the accuracy of the scheme is

Date: October 21, 2021.

¹this is sometimes referred to as a super convergence property

$O(h^\alpha)$. Finally, in [37] the authors consider the wave equation, which is approximated via a summation-by-parts finite difference scheme and is fourth-order accurate at interior points. They show that if the truncation error at interior points is $O(h^4)$ and it is only $O(h^2)$ at near boundary points, then the overall accuracy of the scheme is $O(h^4)$.

In previous papers we studied various compact schemes for the Navier Stokes equations: We describe now the convergence analysis that what was already carried out in previous works. The NS equation in streamfunction form reads

$$(1.2) \quad \partial_t \Delta \psi(\mathbf{x}, t) + C(\psi) = \nu \Delta^2 \psi(\mathbf{x}, t) + f(\mathbf{x}, t),$$

where the convective term is the RHS is $C(\psi) = \nabla^\perp \psi \cdot (\nabla \Delta \psi(\mathbf{x}, t))$. In [2] and in [5] the authors described a fourth-order accurate finite difference scheme for the 2D Navier-Stokes equations. In [3] we have proved convergence for a second-order compact scheme, where the convective term was replaced by a constant coefficients equation. It was shown in [3] that the error in the discrete l_2 norm (including the discrete l_2 norm of the first order derivatives) is bounded by Ch^2 , where h is the mesh size. In [4], we had proven convergence and error bounds for a second-order compact scheme for the full (nonlinear) Navier-Stokes equations. We have shown that the error in the discrete l_2 norm of the first order derivatives is bounded by $Ch^{1.5}$. We don't have yet a convergence proof for the 4th order compact scheme for the NS equations, which we have suggested in [2] and [5]. The convergence analysis carried out in this paper for the Stokes equation may pave out the way for the convergence of a fourth-order scheme for the NS system which was suggested in [5].

The focus of this paper is the way to overcome the first order truncation at near boundary points and prove 4th order error estimate. This numerical "phenomenon" is the main topic of the paper. It is of interest at the computing level, since it gives a direction to better understand "superconvergence" phenomenon observed in scientific computing.

Our objective here is to extend the DBO fourth-order error analysis in the one-dimensional time-independent case, as accomplished in [16, 2], to a discrete-in-space and continuous-in-time approximation of (1.1). The key idea is to represent the error $e(t)$ as an integral over time of the time-evolution operator, which operates on the truncation error. Then, to express $e(t)$ in terms of the DBO error, rather than the truncation error, using the analysis derived in [16, 2]. As a first step we establish the fourth-order convergence for the model problem $\partial_{xxt}u = \partial_x^4 u + f$. Then, we proceed to the convergence analysis for (1.1).

The novelty in this approach is that the exact solution of the system of ordinary differential equations is handled delicately, while postponing to later stages of the proof the application of energy estimates. This is done by dealing with the orders of magnitude of the matrices involved, using their known structure, given in [7, 2, 6]. This approach paves the way for the case of the full NS equations, which may however require an additional analysis which is related to the non-linear convective term. In [4] an error estimate was derived for a second order compact scheme for the NS equations in streamfunction formulation. The error due to the convective term was decomposed to a sum of several terms that may be bounded in terms of h and in terms of the error. Finally, with convergence analysis was concluded using the Gronwall's inequality. We have shown in [4] a sub-optimal error estimate of $O(h^{1.5})$. For the fourth-order compact scheme suggested in [5] a similar decomposition of the error due to the convective term may be applied as well, however the challenge is to prove optimal $O(h^4)$ convergence. This is deferred to a future work.

The quest for fourth order accuracy and beyond for the incompressible Navier-Stokes equations has been an active area of research since more than forty years. In the context of finite differences, we refer to [15, 11, 27, 23]. In the paper of E and Liu [15], the authors suggested a fourth-order essentially compact scheme for the Navier-Stokes equations. The scheme in [15] is designed for the NS equations in vorticity-streamfunction formulation, where the vorticity is evolved in time and the streamfunction is constructed from the vorticity in a separate stage. The boundary condition on the vorticity was chosen via a Briley's formula (see [10], [38]). Its stability was proved in [38]. Moreover, since Briley's formula is formally only third order accurate, a new vorticity boundary formula, formally fourth order accurate on the boundary, has been proposed and analyzed in [39, 40, 28]. Stability and convergence estimates have been provided for the full PDE system. The scheme in [5] uses pure streamfunction formulation, where the streamfunction is evolved directly in time. The boundary conditions are applied solely on the streamfunction and there is no need for vorticity boundary conditions. A convergence analysis for the biharmonic equation in two dimensions was carried out in [30]. It was proved that the standard 13-point finite difference scheme converges to the exact solution with $O(h^2)$ error in case the exact solution u is in $H^4(\Omega) \cap H_0^2(\Omega)$.

In [24] a second-order scheme for the NS equations in streamfunction formulation was designed; numerical results for the driven cavity were performed. In the context of finite element methods (continuous or discontinuous) a primitive formulation (velocity-pressure) of NS equations is invoked. Then, convergence analysis is often available. The precise design of the global discrete operators resulting from the variational form requires a fine tuning, including nontrivial penalization operators. It is possible to use these schemes for the Stokes equation by dropping the convective term. We refer to [13] and the references therein.

The outline of the paper is the following. In Section 2 we introduce the notation of the compact scheme involved on the model equation

$$(1.3) \quad \partial_{xxt}u = \partial_x^4 u + f.$$

Vectors and matrices, which are associated with the compact difference operators are introduced as well. In Section 3 the fourth order convergence for (1.3) is proved. In Section 4, it is shown that a similar result extends to the two dimensional time dependent Stokes problem. In Section 5 optimal bounds, given on norms of several matrices of interest which operate on the truncation vector, are derived. These bounds are needed for the convergence proofs in Sections 3 and 4. Truncation errors, needed for the analysis conducted in Sections 3 and 4, are derived in Section 6. Finally, numerical results are presented in Section 7. They assess the fourth-order accuracy of the scheme.

2. A COMPACT SCHEME FOR THE DIFFUSION EQUATION $\partial_{xxt}u = \partial_x^4 u + f$

In [32] the diffusion equation (1.3) is considered as a model for (1.1). The goal was to analyze the constraints on the solution $u(x, t)$ of (1.3) imposed by the rigid wall conditions $u(0, t) = 0$ and $\partial_x u(0, t) = 0$. The absence of boundary conditions on $\partial_{xx}u(0, t)$ corresponds to the unknown vorticity at the wall for (1.1). In this section, we study equation (1.3) as a model for our convergence analysis. Similar ideas will be applied to (1.1) in Section 4.

2.1. Grid function notation. We consider the problem

$$(2.1) \quad \begin{cases} \partial_{xxt}u = \partial_x^4 u + f, & x \in (0, 1), \quad t \geq 0, \\ u(0, t) = u(1, t) = 0, \quad \partial_x u(0, t) = \partial_x u(1, t) = 0; & u(x, 0) = g(x). \end{cases}$$

In what follows, we assume that the data, $f(x, t)$, $g(x)$ and the solution $u(x, t)$, are regular functions of $(x, t) \in [0, 1] \times [0, t_f]$, where $0 < t_f$ is a fixed time.

The discrete setting in approximating (2.1) is based on the analysis conducted in [2, Part II]. If a uniform grid $x_j = jh$, $h = 1/N$, $j = 0, 1, \dots, N$ is laid out on $[0, 1]$, then we denote by $\mathbf{v} = [\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_{N-1}, \mathbf{v}_N]^\top$ the grid function, defined at the points x_j , $j = 0, 1, \dots, N$ only. Using grid functions, it is understood that $h = 1/N$ becomes a small parameter, which tends to zero. Equivalently, the number of grid points $N = 1/h$ tends to $+\infty$. Finite difference operators act naturally on grid functions.

Throughout this paper we assume that the grid functions \mathbf{v} satisfy $\mathbf{v}_0 = \mathbf{v}_N = 0$. The space of such grid functions is called $l_{h,0}^2$, which is equipped with the norm

$$(2.2) \quad |\mathbf{v}|_h = \sqrt{h \sum_{j=1}^{N-1} |v_j|^2}.$$

To a given function $u(x)$, $0 \leq x \leq 1$, we denote by u^* the grid function defined by

$$(2.3) \quad u_j^* = u(x_j), \quad 1 \leq j \leq N-1,$$

where $N \geq 1$.

Now, we approximate the solution of Equation (2.1), on the uniform grid $x_j = j/N$, $j = 0, 1, \dots, N$. The grid function $\mathbf{v}(t) = [\mathbf{v}_j(t)]$, $j = 0, 1, \dots, N$, which serves as an approximation in space of $u(x, t)$, is defined as the solution of

$$(2.4) \quad \begin{cases} (\partial_t \tilde{\delta}_x^2 \mathbf{v})_j = \delta_x^4 \mathbf{v}_j + \mathbf{f}_j, & j = 1, \dots, N-1, \\ \mathbf{v}_0(t) = \mathbf{v}_N(t) = 0, \quad (\tilde{\delta}_x \mathbf{v})_0(t) = (\tilde{\delta}_x \mathbf{v})_N(t) = 0, & \mathbf{v}_j(0) = g(x_j), \quad j = 0, 1, \dots, N. \end{cases}$$

Hence $t \mapsto \mathbf{v}(t)$ satisfies

$$(2.5) \quad \partial_t \tilde{\delta}_x^2 \mathbf{v} = \delta_x^4 \mathbf{v} + \mathbf{f}^*,$$

where \mathbf{f}^* is the grid function defined by $\mathbf{f}_j = f(x_j, t)$. The following finite difference operators are involved.

- The three point Laplacian δ_x^2 and the centered difference δ_x are defined by

$$(2.6) \quad (\delta_x^2 \mathbf{v})_j = \frac{\mathbf{v}_{j+1} + \mathbf{v}_{j-1} - 2\mathbf{v}_j}{h^2}, \quad (\delta_x \mathbf{v})_j = \frac{\mathbf{v}_{j+1} - \mathbf{v}_{j-1}}{2h}, \quad 1 \leq j \leq N-1.$$

- The *Hermitian derivative* is $\tilde{\delta}_x \mathbf{v} \in l_{h,0}^2$, defined by

$$(2.7) \quad \tilde{\delta}_x \mathbf{v}_j = (\sigma_x^{-1} \delta_x) \mathbf{v}_j, \quad j = 1, \dots, N-1,$$

where

$$(2.8) \quad (\sigma_x \mathbf{w})_j = \frac{1}{6}(\mathbf{w}_{j-1} + 4\mathbf{w}_j + \mathbf{w}_{j+1}), \quad j = 1, \dots, N-1.$$

Here we assume that $\mathbf{w}_0 = \mathbf{w}_N = (\tilde{\delta}_x \mathbf{v})_0 = (\tilde{\delta}_x \mathbf{v})_N = 0$.

- The Discrete Biharmonic Operator (DBO) δ_x^4 is defined by

$$(2.9) \quad (\delta_x^4 \mathbf{v})_j = \frac{12}{h^2} \left((\delta_x \tilde{\delta}_x) \mathbf{v} - \delta_x^2 \mathbf{v} \right)_j, \quad j = 1, \dots, N-1.$$

- The operator $\tilde{\delta}_x^2$ is a discrete Laplacian, which has higher accuracy compared with δ_x^2 . It is defined by

$$(2.10) \quad (\tilde{\delta}_x^2 \mathbf{v})_j = 2(\delta_x^2 \mathbf{v})_j - \left((\delta_x \tilde{\delta}_x) \mathbf{v} \right)_j = (\delta_x^2 \mathbf{v})_j - \frac{h^2}{12} (\delta_x^4 \mathbf{v})_j, \quad j = 1, \dots, N-1.$$

The exact solution $u(x, t)$ of (2.1) is associated with the grid function $u^*(t)$, defined by

$$(2.11) \quad u_j^*(t) = u(x_j, t), \quad 1 \leq j \leq N-1.$$

The truncation error for the discrete equation (2.5) is the grid function $t \mapsto \mathbf{r}(t)$, defined by

$$(2.12) \quad \partial_t \tilde{\delta}_x^2 u^* = \delta_x^4 u^* + \mathbf{f}^* + \mathbf{r}.$$

It is expressed in terms of $u(x, t)$ as

$$(2.13) \quad \mathbf{r}(t) = -\left(\delta_x^4 u(\cdot, t)^* - (\partial_x^4 u(\cdot, t))^* \right) + \left(\partial_t (\tilde{\delta}_x^2 u^*(\cdot, t)) - (\partial_t \partial_x^2 u(\cdot, t))^* \right).$$

Define the error $\mathbf{e}(t) = \mathbf{v}(t) - u^*(t)$. Subtracting (2.12) from (2.5) yields

$$(2.14) \quad \partial_t (-\tilde{\delta}_x^2 \mathbf{e}) + \delta_x^4 \mathbf{e} = \mathbf{r}.$$

In Section 3, we will show an "optimal estimate" for $\mathbf{e}(t)$ and its derivative $\tilde{\delta}_x \mathbf{e}(t)$, meaning that $|\mathbf{e}(t)|_h$ and $|\tilde{\delta}_x \mathbf{e}(t)|_h$ tend to zero as $O(h^4)$.

2.2. Vector notation. Assuming that the grid size N is fixed, then grid functions and finite difference operators reduce to vectors and matrices. Matrix analysis (with fixed dimension $N-1$) is therefore the main tool in what follows. We represent (informally) by \leftrightarrow , the correspondence between finite difference operators and matrices and between grid functions and vectors. The matrices involved are T and K , of order $(N-1) \times (N-1)$.

$$(2.15) \quad T_{i,m} = \begin{cases} 2, & m = i \\ -1, & |m - i| = 1 \\ 0, & |m - i| \geq 2 \end{cases}, \quad K_{i,m} = \begin{cases} 0, & m = i \\ 1, & m - i = 1 \\ -1, & m - i = -1 \end{cases}$$

- We have

$$(2.16) \quad (-\delta_x^2) \leftrightarrow \tilde{T} \triangleq T/h^2, \quad \delta_x \leftrightarrow K/2h.$$

- The matrices corresponding to $\tilde{\delta}_x$ and to σ_x are

$$(2.17) \quad \tilde{\delta}_x \leftrightarrow 3P^{-1}K/h, \quad \sigma_x \leftrightarrow P/6,$$

where $P = 6I - T$.

- We have $\delta_x^4 \hookrightarrow B$ where B is the matrix

$$(2.18) \quad B = \frac{6}{h^4} P^{-1} T^2 + \frac{36}{h^4} (V_1 V_1^\top + V_2 V_2^\top),$$

where $V = [V_1, V_2] \in \mathbb{M}_{N-1,2}$ is defined by

$$(2.19) \quad \begin{cases} V_1 = (\alpha - \beta)^{1/2} P^{-1} \left(\frac{\sqrt{2}}{2} e_1 - \frac{\sqrt{2}}{2} e_{N-1} \right) \\ V_2 = (\alpha + \beta)^{1/2} P^{-1} \left(\frac{\sqrt{2}}{2} e_1 - \frac{\sqrt{2}}{2} e_{N-1} \right) \end{cases}, \quad \begin{cases} \alpha = 2(2 - e_1^\top P^{-1} e_1) \\ \beta = 2e_{N-1}^\top P^{-1} e_1. \end{cases}$$

- The matrix D such that $(-\tilde{\delta}_x^2) \hookrightarrow D$ is

$$(2.20) \quad D = \frac{1}{h^2} \left(T + \frac{1}{2} P^{-1} T^2 \right) + \frac{3}{h^2} (V_1 V_1^\top + V_2 V_2^\top) = \tilde{T} + \frac{h^2}{12} B.$$

- For a fixed time $t > 0$, $E(t), R(t) \in \mathbb{R}^{N-1}$ are defined by

$$(2.21) \quad \mathbf{e}(t) \hookrightarrow E(t), \quad \mathbf{r}(t) \hookrightarrow R(t).$$

2.3. The optimal convergence theorem for the discrete biharmonic operator. The DBO (discrete biharmonic operator) δ_x^4 has been introduced in [3], based on [35]. It appears as the main tool to solve the NS equations in streamfunction form. This discrete operator has a fourth order truncation error except at near boundary points, where the truncation drops to first order. Nevertheless, it has been proved in [16, 2] that this does not prevent a fourth order error estimate. A precise statement of this "optimal convergence theorem" for the DBO operator (2.9), both in grid function and vector formulations, is given in Theorems 2.1 and 2.2.

Theorem 2.1 (Optimal convergence theorem for the DBO). *Assume that the vector \mathbf{r} , which contains the truncation errors, satisfy the following estimates:*

$$(2.22) \quad |(\sigma_x \mathbf{r})_j| \leq Ch^4, \quad j = 2, \dots, N-2,$$

$$|(\sigma_x \mathbf{r})_1| \leq Ch, \quad |(\sigma_x \mathbf{r})_{N-1}| \leq Ch.$$

The DBO operator δ_x^4 is invertible and its inverse is denoted by

$$(2.23) \quad \delta_x^{-4} = (\delta_x^4)^{-1}.$$

Then, the vector $\delta_x^{-4} \mathbf{r} = [(\delta_x^{-4} \mathbf{r})_1, \dots, (\delta_x^{-4} \mathbf{r})_{N-1}]^\top$ satisfies

$$(2.24) \quad |(\delta_x^{-4} \mathbf{r})_j| \leq Ch^4, \quad j = 1, \dots, N-1.$$

Therefore, the following max norm estimate holds

$$(2.25) \quad \max_{1 \leq j \leq N-1} |(\delta_x^{-4} \mathbf{r})_j| \leq Ch^4.$$

The following Theorem translates Theorem 2.1 into vector form.

Theorem 2.2 (Vector form of the optimal convergence theorem for the DBO). *Assume that $R \in \mathbb{R}^{N-1}$ satisfies*

$$(2.26) \quad PR = [O(h), O(h^4), \dots, O(h^4), O(h)]^\top,$$

then the vector $B^{-1}R$ satisfies

$$(2.27) \quad |(B^{-1}R)_j| \leq Ch^4, \quad j = 1, \dots, N-1.$$

In addition,

$$(2.28) \quad \begin{cases} |(P^{-1}B^{-1}R)_j| \leq Ch^4, & j = 2, \dots, N-2, \\ |(P^{-1}B^{-1}R)_j| \leq Ch^5, & j = 1, N-1. \end{cases}$$

Let us recall the main steps in the proof presented in [16, 2]. From (2.18) we deduce that the matrix PBP is expressed as

$$(2.29) \quad PBP = \frac{6}{h^4} T^2 P + \frac{36}{h^4} J J^\top, \quad \text{where } J = PV.$$

This may be rewritten as

$$(2.30) \quad PBP = GH^{-1},$$

where

$$(2.31) \quad G = I + 6JJ^T P^{-1} T^{-2}, \quad H = \frac{h^4}{6} P^{-1} T^{-2}.$$

Thus,

$$(2.32) \quad (PBP)^{-1} = HG^{-1}.$$

Therefore,

$$(2.33) \quad B^{-1} = PHG^{-1}P,$$

which results in

$$(2.34) \quad P^{-1}B^{-1}R = HG^{-1}PR.$$

We have shown in [16] and [2, chap. 10.7] that

$$(2.35) \quad G^{-1}PR = [O(h^2), O(h^4), \dots, O(h^4), O(h^2)]^T.$$

Furthermore, the matrix H may be diagonalized using the spectral basis Z^k of the matrix T , where

$$(2.36) \quad Z_j^k = \left(\frac{2}{N}\right)^{1/2} \sin \frac{kj\pi}{N}, \quad 1 \leq k, j \leq N-1.$$

It results that (see [2, Sect. 10.7, pp. 174 sqq])

$$(2.37) \quad H_{ij} = \begin{cases} O(h), & \text{if } 2 \leq i, j \leq N-2 \\ O(h^2), & \text{if } i \in \{1, N-1\} \text{ or } j \in \{1, N-1\} \\ O(h^3), & \text{if } (i, j) \in \{(1, 1), (1, N-1), (N-1, 1), (N-1, N-1)\}. \end{cases}$$

Therefore

$$(2.38) \quad P^{-1}B^{-1}R = H(G^{-1}PR) = [O(h^5), O(h^4), \dots, O(h^4), O(h^5)]^T.$$

Thus, the estimate (2.28) holds. It yields

$$(2.39) \quad B^{-1}R = P[O(h^5), O(h^4), \dots, O(h^4), O(h^5)]^T = [O(h^4), O(h^4), \dots, O(h^4), O(h^4)]^T.$$

Hence, we conclude that (2.27) is valid.

3. CONVERGENCE ANALYSIS FOR THE EQUATION $\partial_{xxt}u = \partial_x^{(4)}u + f$

In this section we state and prove the fourth-order converge of our scheme (2.4) to the exact solution of the diffusion equation (2.1). We have defined the error $\mathbf{e}(t) = \mathbf{v}(t) - u^*(t)$ and obtained that (see Equation (2.14))

$$(3.1) \quad \partial_t(-\tilde{\delta}_x^2 \mathbf{e}) + \delta_x^4 \mathbf{e} = \mathbf{r}.$$

Our goal is to prove the following bound on $\mathbf{e}(t)$, assuming a priori regularity on $u(x, t)$.

Theorem 3.1. *Let $t_f > 0$ be a fixed time.*

(i) *Suppose that u is a solution to the problem (2.1) having spatial derivatives up to order 8, then the error $\mathbf{e}(t)$, satisfying Equation (3.1), is bounded by*

$$(3.2) \quad \max_{0 \leq t \leq t_f} |(-\tilde{\delta}_x^2)^{1/2} \mathbf{e}(t)|_h \leq C(t_f)h^4.$$

(ii) *The following error estimates holds*

$$(3.3) \quad \max_{0 \leq t \leq t_f} |e(t)|_h \leq C(t_f)h^4,$$

and

$$(3.4) \quad \max_{0 \leq t \leq t_f} |\tilde{\delta}_x e(t)|_h \leq C(t_f)h^4.$$

In (3.2)-(3.3)-(3.4), $C(t_f)$ denotes a constant depending only on $\partial_t^{k_0} \partial_x^{k_1} u(x, t)$, $0 \leq k_0 \leq 2$, $4 \leq k_1 \leq 8$ ($x, t \in [0, 1] \times [0, t_f]$).

Proof. (i) Our approach for estimating the error $\mathbf{e}(t)$ relies on vectors and matrices formulations in \mathbb{R}^{N-1} and $\mathbb{M}_{N-1}(\mathbb{R})$, respectively. Consider the vector function $t \in [0, t_f] \mapsto E(t) \in \mathbb{R}^{N-1}$, where $E(t) \hookrightarrow \mathbf{e}(t)$. We have $E(t) = [\mathbf{e}_1(t), \dots, \mathbf{e}_{N-1}(t)]^\top$ at time t . The claim in Theorem 3.1 is equivalent to

$$(3.5) \quad \max_{0 \leq t \leq t_f} \|D^{1/2}E(t)\|_2 \leq C(t_f)h^{3.5}.$$

We invoke the matrices B and D given in (2.18) and (2.20), where the matrix B represents the biharmonic operator δ_x^4 and the matrix

D corresponds to the operator $-\tilde{\delta}_x^2$. In vector form, Equation (2.14) may be written as

$$(3.6) \quad \partial_t(DE(t)) + B E(t) = R(t).$$

The vector $R(t) = [\mathbf{r}_1(t), \dots, \mathbf{r}_{N-1}(t)]^\top \in \mathbb{R}^{N-1}$ is such that $R(t) \hookrightarrow \mathbf{r}(t)$, where $\mathbf{r}(t)$ is the truncation error in (2.13). It is obtained by applying matrices of fixed size $N-1$ to the vectors $U(t) \hookrightarrow [u(x_1, t), \dots, u(x_{N-1}, t)]^\top$ and $U^k(t) \hookrightarrow [\partial_x^{(k)}(x_1, t), \dots, \partial_x^{(k)}(x_{N-1}, t)]^\top$. Due to the a priori regularity hypothesis on $u(x, t)$, the solution of (2.1), the function $t \mapsto R(t)$ is regular in time, with $R(0) = 0$.

Define

$$(3.7) \quad F(t) = D^{1/2}E(t),$$

or equivalently

$$(3.8) \quad E(t) = D^{-1/2}F(t),$$

then

$$(3.9) \quad \partial_t(D^{1/2}F(t)) + BD^{-1/2}F(t) = R(t).$$

Multiplying both sides by $D^{-1/2}$ we have

$$(3.10) \quad \partial_t F(t) + D^{-1/2}BD^{-1/2}F(t) = D^{-1/2}R(t).$$

Let \tilde{B} be the symmetric matrix defined by

$$(3.11) \quad \tilde{B} = D^{-1/2}BD^{-1/2},$$

then, the time dependent equation for $F(t)$ is

$$(3.12) \quad \partial_t F(t) + \tilde{B}F(t) = D^{-1/2}R(t).$$

Note that the matrix \tilde{B} stands for an approximation of a second order derivative. Since $F(0) = 0$, the Duhamel formula for (3.12) is

$$(3.13) \quad F(t) = \int_0^t e^{-\tilde{B}(t-\rho)} D^{-1/2}R(\rho) d\rho,$$

or equivalently

$$(3.14) \quad F(t) = \int_0^t e^{-\rho\tilde{B}} D^{-1/2}R(t-\rho) d\rho.$$

For $\rho \in [0, t]$, $t \leq T$, we rewrite the integrand of (3.14) as

$$(3.15) \quad e^{-\rho\tilde{B}} D^{-1/2}R(t-\rho) = e^{-\rho\tilde{B}} \tilde{B}^{-1} D^{1/2}R(t-\rho) = e^{-\rho\tilde{B}} \tilde{B} D^{1/2}B^{-1}R(t-\rho).$$

Therefore, (3.14) is expressed as

$$(3.16) \quad F(t) = \int_0^t (e^{-\rho\tilde{B}} \tilde{B}) \left(D^{1/2}B^{-1}R(t-\rho) \right) d\rho.$$

Here the main point is that the integrand involves the error $B^{-1}R$, which is $O(h^4)$ and not the truncation error R . Since \tilde{B} is a symmetric definite positive, it may be diagonalized via an orthogonal matrix Q . We denote the eigenvalues of \tilde{B} by

$$(3.17) \quad 0 < \lambda_1 < \dots < \lambda_{N-1}.$$

Let $\tilde{\Lambda}(\rho)$ be

$$(3.18) \quad \tilde{\Lambda}(\rho) = \text{diag} \{ e^{-\rho\lambda_1} \lambda_1, \dots, e^{-\rho\lambda_{N-1}} \lambda_{N-1} \},$$

where $\lambda_i, i = 1, \dots, N - 1$ are the eigenvalues of $\tilde{B} = D^{-1/2} B D^{-1/2}$. Then,

$$(3.19) \quad e^{-\rho \tilde{B}} \tilde{B} = Q \tilde{\Lambda}(\rho) Q^\top.$$

Inserting Equation (3.19) in (3.16), we have

$$(3.20) \quad F(t) = \int_0^t Q \tilde{\Lambda}(\rho) Q^\top D^{1/2} B^{-1} R(t - \rho) d\rho.$$

Notice that

$$(3.21) \quad \tilde{\Lambda}(\rho) = -\frac{d}{d\rho} \Lambda(\rho),$$

where

$$(3.22) \quad \Lambda(\rho) = \text{diag}\{e^{-\rho\lambda_1}, \dots, e^{-\rho\lambda_{N-1}}\}.$$

Inserting (3.21)-(3.22) into (3.20), we have

$$(3.23) \quad F(t) = -\int_0^t \frac{d}{d\rho} (Q \Lambda(\rho) Q^\top) D^{1/2} B^{-1} R(t - \rho) d\rho.$$

Integration by parts yields

$$(3.24) \quad F(t) = \left[-(Q \Lambda(\rho) Q^\top) D^{1/2} B^{-1} R(t - \rho) \right]_{\rho=0}^t + \int_0^t (Q \Lambda(\rho) Q^\top) D^{1/2} B^{-1} \left(\frac{d}{d\rho} R(t - \rho) \right) d\rho.$$

We decompose $F(t)$ as $F(t) = F^{(1)}(t) + F^{(2)}(t)$, where

$$(3.25) \quad \begin{cases} F^{(1)}(t) = -\left[(Q \Lambda(\rho) Q^\top) D^{1/2} B^{-1} R(t - \rho) \right]_{\rho=0}^t, \\ F^{(2)}(t) = \int_0^t (Q \Lambda(\rho) Q^\top) D^{1/2} B^{-1} \left(\frac{d}{d\rho} R(t - \rho) \right) d\rho. \end{cases}$$

The first term is

$$(3.26) \quad \begin{aligned} F^{(1)}(t) &= (Q \Lambda(0) Q^\top) D^{1/2} B^{-1} R(t) - (Q \Lambda(t) Q^\top) D^{1/2} B^{-1} R(0) \\ &= (Q \Lambda(0) Q^\top) D^{1/2} B^{-1} R(t). \end{aligned}$$

Refer now to Section 6, Corollary 6.6, which shows that $PR(t) = [O(h), O(h^4), \dots, O(h^4), O(h)]^\top$ with explicit representation of each component in terms of $x \mapsto u(\cdot, t)$, $x \mapsto \partial_t^{k_0} \partial_x^{k_1} u(\cdot, t)$, with $0 \leq k_0 \leq 1$, $4 \leq k_1 \leq 8$. Then, it will result from Lemma 5.2 that the components of $D^{1/2} B^{-1} R(t)$ are $O(h^4)$. This gives

$$(3.27) \quad \|D^{1/2} B^{-1} R(t)\|_2 \leq C_1(u, t) h^{3.5},$$

where $C_1(u, t)$ can be expressed as

$$(3.28) \quad C_1(u, t) = C' \max_{\substack{0 \leq k_0 \leq 1 \\ 4 \leq k_1 \leq 8}} \max_{x \in [0, 1]} |\partial_t^{k_0} \partial_x^{k_1} u(x, t)|$$

and C' a universal constant. In addition, $\|Q\|_2 = \|Q^\top\|_2 = 1$ and $\|\Lambda(\rho)\|_2 \leq 1$, for $\rho \geq 0$. Thus, we will have

$$(3.29) \quad \|F^{(1)}(t)\|_2 \leq \|D^{1/2} B^{-1} R(t)\|_2 \leq C_1(u, t) h^{3.5}.$$

We turn now to $F^{(2)}(t)$. We have

$$(3.30) \quad \|F^{(2)}(t)\|_2 \leq \max_{0 \leq \rho \leq t} \|D^{1/2} B^{-1} \frac{dR}{d\rho}(\rho)\|_2.$$

By continuity, there exists $\bar{\rho} \in [0, t]$ such that $\max_{0 \leq \rho \leq t} \|D^{1/2} B^{-1} \frac{dR}{d\rho}(\rho)\|_2 = \|D^{1/2} B^{-1} \frac{dR}{d\rho}(\bar{\rho})\|_2$. Therefore

$$(3.31) \quad \|F^{(2)}\|_2 \leq C_1(\partial_t u, \bar{\rho}) h^{3.5}$$

Combining Equations (3.24), (3.29) and (3.31), we conclude that

$$(3.32) \quad \|F(t)\|_2 \leq \|F^{(1)}(t)\|_2 + \|F^{(2)}(t)\|_2 \leq C(t) h^{3.5},$$

where

$$(3.33) \quad C(t) = C'' \max_{\substack{0 \leq k_0 \leq 2 \\ 4 \leq k_1 \leq 8}} \max_{x \in [0,1]} |\partial_t^{k_0} \partial_x^{k_1} u(x, t)|.$$

Using $F(t) = D^{1/2}E(t)$, we have $\|D^{1/2}E(t)\|_2 \leq C(t)h^{3.5}$. Hence we conclude that

$$(3.34) \quad |(-\tilde{\delta}_x^2)^{1/2}\mathbf{e}(t)|_h \leq C(t)h^4.$$

This gives (3.2), by taking the maximum over $t \in [0, t_f]$.

(ii) The function $t \in \mathbb{R}_+ \mapsto t^{1/2}$ is operator monotone (see [9, chap.5]). This means that for positive definite matrices A and B , we have

$$(3.35) \quad A \succ B \Rightarrow A^{1/2} \succ B^{1/2},$$

where \succ stands for the ordering of positive definite matrices (see [21, chap. 7.7]). By (2.20)

$$(3.36) \quad D = \tilde{T} + \frac{h^2}{12}B \succ \tilde{T} \succ I.$$

Thus, using (3.35), we have

$$(3.37) \quad D^{1/2} \succ \tilde{T}^{1/2} \succ I.$$

We have used that the eigenvalues of \tilde{T} are $\tilde{\lambda}_k = 4 \sin^2(k\pi/2N)/h^2 > 1$, $1 \leq k \leq N-1$. Therefore,

$$(3.38) \quad |\mathbf{e}(t)|_h \leq |(-\delta_x^2)^{1/2}\mathbf{e}(t)|_h \leq |(-\tilde{\delta}_x^2)^{1/2}\mathbf{e}(t)|_h \leq C(t)h^4.$$

Consider next the bound on the discrete derivative $\tilde{\delta}_x$. By the definition of $\tilde{\delta}_x$ in (2.7) and by the definitions of the one-sided discrete first-order derivatives $\delta_x^+ e_j = (e_{j+1} - e_j)/2$ and $\delta_x^- e_j = (e_j - e_{j-1})/2$, we have

$$(3.39) \quad \sigma_x \tilde{\delta}_x \mathbf{e}(t) = \delta_x \mathbf{e}(t) = \frac{1}{2} \left(\delta_x^+ \mathbf{e}(t) + \delta_x^- \mathbf{e}(t) \right).$$

Then,

$$(3.40) \quad \begin{aligned} |\sigma_x \tilde{\delta}_x \mathbf{e}(t)|_h^2 &= \frac{1}{4} |\delta_x^+ \mathbf{e}(t) + \delta_x^- \mathbf{e}(t)|_h^2 \leq \frac{1}{2} (|\delta_x^+ \mathbf{e}(t)|_h^2 + |\delta_x^- \mathbf{e}(t)|_h^2) \\ &= (-\delta_x^2 \mathbf{e}(t), \mathbf{e}(t))_h = |(-\delta_x^2)^{1/2} \mathbf{e}(t)|_h^2 \\ &\leq C(t)^2 h^8. \end{aligned}$$

In the last inequality we have invoked (3.38). In addition, σ_x^{-1} is a uniformly (with h) bounded operator, hence

$$(3.41) \quad |\tilde{\delta}_x \mathbf{e}(t)|_h \leq |\sigma_x^{-1}|_h |\sigma_x \tilde{\delta}_x \mathbf{e}(t)|_h \leq C(t)h^4.$$

The inequalities (3.3)-(3.4) are deduced by taking the maximum of all constants of the form $C(t)$ over $t \in [0, t_f]$. \blacksquare

4. CONVERGENCE ANALYSIS FOR THE EQUATION $\partial_t \Delta u = \Delta^2 u + f$

Consider the time dependent Stokes problem

$$(4.1) \quad \begin{cases} \partial_t \Delta u = \Delta^2 u + f, & (x, y) \in [0, 1] \times [0, 1], & t \geq 0, \\ u(0, y, t) = u(1, y, t) = 0, & u_x(0, y, t) = u_x(1, y, t) = 0, \\ u(x, 0, t) = u(x, 1, t) = 0, & u_y(x, 0, t) = u_y(x, 1, t) = 0, \\ u(x, y, 0) = g(x, y), & 0 \leq x, y \leq 1. \end{cases}$$

As in Section 2, we assume that the functions $f(x, y, t)$, $g(x, y)$ and the solutions $u(x, y, t)$ are regular functions of their variables $(x, y, t) \in [0, 1] \times [0, 1] \times [0, t_f]$, where $0 < t_f$ is a fixed time.

Define the grid function $t \mapsto \mathbf{v}_{j,k}(t)$, $j, k = 0, 1, \dots, N$, which serves as an approximation of u , to be the solution of

$$(4.2) \quad \begin{cases} \partial_t \tilde{\Delta}_h \mathbf{v}_{j,k} = \tilde{\Delta}_h^2 \mathbf{v}_{j,k} + \mathbf{f}_{j,k}, & j, k = 1, \dots, N-1, \\ \mathbf{v}_{0,k}(t) = \mathbf{v}_{N,k}(t) = 0, & (\tilde{\delta}_x \mathbf{v})_{0,k}(t) = (\tilde{\delta}_x \mathbf{v})_{N,k}(t) = 0, & k = 0, \dots, N, \\ \mathbf{v}_{j,0}(t) = \mathbf{v}_{j,N}(t) = 0, & (\tilde{\delta}_y \mathbf{v})_{j,0}(t) = (\tilde{\delta}_y \mathbf{v})_{j,N}(t) = 0, & j = 0, \dots, N, \\ \mathbf{v}_{j,k}(0) = g_{j,k}, & j, k = 0, \dots, N. \end{cases}$$

Here, for $1 \leq j, k \leq N-1$,

$$(4.3) \quad \begin{cases} (\delta_x^4 \mathbf{v})_{j,k} = \frac{12}{h^2} (\delta_x \tilde{\delta}_x \mathbf{v} - \delta_x^2 \mathbf{v})_{j,k}, & (\delta_y^4 \mathbf{v})_{j,k} = \frac{12}{h^2} (\delta_y \tilde{\delta}_y \mathbf{v} - \delta_y^2 \mathbf{v})_{j,k}, \\ (\sigma_x \tilde{\delta}_x \mathbf{v})_{j,k} = (\delta_x \mathbf{v})_{j,k}, & (\sigma_y \tilde{\delta}_y \mathbf{v})_{j,k} = (\delta_y \mathbf{v})_{j,k}, \\ (\sigma_x \mathbf{w})_{j,k} = \frac{1}{6} (\mathbf{w}_{j-1,k} + 4\mathbf{w}_{j,k} + \mathbf{w}_{j+1,k}), & (\sigma_y \mathbf{w})_{j,k} = \frac{1}{6} (\mathbf{w}_{j,k-1} + 4\mathbf{w}_{j,k} + \mathbf{w}_{j,k+1}). \end{cases}$$

In addition,

$$(4.4) \quad \tilde{\Delta}_h^2 \mathbf{v} = \delta_x^4 \mathbf{v} + \delta_y^4 \mathbf{v} + 2 \left(\delta_x^2 \delta_y^2 - \frac{h^2}{12} (\delta_x^4 \delta_y^2 + \delta_y^4 \delta_x^2) \right) \mathbf{v},$$

and

$$(4.5) \quad \tilde{\Delta}_h \mathbf{v} = (\tilde{\delta}_x^2 + \tilde{\delta}_y^2) \mathbf{v}, \quad \tilde{\delta}_x^2 \mathbf{v} = 2\delta_x^2 \mathbf{v} - \delta_x \tilde{\delta}_x \mathbf{v}, \quad \tilde{\delta}_y^2 \mathbf{v} = 2\delta_y^2 \mathbf{v} - \delta_y \tilde{\delta}_y \mathbf{v}.$$

The exact solution u satisfies

$$(4.6) \quad \partial_t \tilde{\Delta}_h u_{j,k}^* = \tilde{\Delta}_h^2 u_{j,k}^* + \mathbf{f}_{j,k} + \boldsymbol{\tau}_{j,k}.$$

where $\boldsymbol{\tau}$ is the truncation error. By Taylor expansions, the components of the truncation error $\boldsymbol{\tau}$ may be written as (see [2] Proposition 10.8)

$$(4.7) \quad \begin{cases} \sigma_x \boldsymbol{\tau}_{j,k} = O(h^4), & \sigma_y \boldsymbol{\tau}_{j,k} = O(h^4) & j, k = 2, \dots, N-2, \\ \sigma_x \boldsymbol{\tau}_{1,k} = O(h), & \sigma_x \boldsymbol{\tau}_{N-1,k} = O(h), & k = 1, \dots, N, \\ \sigma_y \boldsymbol{\tau}_{j,1} = O(h), & \sigma_y \boldsymbol{\tau}_{j,N-1} = O(h), & j = 1, \dots, N. \end{cases}$$

with explicit representations similar to the ones in Section 3 (dimension 1). Define the error

$$(4.8) \quad \boldsymbol{\epsilon} = \mathbf{v} - u^*.$$

By subtracting (4.6) from (4.2), we have

$$(4.9) \quad \partial_t (-\tilde{\Delta}_h) \boldsymbol{\epsilon} + \tilde{\Delta}_h^2 \boldsymbol{\epsilon} = \boldsymbol{\tau}.$$

We relate the grid function $\mathbf{v}_{j,k}$, $j, k = 1, \dots, N-1$ with the column vector $V = \mathbf{vec}(\mathbf{v})$ defined by

$$(4.10) \quad V = [\mathbf{v}_{1,1}, \dots, \mathbf{v}_{N-1,1}, \mathbf{v}_{1,2}, \dots, \mathbf{v}_{N-1,2}, \dots, \mathbf{v}_{1,N-1}, \dots, \mathbf{v}_{N-1,N-1}]^\top \in \mathbb{R}^{(N-1)^2}.$$

The bottom ordering of the vector $V \in \mathbb{R}^{(N-1)^2}$ is obtained by letting the index j vary first while keeping k fixed, then vary the index k ². Then, we relate the two-dimensional finite difference operators with matrix operators of size $(N-1) \times (N-1)$ (for $N \geq 2$) acting on a vector V . Most of those operators are obtained as Kronecker products of $(N-1) \times (N-1)$ matrices. Recall that the Kronecker product of the matrices $G \in \mathbb{M}_{m,n}$ and $H \in \mathbb{M}_{p,q}$ is the matrix $G \otimes H \in \mathbb{M}_{mp,nq}$ defined by

$$(4.11) \quad G \otimes H = \begin{bmatrix} g_{1,1}H & g_{1,2}H & \dots & g_{1,n}H \\ \dots & \dots & \dots & \dots \\ g_{m,1}H & g_{m,2}H & \dots & g_{m,n}H \end{bmatrix}.$$

²This is the standard definition of the operator \mathbf{vec} , see [22, Ch. 4.2, p. 244]

Let the matrix B represent the biharmonic discrete operator δ_x^4 in one dimension and the matrix \tilde{T} represent $-\delta_x^2$ (or $-\delta_y^2$) in one dimension. Then, $I \otimes B$ and $B \otimes I$ represent the biharmonic operators δ_x^4 and δ_y^4 , respectively. Similarly, $I \otimes \tilde{T}$ and $\tilde{T} \otimes I$ represents the operator $-\delta_x^2$ and $-\delta_y^2$, respectively, in the 2D setting. In addition,

$$(4.12) \quad R(t) = [\mathbf{r}_{1,1}, \dots, \mathbf{r}_{N-1,1}, \mathbf{r}_{1,2}, \dots, \mathbf{r}_{N-1,2}, \dots, \mathbf{r}_{1,N-1}, \dots, \mathbf{r}_{N-1,N-1}]^\top \in \mathbb{R}^{(N-1)^2}$$

is related to the truncation error.

As in Section 3, the function $t \mapsto R(t) = \mathbf{vec}(\mathbf{r}(t))$ is regular, since it is obtained by applying linear operators to the function $t \mapsto u(\cdot, t)$ and its derivatives up to order 8, (see Section 6). Invoking Theorem 2.2 we have

Corollary 4.1. *Let $R(t) = R^{(1)}(t) + R^{(2)}(t) \in \mathbb{R}^{(N-1)^2}$, where $R^{(1)} = [R_1^{(1)}; \dots; R_{N-1}^{(1)}]$ and $R^{(2)} = [R_1^{(2)}; \dots; R_{N-1}^{(2)}]$, respectively. Here*

$$(4.13) \quad \begin{aligned} R_1^{(1)} &= [\mathbf{r}_{1,1}, 0, \dots, 0, \mathbf{r}_{N-1,1}]^\top, & R_1^{(2)} &= [0, \mathbf{r}_{2,1}, \dots, \mathbf{r}_{N-2,1}, 0]^\top, \\ R_j^{(1)} &= [\mathbf{r}_{1,j}, \dots, \mathbf{r}_{N-1,j}]^\top, \quad j = 2, \dots, N-2, & R_j^{(2)} &= [0, \dots, 0]^\top, \quad j = 2, \dots, N-2, \\ R_{N-1}^{(1)} &= [\mathbf{r}_{1,N-1}, 0, \dots, 0, \mathbf{r}_{N-1,N-1}]^\top, & R_{N-1}^{(2)} &= [0, \mathbf{r}_{2,N-1}, \dots, \mathbf{r}_{N-2,N-1}, 0]^\top. \end{aligned}$$

Then,

$$(4.14) \quad \max_{1 \leq m \leq (N-1)^2} |((I \otimes B^{-1})R^{(1)}(t))_m| \leq Ch^4, \quad 0 < t < t_f,$$

where $I \otimes B^{-1}$ represents the operator δ_x^{-4} , and

$$(4.15) \quad \max_{1 \leq m \leq (N-1)^2} |((B^{-1} \otimes I)R^{(2)}(t))_m| \leq Ch^4, \quad 0 < t < t_f,$$

where $(B^{-1} \otimes I)$ represents the operator δ_y^{-4} .

Proof. Using the definition of a Kronecker product, we have

$$(4.16) \quad I \otimes B = \begin{bmatrix} B & \mathbf{0} & \dots & \dots & \mathbf{0} \\ \mathbf{0} & B & \mathbf{0} & \dots & \mathbf{0} \\ \dots & & & & \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & B \end{bmatrix}, \quad (I \otimes B)^{-1} = \begin{bmatrix} B^{-1} & \mathbf{0} & \dots & \dots & \mathbf{0} \\ \mathbf{0} & B^{-1} & \mathbf{0} & \dots & \mathbf{0} \\ \dots & & & & \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & B^{-1} \end{bmatrix}.$$

Therefore, $(I \otimes B^{-1})R^{(1)}(t) = \left[B^{-1}R_1^{(1)}(t), B^{-1}R_2^{(1)}(t), \dots, B^{-1}R_{N-2}^{(1)}(t), B^{-1}R_{N-1}^{(1)}(t) \right]^\top$. By the optimal convergence Theorem 2.1

$$(4.17) \quad \max_{1 \leq m \leq (N-1)^2} |((I \otimes B^{-1})R^{(1)}(t))_m| \leq Ch^4, \quad 0 < t < t_f.$$

Hence (4.14) holds. By a similar proof (4.15) holds. ■

We return now to the error analysis.

Theorem 4.2. *Suppose that u , the solution to the system (4.1), has derivatives up to order 8, then the error \mathbf{e} is bounded by*

$$(4.18) \quad \max_{0 \leq t \leq t_f} |(-\tilde{\Delta}_h)^{1/2} \mathbf{e}(t)|_h \leq C(t_f)h^4, \quad 0 < t < t_f.$$

Here $-\tilde{\Delta}_h \mathbf{g}(t) = -(\tilde{\delta}_x^2 + \tilde{\delta}_y^2) \mathbf{g}(t)$ is a fourth-order approximation to the minus Laplacian operator, and $C(t_f)$ denotes a constant depending only on $\partial_t^{k_0} u$, $\partial_x^{k_1} u$, $\partial_y^{k_2} u$, $t \in [0, t_f]$, $(x, y) \in [0, 1] \times [0, 1]$, with $0 \leq k_0 \leq 2$, $4 \leq k_1, k_2 \leq 8$.

Proof. Let $R(t)$ represents the vector of the truncation error. Define $E(t) \triangleq \mathbf{vec}(\mathbf{e}(t))$ as the vector containing the components of the error at time t

$$(4.19) \quad \mathbf{e} = \left[\mathbf{e}_{1,1}, \dots, \mathbf{e}_{N-1,1}; \mathbf{e}_{1,2}, \dots, \mathbf{e}_{N-1,2}; \dots; \mathbf{e}_{1,N-1}, \dots, \mathbf{e}_{N-1,N-1} \right]^\top \in \mathbb{R}^{(N-1)^2}.$$

The operator $\tilde{\Delta}_h^2$ may be represented by the matrix A of size $(N-1)^2 \times (N-1)^2$ (refer to Section 2.2 for the matrix notation.)

$$(4.20) \quad A = I \otimes B + B \otimes I + 2 \left[(\tilde{T} \otimes I)(I \otimes \tilde{T}) + \frac{h^2}{12}(B \otimes I)(I \otimes \tilde{T}) + \frac{h^2}{12}(I \otimes B)(\tilde{T} \otimes I) \right] \succ 0.$$

We turn now to the matrix representation of the error. The operator $-\tilde{\Delta}_h$ may be represented by the matrix M of size $(N-1)^2 \times (N-1)^2$ (see [7]), where

$$(4.21) \quad M = I \otimes \tilde{T} + \tilde{T} \otimes I + \frac{h^2}{12}(I \otimes B) + \frac{h^2}{12}(B \otimes I).$$

Therefore, in vector notation, Equation (4.9) may be written as

$$(4.22) \quad \partial_t(ME(t)) + A E(t) = R(t).$$

Define

$$(4.23) \quad F(t) = M^{1/2}E,$$

Proceeding parallel to (3.6-3.12), we obtain that

$$(4.24) \quad F(t) = \int_0^t e^{-\tilde{A}\rho} M^{-1/2} R(t-\rho) d\rho.$$

Multiplying $M^{-1/2}R(\rho)$ from the left by $\tilde{A} (\tilde{A})^{-1}$, we have

$$(4.25) \quad F(t) = \int_0^t [e^{-\tilde{A}\rho} \tilde{A}] [\tilde{A}^{-1} M^{-1/2} R(t-\rho)] d\rho.$$

where the matrix $\tilde{A} = M^{-1/2}AM^{-1/2}$ is a symmetric positive-definite matrix. Therefore, $e^{-\tilde{A}\rho} \tilde{A}$ is also a symmetric matrix positive-definite matrix. It may be diagonalized by a unitary matrix U . Thus,

$$(4.26) \quad e^{-\tilde{A}\rho} \tilde{A} = U \tilde{\Lambda}(\rho) U^\top,$$

where

$$(4.27) \quad \tilde{\Lambda}(\rho) = \text{diag}\{e^{-\rho\lambda_1}\lambda_1, \dots, e^{-\rho\lambda_{(N-1)^2}}\lambda_{(N-1)^2}\}.$$

Here, λ_k , $k = 1, \dots, (N-1)^2$ are the eigenvalues of \tilde{A} , which are positive. Therefore,

$$(4.28) \quad F(t) = \int_0^t U \tilde{\Lambda}(\rho) U^\top [\tilde{A}^{-1} M^{-1/2} R(t-\rho)] d\rho.$$

Decomposing $R(t)$ into $R(t) = R^{(1)}(t) + R^{(2)}(t)$, where $R^{(1)}(t)$ and $R^{(2)}(t)$ are defined in (4.13)), we have

$$(4.29) \quad F(t) = \int_0^t [U \tilde{\Lambda}(\rho) U^\top] [\tilde{A}^{-1} M^{-1/2}] [R^{(1)}(t-\rho) + R^{(2)}(t-\rho)] d\rho.$$

Invoking $\tilde{A}^{-1} = M^{1/2}A^{-1}M^{1/2}$ and decomposing $F(t)$ in the sum $F(t) = F^{(1)}(t) + F^{(2)}(t)$ we have

$$(4.30) \quad \begin{cases} F^{(1)}(t) = \int_0^t [U \tilde{\Lambda}(\rho) U^\top] [M^{1/2}A^{-1} R^{(1)}(t-\rho)] d\rho, \\ F^{(2)}(t) = \int_0^t [U \tilde{\Lambda}(\rho) U^\top] [M^{1/2}A^{-1} R^{(2)}(t-\rho)] d\rho. \end{cases}$$

We show that $\|F^{(1)}(t)\|_2 \leq Ch^3$. The equality $\|F^{(2)}(t)\|_2 \leq Ch^3$ follows similarly. Notice that

$$(4.31) \quad \tilde{\Lambda}(\rho) = -\frac{d}{d\rho}\Lambda(\rho),$$

where

$$(4.32) \quad \Lambda(\rho) = \text{diag}\{e^{-\rho\lambda_1}, \dots, e^{-\rho\lambda_{(N-1)^2}}\}.$$

Inserting (4.31)-(4.32) into (4.30), we have

$$(4.33) \quad F^{(1)}(t) = -\int_0^t \frac{d}{d\rho}(U\Lambda(\rho) U^\top) [M^{1/2} A^{-1} R^{(1)}(t-\rho)] d\rho.$$

Integration by parts yields

$$(4.34) \quad F^{(1)}(t) = [-(U\Lambda(\rho) U^\top) M^{1/2} A^{-1} R^{(1)}(t-\rho)]_0^t + \int_0^t (U\Lambda(\rho) U^\top) M^{1/2} A^{-1} \frac{d}{d\rho} R^{(1)}(t-\rho) d\rho.$$

We decompose $F^{(1)}(t)$ to $F^{(1,1)}(t) + F^{(1,2)}(t)$, where

$$(4.35) \quad \begin{cases} F^{(1,1)}(t) = \left[-(U\Lambda(\rho) U^\top) M^{1/2} A^{-1} R^{(1)}(t-\rho) \right]_{\rho=0}^t \\ F^{(1,2)}(t) = \int_0^t (U\Lambda(\rho) U^\top) M^{1/2} A^{-1} \frac{d}{d\rho} R^{(1)}(t-\rho) d\rho. \end{cases}$$

We have

$$(4.36) \quad \begin{aligned} F^{(1,1)}(t) &= (U\Lambda(0) U^\top) M^{1/2} A^{-1} R^{(1)}(t) - (U\Lambda(t) U^\top) M^{1/2} A^{-1} R^{(1)}(0) \\ &= (U\Lambda(0) U^\top) M^{1/2} A^{-1} R^{(1)}(t). \end{aligned}$$

We will show in Lemma 5.3 that $\|M^{1/2} A^{-1} R^{(1)}(t-\rho)\|_2 \leq C(t)h^3$. In addition $\|\Lambda(\rho)\|_2 \leq 1$ for $\rho \geq 0$ and $\|U\|_2 = \|U^\top\|_2 = 1$. Therefore,

$$(4.37) \quad \|F^{(1,1)}(t)\|_2 \leq C(t)h^3.$$

We turn now to $F^{(1,2)}(t)$. Notice that $R^{(1)}(\rho)$ (and therefore $\frac{d}{d\rho} R^{(1)}(t-\rho)$) contains fifth and eighth-order derivatives of u with respect to x, y at some intermediate points (\bar{x}, \bar{y}, ρ) , where $(\bar{x}, \bar{y}) \in [0, 1] \times [0, 1]$. Therefore, the time-derivative $\frac{d}{d\rho} R^{(1)}(t-\rho)$ contains first order time-derivative of fifth and eighth-order spatial derivatives of u , where they are all assumed to be bounded.

Since the components of $\frac{d}{d\rho} R^{(1)}(t-\rho)$ are of the same order (as powers of h) as the components of $R^{(1)}(t-\rho)$, we may apply Lemma 5.3 to $\tilde{R} = \frac{d}{d\rho} R^{(1)}(t-\rho)$ and obtain

$$\|M^{1/2} A^{-1} \frac{d}{d\rho} R^{(1)}(t-\rho)\|_2 \leq C(t)h^3.$$

In addition, since $\|U\|_2 = \|U^\top\|_2 = 1$ and $\|\Lambda(\rho)\|_2 \leq 1$, for $0 \leq \rho \leq t$, therefore

$$(4.38) \quad \|F^{(1,2)}\|_2 \leq t \max_{0 \leq \rho \leq t} \|M^{1/2} A^{-1} \frac{d}{d\rho} R^{(1)}(\rho)\|_2 \leq C(t)h^3.$$

From (4.37)-(4.38), we conclude that

$$(4.39) \quad \|F^{(1)}(t)\|_2 \leq C(t)h^3.$$

Similarly,

$$(4.40) \quad \|F^{(2)}(t)\|_2 \leq C(t)h^3.$$

Therefore, $\|F(t)\|_2 \leq C(t)h^3$. Noting that $F(t) = M^{-1/2}E(t)$, we conclude that

$$(4.41) \quad |(-\tilde{\Delta}_h)^{1/2} \mathbf{e}(t)|_h \leq C(t)h^4. \quad \blacksquare$$

5. FOURTH ORDER ESTIMATES FOR THE MATRICES $D^{1/2}B^{-1}R$ AND $M^{1/2}A^{-1}\bar{R}$

In this section we provide fourth-order estimates for $\tilde{T}^{1/2}B^{-1}R$, $\tilde{T}^{1/2}B^{-1}\bar{R}$ and $M^{1/2}A^{-1}\bar{R}$, where R and \bar{R} is the vector which contains the truncation errors in 1D and 2D, respectively. In subsections 5.1, 5.2 and 5.3 we provide a bounds for $\tilde{T}^{1/2}B^{-1}R$, $D^{1/2}B^{-1}R$ and $M^{1/2}A^{-1}\bar{R}$, respectively.

Consider the grid function $\mathbf{r}(t)$ in (2.14) with its components $R_i(t)$. The vector $R(t)$ is the source term in (3.6). In fact, it appears in the expression $D^{1/2}B^{-1}R(t)$, which is part of the integrand of Equation (3.16) In terms of grid functions, it consists of estimating

$$(5.1) \quad (-\tilde{\delta}_x^2)^{1/2}(\delta_x^4)^{-1}\mathbf{r}$$

as a function of h . However, Theorem 2.2 provides only the estimate

$$(5.2) \quad |(\delta_x^4)^{-1}\mathbf{r}| \leq Ch^4,$$

which is not sufficient for our purpose. The goal in this section is to prove a sharper estimate

$$(5.3) \quad |(-\tilde{\delta}_x^2)^{1/2}(\delta_x^4)^{-1}\mathbf{r}| \leq Ch^4.$$

5.1. Estimate for $\tilde{T}^{1/2}B^{-1}R$. The proof of Theorem 2.3 is based on accurate representations of the vector $B^{-1}R$. Here, instead of estimating $D^{1/2}B^{-1}R$, we consider, as a first step, estimating $\tilde{T}^{1/2}B^{-1}R$.

Lemma 5.1. *Assume that the vector $R \in \mathbb{R}^{N-1}$ satisfies*

$$(5.4) \quad PR = [Ch, Ch^4, \dots, Ch^4, Ch]^\top,$$

then all components of $(\tilde{T}^{1/2}B^{-1})R$ are $O(h^4)$, i.e.,

$$(5.5) \quad |(\tilde{T}^{1/2}B^{-1}R)_i| \leq Ch^4, \quad 1 \leq i \leq N-1.$$

Proof. The proof proceeds along the same lines of Theorem 2.2. We replace the matrix B^{-1} by the matrix $\tilde{T}^{1/2}B^{-1}$. By Equation (2.32)

$$(5.6) \quad P^{-1}B^{-1}P^{-1} = HG^{-1}.$$

For $\tilde{T}^{1/2}B^{-1}$ we have

$$(5.7) \quad P^{-1}\tilde{T}^{1/2}B^{-1}P^{-1} = \tilde{H}G^{-1},$$

where $\tilde{H} = \tilde{T}^{1/2}H \in \text{Span}(T)$, (since $\tilde{T} = T/h^2$). Thus, (see (2.31) for the definition of H)

$$(5.8) \quad \tilde{H} = \frac{h^3}{6}P^{-1}T^{-3/2}.$$

As in [16, 2], we bound the elements of \tilde{H} in terms of h . Similar to Equation (10.140) in [2], we have for \tilde{H}

$$(5.9) \quad \tilde{H}_{i,k} = \sum_{j=1}^{N-1} h^3 \frac{2}{N} \frac{\sin(\frac{ij\pi}{N}) \sin(\frac{jk\pi}{N})}{\sin^3(\frac{j\pi}{2N})(6 - 4\sin^2(\frac{j\pi}{2N}))}.$$

In particular

$$(5.10) \quad \tilde{H}_{i,1} = \sum_{j=1}^{N-1} 2h^4 \frac{\sin(\frac{ij\pi}{N}) \sin(\frac{j\pi}{N})}{\sin^3(\frac{j\pi}{2N})(6 - 4\sin^2(\frac{j\pi}{2N}))}.$$

Since for $0 \leq x \leq \pi/2$, we have $\frac{2}{\pi}|x| \leq |\sin x| \leq |x|$ and $2 \leq 6 - 4\sin^2(\frac{j\pi}{2N}) \leq 6$. This gives

$$(5.11) \quad 0 \leq \tilde{H}_{1,1} = \tilde{H}_{N-1,N-1} \leq C \sum_{j=1}^{N-1} \left| \frac{h^4(jh)^2}{(jh)^3} \right| \leq Ch^3 |\ln(h)|,$$

and

$$(5.12) \quad |\tilde{H}_{1,N-1}| = |\tilde{H}_{N-1,1}| \leq Ch^3 |\ln(h)|.$$

In addition,

$$(5.13) \quad |\tilde{H}_{i,1}| \leq C \sum_{j=1}^{N-1} \frac{h^4 j h}{(jh)^3} \leq Ch^2,$$

and

$$(5.14) \quad |\tilde{H}_{i,k}| \leq C \sum_{j=1}^{N-1} \frac{h^4}{(jh)^3} \leq Ch, \quad i, k = 2, \dots, N-2.$$

Therefore,

$$(5.15) \quad \tilde{H} = \begin{bmatrix} O(h^3 |\ln(h)|) & O(h^2) & \dots & O(h^2) & O(h^3 |\ln(h)|) \\ O(h^2) & O(h) & \dots & O(h) & O(h^2) \\ \vdots & \vdots & \dots & \dots & \vdots \\ O(h^2) & O(h) & \dots & O(h) & O(h^2) \\ O(h^3 |\ln(h)|) & O(h^2) & \dots & O(h^2) & O(h^3 |\ln(h)|) \end{bmatrix}.$$

Invoking the estimates in (10.162), (10.163) in [2] and in [16], we have

$$(5.16) \quad G^{-1}R(t) = [O(h^2), O(h^4), \dots, O(h^4), O(h^2)]^\top.$$

Therefore,

$$(5.17) \quad \tilde{H}G^{-1}R(t) = \begin{bmatrix} O(h^3|\ln(h)|) & O(h^2) & \dots & O(h^2) & O(h^3|\ln(h)|) \\ O(h^2) & O(h) & \dots & O(h) & O(h^2) \\ \vdots & \vdots & \dots & \dots & \vdots \\ O(h^2) & O(h) & \dots & O(h) & O(h^2) \\ O(h^3|\ln(h)|) & O(h^2) & \dots & O(h^2) & O(h^3|\ln(h)|) \end{bmatrix} \begin{bmatrix} O(h^2) \\ O(h^4) \\ \vdots \\ O(h^4) \\ O(h^2) \end{bmatrix} = \begin{bmatrix} O(h^5 \ln(h)) \\ O(h^4) \\ \vdots \\ O(h^4) \\ O(h^5 \ln(h)) \end{bmatrix}.$$

Thus, the components of the matrix $\tilde{T}^{1/2}B^{-1}R(t)$ are

$$(5.18) \quad (\tilde{T}^{1/2}B^{-1}R(t))_j = (\tilde{H}G^{-1}\tilde{R}(t))_j = O(h^4), \quad 1 \leq j \leq N-1.$$

5.2. Estimate for $D^{1/2}B^{-1}R$.

Lemma 5.2. *Let D be*

$$(5.19) \quad D = \tilde{T} + \frac{h^2}{12}B.$$

Assume that the vector $R \in \mathbb{R}^{N-1}$ satisfies

$$(5.20) \quad PR = [Ch, Ch^4, \dots, Ch^4, Ch]^\top,$$

then

$$(5.21) \quad \|D^{1/2}B^{-1}R\|_2 \leq Ch^{3.5}.$$

Proof. Denoting $U = D^{1/2}B^{-1}R$, $U^\top U$ is expressed as

$$(5.22) \quad U^\top U = R^\top B^{-1} D B^{-1}R.$$

Using (5.19)-(5.22) gives

$$(5.23) \quad U^\top U = R^\top B^{-1} \tilde{T} B^{-1}R + (h^2/12) R^\top B^{-1}R.$$

Define

$$(5.24) \quad J = R^\top B^{-1} \tilde{T} B^{-1}R,$$

$$K = ((h^2/12) R)^\top B^{-1}R.$$

We look first at the second term K in $U^\top U$. By Theorem 2.2 $B^{-1}R = [O(h^4), O(h^4), \dots, O(h^4), O(h^4)]$, we find that

$$(5.25) \quad \begin{aligned} (h^2/12)R &= [O(h^3), O(h^6), \dots, O(h^6), O(h^3)], \\ B^{-1}R(t) &= [O(h^4), O(h^4), \dots, O(h^4), O(h^4)]. \end{aligned}$$

Therefore,

$$(5.26) \quad K = O(h^7).$$

As for the first term J in $U^\top U$, denoting by U_1

$$(5.27) \quad U_1 = (\tilde{T})^{1/2} B^{-1}R,$$

we have

$$(5.28) \quad J = U_1^\top U_1 = R^\top B^{-1} \tilde{T} B^{-1}R.$$

It has been shown in (5.18) that

$$(5.29) \quad U_1 = [O(h^4), O(h^4), \dots, O(h^4), O(h^4)].$$

Therefore,

$$(5.30) \quad J = U_1^\top U_1 = O(h^7).$$

Combining (5.26) and (5.29) yields

$$(5.31) \quad U^\top U = J + K = O(h^7).$$

■

5.3. Estimate for $M^{1/2}A^{-1}\tilde{R}$. We prove now that $\|M^{1/2}A^{-1}R^{(1)}(t-\rho)\|_2 \leq Ch^3$. To simplify the notation, let us define $\tilde{R} = R^{(1)}(t-\rho)$ [a similar Lemma is true for $\tilde{R} = \frac{d}{d\rho}R^{(1)}(t-\rho)$].

Lemma 5.3. *Let M be (see (4.21))*

$$(5.32) \quad M = I \otimes D + D \otimes I = I \otimes \tilde{T} + \frac{h^2}{12}(I \otimes B) + \tilde{T} \otimes I + \frac{h^2}{12}(B \otimes I)$$

and let $\tilde{R} = R^{(1)}$, where $R^{(1)}$ is defined in (4.13). Then,

$$(5.33) \quad \|M^{1/2}A^{-1}\tilde{R}\|_2 \leq C(t)h^3.$$

Proof. Consider now

$$(5.34) \quad (\tilde{R})^\top A^{-1}M A^{-1} \tilde{R}.$$

Invoking (4.21), we have

$$(5.35) \quad A^{-1}M A^{-1} = A^{-1}[I \otimes D + D \otimes I] A^{-1}.$$

Splitting the sum above to two terms $I \otimes D$ and $D \otimes I$, we have

$$(5.36) \quad A^{-1}M A^{-1} = A^{-1}(I \otimes D)A^{-1} + A^{-1}(D \otimes I)A^{-1} = P_1 + P_2,$$

where

$$(5.37) \quad \begin{aligned} P_1 &= A^{-1}(I \otimes D)A^{-1}, \\ P_2 &= A^{-1}(D \otimes I)A^{-1}. \end{aligned}$$

First, we consider P_1 . Since the elements of $(I \otimes D)^{1/2}A^{-1}\tilde{R}$ are smaller or equal in terms of orders of h compared to $(I \otimes D)^{1/2}(I \otimes B^{-1})\tilde{R}$, we will bound the latter. Thus, we want to bound now $\tilde{R}^\top (I \otimes B^{-1})(I \otimes D)(I \otimes B^{-1})\tilde{R}$.

Recall that

$$(5.38) \quad D = \tilde{T} + \frac{h^2}{12}B.$$

We will show that

$$(5.39) \quad \|(\tilde{R})^\top P_1 \tilde{R}\|_2 \leq C\|(\tilde{R})^\top [(I \otimes B)^{-1}(I \otimes D)(I \otimes B)^{-1}] \tilde{R}\|_2 \leq Ch^6.$$

This is equivalent to showing that

$$(5.40) \quad \|(I \otimes D)^{1/2}(I \otimes B)^{-1}\tilde{R}\|_2 \leq C(t)h^3.$$

It is shown in Lemma 5.2

$$\|D^{1/2}B^{-1}R\|_2 \leq Ch^{3.5},$$

therefore

$$(5.41) \quad \|(I \otimes D)^{1/2}(I \otimes B)^{-1}\tilde{R}\|_2 \leq Ch^3.$$

Thus, we conclude that

$$(5.42) \quad (\tilde{R})^\top P_1 \tilde{R} \leq C(t)h^6.$$

Second, we consider P_2 , defined in (5.37). In Lemma 5.4 we prove that

$$(5.43) \quad P_2 = (\tilde{R})^\top P_2 \tilde{R} \leq C(t)h^6.$$

Therefore, from $A^{-1}M A^{-1} = P_1 + P_2$, we will conclude that

$$(5.44) \quad \|M^{1/2}A^{-1}\tilde{R}\|_2 \leq Ch^3. \quad \blacksquare$$

Lemma 5.4. *We show that $(R^{(1)})^\top P_2 R^{(1)} \leq Ch^6$, where P_2 is defined*

$$(5.45) \quad P_2 = A^{-1}(D \otimes I)A^{-1},$$

and $R^{(1)}$ is defined in (4.13).

Proof. Here we look at $R^{(1)} = R^{(1,1)} + R^{(1,2)}$, where

$R^{(1,1)} = [R_1^{(1,1)}; \dots; R_{N-1}^{(1,1)}]$ and $R^{(1,2)} = [R_1^{(1,2)}; \dots; R_{N-1}^{(1,2)}]$, respectively, where

$$(5.46) \quad \begin{aligned} R_1^{(1,1)} &= [\mathbf{r}_{1,1}, 0, \dots, 0, \mathbf{r}_{N-1,1}]^\top, & R_1^{(1,2)} &= [0, 0, \dots, 0, 0]^\top, \\ R_j^{(1,1)} &= [\mathbf{r}_{1,j}, 0, \dots, 0, \mathbf{r}_{N-1,j}]^\top, \quad j = 2, \dots, N-2, & R_j^{(1,2)} &= [0, \mathbf{r}_{2,j}, \dots, \mathbf{r}_{N-2,j}, 0]^\top, \quad j = 2, \dots, N-2, \\ R_{N-1}^{(1,1)} &= [\mathbf{r}_{1,N-1}, 0, \dots, 0, \mathbf{r}_{N-1,N-1}]^\top. & R_{N-1}^{(1,2)} &= [0, 0, \dots, 0, 0]^\top. \end{aligned}$$

First, we consider $(\tilde{R})^\top P_2 \tilde{R}$, where we take $\tilde{R} = R^{(1,1)}$. Note that $P_2 = A^{-1}(D \otimes I)A^{-1}$.

Consider now

$$(5.47) \quad (\tilde{R})^\top P_2 \tilde{R} = \tilde{R}^\top A^{-1}(D \otimes I)A^{-1} \tilde{R}.$$

We want to determine the order with respect to h of the components of $(D^{1/2} \otimes I)A^{-1}\tilde{R}$. Notice that the elements in \tilde{R} related to $i = 2, \dots, N-2$ are zero for $j = 1, \dots, N-1$. Consider the operator $(D^{1/2} \otimes I)(I \otimes B)^{-1}$. This operator approximates $\partial_y(\delta_x^4)^{-1}$. Thus, for $i = 2, \dots, N-2$, $j = 1, \dots, N-1$ we have that $(D^{1/2} \otimes I)A^{-1}\tilde{R}$ is zero. Consider now $i = 1, N-1$. The elements of \tilde{R} which are related to $i = 1, N-1$ are $O(h)$ for any $j = 1, \dots, N-1$. Note that the elements of $A^{-1}\tilde{R}$ are smaller or equal compared to the magnitudes of the elements of $(I \otimes B)^{-1}\tilde{R}$, therefore we consider now $(I \otimes B)^{-1}\tilde{R}$. We have that $(I \otimes B^{-1}) = (I \otimes H)(I \otimes G^{-1})$. We also have that $G_{1,1}^{-1}, G_{N-1,1}^{-1}, G_{1,N-1}^{-1}, G_{N-1,N-1}^{-1}$ are $O(h)$, therefore the elements of $(I \otimes G^{-1})\tilde{R}$ related to $i = 1, N-1$ are $O(h^2)$ for $j = 1, \dots, N-1$. In addition, we have that $H_{1,1}, H_{N-1,1}, H_{1,N-1}, H_{N-1,N-1}$ are $O(h^3)$. Therefore, the elements of $(I \otimes B^{-1})\tilde{R}$ which are related to $i = 1, N-1$ are $O(h^5)$ for all $j = 1, \dots, N-1$. Operating on $(I \otimes B)^{-1}\tilde{R}$ with $(D^{1/2} \otimes I)$ results in elements which are $O(h^4)$ for $i = 1, N-1$, $j = 1, \dots, N-1$.

The vector $(D^{1/2} \otimes I)(I \otimes B)^{-1}\tilde{R}$ contains only $2(N-1)$ non zero terms, which are $O(h^4)$. Therefore,

$$(5.48) \quad \|(D^{1/2} \otimes I)(I \otimes B)^{-1}\tilde{R}\|_2 = O(h^3).$$

Thus,

$$(5.49) \quad \tilde{R}^\top (I \otimes B)^{-1}(D \otimes I)(I \otimes B)^{-1}\tilde{R} = O(h^6).$$

Therefore, we have

$$(5.50) \quad \tilde{R}^\top P_2 \tilde{R} \leq Ch^6.$$

As for

$$(5.51) \quad (\tilde{R}^{(1,2)})^\top A^{-1}(D \otimes I)A^{-1}\tilde{R}^{(1,2)},$$

since $(D \otimes I)^{1/2}A^{-1}$ is a bounded operator, and since the non-zero components of $\tilde{R}^{(1,2)}$ are $O(h^4)$, we conclude that

$$(5.52) \quad (\tilde{R}^{(1,2)})^\top P_2 \tilde{R}^{(1,2)} = O(h^6).$$

Combining (5.50)-(5.52), we conclude that

$$(5.53) \quad (\tilde{R}^{(1)})^\top P_2 \tilde{R}^{(1)} \leq C h^6. \quad \blacksquare$$

6. TRUNCATION ERROR REPRESENTATION

Let $u(x, t)$ be a regular function with associated grid function $u^*(t) = [u(x_1, t), u(x_2, t), \dots, u(x_{N-1}, t)]$. The truncation error $\mathbf{r}(u)(t)$, denoted for simplicity also by $\mathbf{r}(t)$, is expressed as

$$(6.1) \quad \mathbf{r}(t) = (\mathcal{L}u(\cdot, t))^* - \mathcal{L}_h u^*(t),$$

where \mathcal{L} the continuous operator $\mathcal{L} = -\partial_{txx} + \partial_x^4$ and \mathcal{L}_h the semi-discrete operator $\mathcal{L}_h = \partial_t(-\tilde{\delta}_x^2) + \delta_x^4$, with Dirichlet boundary conditions in both cases. Here, our goal is to derive an expression of the form

$$(6.2) \quad \sigma_x \mathbf{r}(t) = [O(h), O(h^4), \dots, O(h^4), O(h)]^\top,$$

with an effective *representation* of the $O(h^\alpha)$ components in terms of a spatial multi-point Taylor expansion of the functions $\partial_x^{(k)} u(x, t)$ and $\partial_t \partial_x^{(k)} u(x, t)$, (with fixed time t). Indeed, such a representation is more explicit than the formal identity (6.2). This form is used to accurately handling the estimates (3.29) and (3.31) in Section 3.

The analysis hereafter is elementary and relies on the standard Taylor-Lagrange expansions and convex combinations of derivatives. The following abbreviation is used for convenience. For a given regular function $x \in (0, 1) \mapsto u(x)$, we denote

$$(6.3) \quad \begin{cases} u^k \triangleq \partial_x^{(k)} u, \\ u^{k,*} \triangleq (\partial_x^{(k)} u)^*, \quad u_i^{k,*} \triangleq (\partial_x^{(k)} u)(x_i). \end{cases}$$

The parenthesis is omitted when specifying an operation evaluated at a grid point x_i . For example, we denote $\delta_x u_i^* \triangleq (\delta_x u^*)_i$. The truncation error $\mathbf{r}(t)$ is decomposed as $\mathbf{r}(t) = -\mathbf{r}_a(t) + \mathbf{r}_b(t)$ with

$$(6.4) \quad \begin{cases} \mathbf{r}_a(t) = \delta_x^4(u(\cdot, t)^*) - u^{4,*}(t), \\ \mathbf{r}_b(t) = \tilde{\delta}_x^2(\partial_t u)^*(\cdot, t) - (\partial_t u)^{2,*}(t). \end{cases}$$

Concerning the Hermitian derivative, the following upper bound of the truncation error for the Hermitian derivative has been derived in [2] Lemma 10.1.

$$(6.5) \quad |\tilde{\delta}_x u_i^* - u_i^{1,*}| \leq Ch^4 \|u^5\|_{\infty, (0,1)},$$

with the constant $C = 1/60$. In [8], a representation of this truncation is given using the basis Z^k in (2.36). In the following Lemma, we give a simpler representation of the truncation error.

Lemma 6.1. *Let $u(x)$ be a regular function. The components of the truncation error $\tau(u) = \tilde{\delta}_x u^* - u^{1,*}$ can be expressed in the form*

$$(6.6) \quad \tau(u)_i = \tilde{\delta}_x u_i^* - u_i^{1,*} = h^4 \left(\gamma_i^1 u^5(\xi_i^1) - \gamma_i^2 u^5(\xi_i^2) \right), \quad 1 \leq i \leq N-1,$$

where

- The scalar $\xi_i^1, \xi_i^2 \in (0, 1)$.
- The scalars γ_i^1, γ_i^2 are fixed positive constants, independent of $u(x)$, and satisfying for all $i \in \llbracket 1, N-1 \rrbracket$ and N :

$$(6.7) \quad \gamma_i^1 \leq 0.0547, \quad \gamma_i^2 \leq 0.0506.$$

Remark 6.2. *The scalars ξ_i^1 and ξ_i^2 are not necessarily close to x_i . This corresponds to the fact that $\mathbf{u} \mapsto \tilde{\delta}_x \mathbf{u}$ is a "non local" operator. They depend on $u(x)$.*

Proof. The 4th order Hermitian derivative $\tilde{\delta}_x u_j$, $1 \leq j \leq N-1$, is defined in (2.7). It is a non local finite difference operator. This non-locality gives that the truncation error results in a multi-point Taylor representation. The multi-point can be compacted using convex combination of derivatives into a two point representation. The proof is as follows. In vector form, we have $\sigma_x^{-1} \hookrightarrow [\alpha_{i,j}] \in \mathbb{M}_{N-1}$ with

$$(6.8) \quad [\alpha_{i,j}] = \begin{bmatrix} \frac{2}{3} & \frac{1}{6} & 0 & \dots & 0 \\ \frac{1}{6} & \frac{2}{3} & \frac{1}{6} & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & \dots & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \\ 0 & \dots & 0 & \frac{1}{6} & \frac{2}{3} \end{bmatrix}^{-1}.$$

The matrix $[\alpha_{i,j}]$ is the inverse of a tridiagonal Toeplitz positive definite matrix. According to [29], the coefficients $\alpha_{i,j}$ are

$$(6.9) \quad \alpha_{i,j} = (-1)^{i+j} \beta_{i,j}, \quad \beta_{i,j} = |\alpha_{i,j}|, \quad 1 \leq i, j \leq N-1,$$

where

$$(6.10) \quad \begin{cases} \beta_{i,j} = 6 \frac{(\lambda^i - \mu^i)(\lambda^{N-j} - \mu^{N-j})}{(\lambda - \mu)(\lambda^N - \mu^N)}, & 1 \leq i \leq j \leq N-1, \\ \beta_{i,j} = \beta_{j,i}, & 1 \leq i \leq j \leq N-1. \end{cases}$$

The constants λ and μ are $\lambda = 2 + \sqrt{3} \simeq 3.73$, $\mu = 2 - \sqrt{3} \simeq 0.28$. The matrix $[\alpha_{i,j}]$ is full. The $\beta_{i,j}$ behave as

$$(6.11) \quad 0 < \beta_{i,1} < \beta_{i,2} < \dots, \beta_{i,i-1} < \beta_{i,i} > \beta_{i,i+1} \dots > \beta_{i,N-1} > 0, \quad 1 \leq i \leq N-1,$$

with

$$(6.12) \quad \begin{cases} 0 < \beta_{i,i} < \frac{6}{\lambda - \mu} = \sqrt{3}, \\ \frac{\beta_{i,j}}{\beta_{i,j+1}} > \lambda > 1, & i \leq j \leq N-1, \quad 1 \leq i \leq N-1, \\ \frac{\beta_{i,j}}{\beta_{i,j-1}} > \lambda > 1, & 1 \leq j \leq i-1, \quad 1 \leq i \leq N-1. \end{cases}$$

The truncation error τ of $\tilde{\delta}_x \mathbf{u}$ satisfies

$$(6.13) \quad (\sigma_x \tau)_i = \delta_x u_i^* - \sigma_x u_i^{1,*}, \quad 1 \leq i \leq N-1.$$

Taylor-Lagrange expansions of u_{i+1}^* and u_{i-1}^* give that

$$(6.14) \quad \delta_x u_i^* = u_i^{1,*} + h^2 \frac{1}{3!} u_i^{3,*} + h^4 \frac{1}{5!} u^5(\eta_i^1), \quad \eta_i^1 \in (x_{i-1}, x_{i+1}).$$

By definition of σ_x is $\sigma_x u_i^{1,*} = u_i^{1,*} + \frac{h^2}{3!} \delta_x^2 u_i^{1,*}$. Furthermore, by Taylor expansion

$$(6.15) \quad \delta_x^2 u_i^{1,*} = u_i^{3,*} + h^2 \frac{2}{4!} u^5(\eta_i^2), \quad \eta_i^2 \in (x_{i-1}, x_{i+1}).$$

It results from (6.13)-(6.15) that

$$(6.16) \quad \begin{aligned} (\sigma_x \tau)_i &= \delta_x u_i - \sigma_x u_i^{1,*} \\ &= u_i^{1,*} + h^2 \frac{1}{3!} u_i^{3,*} + h^4 \frac{1}{5!} u^5(\eta_i^1) - \left(I + h^2 \frac{1}{3!} \delta_x^2 \right) u_i^{1,*} \\ &= h^4 \left(\frac{1}{5!} u^5(\eta_i^1) - \frac{1}{2 \cdot 3!^2} u^5(\eta_i^2) \right). \end{aligned}$$

Multiplying the vector $[(\sigma_x \tau)_1, \dots, (\sigma_x \tau)_{N-1}]^T$ by the matrix $[\alpha_{i,j}]$ in (6.8) yields

$$(6.17) \quad \tau_i = h^4 \sum_{j=1}^{N-1} (-1)^{i+j} \beta_{i,j} \left(\frac{1}{5!} u^5(\eta_j^1) - \frac{1}{2 \cdot 3!^2} u^5(\eta_j^2) \right).$$

The sum in (6.17) is the difference of two sums, both with positive coefficients. Defining the two sets of indices $J_{\text{even}}(i) = \{1 \leq j = i + 2p \leq N-1\}$ and $J_{\text{odd}}(i) = \{1 \leq j = i + 2p + 1 \leq N-1\}$, $p \in \mathbb{Z}$, one has

$$(6.18) \quad \begin{aligned} \tau_i &= h^4 \left[\frac{1}{5!} \sum_{j \in J_{\text{even}}(i)} \beta_{i,j} u^5(\eta_j^1) + \frac{1}{2 \cdot 3!^2} \sum_{j \in J_{\text{odd}}(i)} \beta_{i,j} u^5(\eta_j^2) \right] \\ &\quad - h^4 \left[\frac{1}{5!} \sum_{j \in J_{\text{odd}}(i)} \beta_{i,j} u^5(\eta_j^1) + \frac{1}{2 \cdot 3!^2} \sum_{j \in J_{\text{even}}(i)} \beta_{i,j} u^5(\eta_j^2) \right]. \end{aligned}$$

The first and second term can be collected in one single term by convex combination of the derivatives. Define

$$(6.19) \quad \begin{cases} \gamma_i^1 = \frac{1}{5!} \sum_{j \in J_{\text{even}}(i)} \beta_{i,j} + \frac{1}{2 \cdot 3!^2} \sum_{j \in J_{\text{odd}}(i)} \beta_{i,j}, \\ \gamma_i^2 = \frac{1}{5!} \sum_{j \in J_{\text{odd}}(i)} \beta_{i,j} + \frac{1}{2 \cdot 3!^2} \sum_{j \in J_{\text{even}}(i)} \beta_{i,j}. \end{cases}$$

Due to (6.12), the scalars γ_i^1, γ_i^2 are well defined and uniformly bounded (independently of i and N). This gives that there exists $\xi_i^1, \xi_i^2 \in (0, 1)$ such that (6.6) holds. \blacksquare

We next derive representation formulas for $\mathbf{r}_a(t)$ and $\mathbf{r}_b(t)$ in (6.4) having a form similar to (6.6: they use a set of points in $(0, 1)$ depending of the function $u(x)$.

Lemma 6.3 (Truncation term τ_a). *For any regular function $u(x)$, consider the truncation error $\mathbf{r}_a(u)$ defined by*

$$(6.20) \quad \mathbf{r}_a(u) = \delta_x^4 u^* - u^{4,*}.$$

The components of $\sigma_x(\mathbf{r}_a)$ are of the following orders of h

$$(6.21) \quad \sigma_x \mathbf{r}_a = [O(h), O(h^4), \dots, O(h^4), O(h)]^\top.$$

The structure of (6.21) is as follows. There exists positive constants β_i^k uniformly bounded independently of i and N , such that

- At the internal points x_i , $2 \leq i \leq N - 2$,

$$(6.22) \quad (\sigma_x \mathbf{r}_a)_i = h^4 \left(\beta_i^1 u^8(\xi_i^1) - \beta_i^2 u^8(\xi_i^2) \right),$$

with constants $\beta_i^1 = \frac{2^6}{7!} \simeq 0.0127$ and $\beta_i^2 = \frac{7!}{7!} \simeq 0.0141$ and for some $\xi_i^1 \in (x_{i-2}, x_{i+2})$, $\xi_i^2 \in (x_{i-1}, x_{i+1})$.

- At the near boundary point x_1 , we have

$$(6.23) \quad (\sigma_x \mathbf{r}_a)_i = h \left(\beta_1^3 u^5(\xi_i^3) - \beta_1^4 u^5(\xi_i^4) \right) + h^2 \left(\beta_1^5 u^6(\xi_i^5) - \beta_1^6 u^6(\xi_i^6) \right),$$

for some $\xi_i^3, \xi_i^4, \xi_i^5, \xi_i^6 \in (0, x_3)$. A similar expression holds at the point x_{N-1} .

Remark 6.4. Note that (6.22) and (6.23) contain a finite number of terms, which are uniformly bounded in terms of derivatives of $u(x)$. The points ξ_i^j depend on the function $u(x)$.

Proof. For the truncation analysis at interior points $2 \leq i \leq N - 2$, we refer to [2, chap. 10.4, pp. 160 sqq]. It is based on Taylor-Lagrange expansions at point x_i of

$$(6.24) \quad (\sigma_x \delta_x^4) \mathbf{v}_i = (\delta_x^2)^2 \mathbf{v}_i = \frac{1}{h^4} (\mathbf{v}_{i+2} - 4\mathbf{v}_{i+1} + 6\mathbf{v}_i - 4\mathbf{v}_{i-1} + \mathbf{v}_{i-2}), \quad 2 \leq i \leq N - 2,$$

and of

$$(6.25) \quad \sigma_x u_i^{4,*} = \frac{1}{6} u_{i-1}^{4,*} + \frac{2}{3} u_i^{4,*} + \frac{1}{6} u_{i+1}^{4,*}, \quad 2 \leq i \leq N - 2.$$

Subtracting these two expansions results in (6.22).

Next consider x_1 and x_2 . We have

$$(6.26) \quad \begin{aligned} (\sigma_x \mathbf{r}_a)_1 &= \sigma_x \delta_x^4 u_1^* - \sigma_x u_1^{4,*} \\ &= \frac{2}{3} (\delta_x^4 u_1^* - u_1^{4,*}) + \frac{1}{6} (\delta_x^4 u_2^* - u_2^{4,*}). \end{aligned}$$

We have

$$(6.27) \quad (\delta_x^4 u^*)_1 = \frac{12}{h^2} \left((\delta_x \tilde{\delta}_x) u_1^* - \delta_x^2 u_1^* \right), \quad (\delta_x^4 u^*)_2 = \frac{12}{h^2} \left(\delta_x \tilde{\delta}_x u_2^* - \delta_x^2 u_2^* \right).$$

For all regular functions $v(x)$, the Taylor-Lagrange formula for all $1 \leq i \leq N - 1$ gives

$$(6.28) \quad \delta_x^2 v_i^* = v_i^{2*} + h^2 \frac{2}{4!} v_i^{4,*} + h^4 \frac{2}{6!} v_i^6(\eta_i^9), \quad \eta_i^9 \in (x_{i-1}, x_{i+1}).$$

In particular, for $i = 1$,

$$(6.29) \quad \delta_x^2 u_1^* = u_1^{2,*} + h^2 \frac{2}{4!} u_1^{4,*} + h^4 \frac{2}{6!} u_1^6(\eta_1^9), \quad \eta_1^9 \in (x_0, x_2).$$

Consider in (6.27) the term $(\delta_x \tilde{\delta}_x) u_1^*$. By Lemma 6.1 and by (6.14) (recalling that $u_0^{1,*} = \tilde{\delta}_x u_0^* = 0$), we have

$$(6.30) \quad \begin{aligned} (\delta_x \tilde{\delta}_x) u_1^* &= \frac{1}{2h} (\tilde{\delta}_x u_2^* - \tilde{\delta}_x u_0^*) = \frac{1}{2h} \tilde{\delta}_x u_2^* = \frac{1}{2h} \left(u_2^{1,*} - u_0^{1,*} + h^4 (\gamma_2^1 u^5(\xi_2^1) - \gamma_2^2 u^5(\xi_2^2)) \right) \\ &= u_1^{2,*} + h^2 \frac{1}{3!} u_1^{4,*} + h^3 \frac{1}{2} \left(\gamma_2^1 u^5(\xi_2^1) - \gamma_2^2 u^5(\xi_2^2) \right) + h^4 \frac{1}{5!} u^6(\eta_1^{10}), \quad \xi_2^1, \xi_2^2 \in (0, 1), \quad \eta_1^{10} \in (x_0, x_2). \end{aligned}$$

By (6.27), it results that

$$(6.31) \quad \delta_x^4 u_1^* - u_1^{4,*} = 6h \left(\gamma_2^1 u^5(\xi_2^1) - \gamma_2^2 u^5(\xi_2^2) \right) + h^2 \left(\frac{12}{5!} u^6(\eta_1^{10}) - \frac{24}{6!} u^6(\eta_1^9) \right).$$

For $i = 2$, we have for some $\eta_2^{11} \in (x_1, x_3)$, $\delta_x^2 u_2^* = u_2^{2,*} + h^2 \frac{2}{4!} u_2^{4,*} + h^4 \frac{2}{6!} u^6(\eta_2^{11})$ and for some $\eta_2^{12} \in (x_1, x_3)$,

$$(6.32) \quad (\delta_x \tilde{\delta}_x) u_2^* = u_2^{2,*} + h^2 \frac{1}{3!} u_2^{4,*} + h^3 \frac{1}{2} \left((\gamma_3^1 u^5(\xi_3^1) + \gamma_1^1 u^5(\xi_1^1)) - (\gamma_3^2 u^5(\xi_3^2) + \gamma_1^2 u^5(\xi_1^2)) \right) + h^4 \frac{1}{5!} u^6(\eta_2^{12}).$$

Thus, by (6.27)

$$(6.33) \quad \delta_x^4 u_2^* - u_2^{4,*} = 6h \left((\gamma_1^1 u^5(\xi_1^1) + \gamma_3^1 u^5(\xi_3^1)) - (\gamma_1^2 u^5(\xi_1^2) + \gamma_3^2 u^5(\xi_3^2)) \right) + h^2 \left(\frac{12}{5!} u^6(\eta_2^{12}) - \frac{24}{6!} u^6(\eta_2^{11}) \right).$$

It results from (6.31) and (6.33) that $(\sigma_x \mathbf{r}_a)_1$, given in (6.26), satisfies (6.23) with $\beta_1^3 = \gamma_1^1 + 4\gamma_2^1 + \gamma_3^1$, $\beta_1^4 = \gamma_1^2 + 4\gamma_2^2 + \gamma_3^2$, $\beta_1^5 = 1/12$, $\beta_1^6 = 1/36$. The proof for $i = N - 1$ is similar. \blacksquare

Consider now the term $\mathbf{r}_b(t)$ in (6.4). Denoting $v(x) = \partial_t u(x, t)$, we have

$$(6.34) \quad \mathbf{r}_b(t) = \tilde{\delta}_x^2 v^*(t) - v^{2,*}(t).$$

Lemma 6.5 (Truncation term τ_b). *For fixed $t > 0$, the truncation error $\mathbf{r}_b(t)$ satisfies*

$$(6.35) \quad \sigma_x \mathbf{r}_b = [O(h^2), O(h^4), \dots, O(h^4), O(h^2)]^\top.$$

There exists positive constants $\beta_i^{\prime, k}$ uniformly bounded with respect to N such that the components of $\sigma_x \mathbf{r}_b$ are represented as follows.

- At internal points x_i , i.e., at x_i , for $2 \leq i \leq N - 2$,

$$(6.36) \quad \sigma_x \mathbf{r}_{b,i} = h^4 \left(\beta_i^{\prime,1} v^6(\xi_i^{\prime,1}) \right) + h^6 \left(\beta_i^{\prime,2} v^8(\xi_i^{\prime,2}) - \beta_i^{\prime,3} v^8(\xi_i^{\prime,3}) \right),$$

for some $\xi_i^{\prime,1}, \xi_i^{\prime,2} \in (x_{i-2}, x_{i+2})$, $\xi_i^{\prime,3} \in (x_{i-1}, x_{i+1})$.

- At near boundary points x_i with at $i = 1$ or $i = N - 1$, we have

$$(6.37) \quad \sigma_x \mathbf{r}_{b,i} = h^2 \left(\beta_i^{\prime,4} v^4(\xi_i^{\prime,4}) \right) + h^3 \left(\beta_i^{\prime,5} v^5(\xi_i^{\prime,5}) - \beta_i^{\prime,6} v^5(\xi_i^{\prime,6}) \right) + h^4 \left(\beta_i^{\prime,7} v^6(\xi_i^{\prime,7}) - \beta_i^{\prime,8} v^6(\xi_i^{\prime,8}) \right)$$

for some $\xi_i^{\prime,4}, \xi_i^{\prime,5}, \xi_i^{\prime,6}, \xi_i^{\prime,7} \in (0, 1)$.

We skip the proof of Lemma 6.5. Note that in (6.36) and (6.37), the points $\xi_i^{\prime, j}$ depend on $v(x)$.

Consider now a regular function $u(x, t)$ and $v(x, t) = \partial_t u(x, t)$. Consider a fixed time $\bar{t} > 0$, and the truncation $\mathbf{r}(\bar{t}) \triangleq \mathbf{r}(u(x, \bar{t}))$. It results from Lemma 6.3 and Lemma 6.5 that $\sigma_x \mathbf{r}(\bar{t})$ can be expressed as follows.

Corollary 6.6. *Suppose that u is a solution to the problem (2.1) such that $x \mapsto u^k(\cdot, t), (\partial_t u)^k(\cdot, t)$ are continuous for $0 \leq k \leq 8$, then the truncation error $\mathbf{r}(t)$ satisfies an identity of the form*

$$(6.38) \quad \sigma_x \mathbf{r}(\bar{t}) = [O(h), O(h^4), \dots, O(h^4), O(h)]^\top$$

with Taylor representation expressed as

$$(6.39) \quad \begin{aligned} (\sigma_x \mathbf{r}(\bar{t}))_i &= h^4 \left[\sum_{j=0}^2 C_{i,j} u^8(y_{i,j}, \bar{t}) + D_{i,j} v^6(z_{i,j}, \bar{t}) \right] + h^6 \left[\sum_{j=0}^2 C'_{i,j} v^8(z'_{i,j}, \bar{t}) \right] = O(h^4), \quad i = 2, \dots, N-2, \\ (\sigma_x \mathbf{r}(\bar{t}))_1 &= h \left[\sum_{j=0}^2 C_{1,j} u^5(y_{1,j}, \bar{t}) \right] + h^2 \left[\sum_{j=0}^2 C'_{1,j} u^6(y'_{1,j}, \bar{t}) + D'_{1,j} v^4(z'_{1,j}, \bar{t}) \right] \\ &\quad + h^3 \left[\sum_{j=0}^2 D''_{1,j} v^5(z''_{1,j}, \bar{t}) \right] + h^4 \left[\sum_{j=0}^2 D'''_{1,j} v^6(z'''_{1,j}, \bar{t}) \right] = O(h), \end{aligned}$$

with a similar expression for $(\sigma_x \mathbf{r}(\bar{t}))_{N-1}$. All the constants $C_{i,j}, D_{i,j}, \dots$ are fixed. The points $y_{i,j}, y'_{i,j}, \dots$ belong to the interval $(0, 1)$. They a priori depend on the function $x \mapsto u(x, \bar{t})$.

Coming back to (3.26), we consider the particular case where $u(x, t)$ is the solution of (2.1). The vector $R(t)$ in (3.26) satisfies $PR(\bar{t}) \hookrightarrow \sigma_x \mathbf{r}(\bar{t})$ with components given in (6.39) with $\bar{t} = t$. From Corollary 6.6, we conclude that the vector $PR(t)$ satisfies the hypothesis (5.20) of Lemma 5.1.

Similarly, we have

$$(6.40) \quad P \frac{d}{dt} R(\bar{t}) = [O(h), O(h^4), \dots, O(h^4), O(h)].$$

with an analog explicit representation.

Remark 6.7. The truncation analysis for the 2D case proceeds in a similar fashion.

7. NUMERICAL EXAMPLES

In the previous sections we have considered compact approximations in space of the equations (2.1) and (4.1). In this section we display numerical results for the equations above, using the fourth-order discrete operators for the approximation in the spatial direction. The temporal integration of these equations is carried out via the method of lines. We invoke an IMEX2 scheme of Crank-Nicolson type. For example, for the two-dimensional time-dependent Stokes equation $\partial_t \Delta u = \Delta^2 u + f$ we invoke the following IMEX scheme

$$(7.1) \quad \frac{(\tilde{\Delta}_h u_{i,j})^{n+1/2} - (\tilde{\Delta}_h u_{i,j})^n}{\Delta t/2} = \frac{1}{2} (\tilde{\Delta}_h^2 u_{i,j}^{n+1/2} + \tilde{\Delta}_h^2 u_{i,j}^n) + f_{i,j}^{n+1/4},$$

$$(7.2) \quad \frac{(\tilde{\Delta}_h u_{i,j})^{n+1} - (\tilde{\Delta}_h u_{i,j})^n}{\Delta t} = +\frac{1}{2} (\tilde{\Delta}_h^2 u_{i,j}^{n+1} + \tilde{\Delta}_h^2 u_{i,j}^n) + f_{i,j}^{n+1/2}.$$

This scheme is if second order in time, therefore in order to achieve fourth-order overall accuracy in time and space, we pick $\Delta t = Ch^2$.

7.1. Example 1. Consider the solution $u(x, t) = e^{-t} \sin(\pi x)/\pi^2$ of the problem

$$(7.3) \quad \begin{cases} \partial_{xxt} u = \partial_x^4 u + f(x, t), & 0 < x < 1, \quad t > 0, \\ u(0, t) = 0, \quad u_x(0, t) = e^{-t}/\pi, & t > 0, \\ u(1, t) = 0, \quad u_x(1, t) = -e^{-t}/\pi, & t > 0, \\ u(x, 0) = \sin(\pi x)/\pi^2, & 0 \leq x \leq 1. \end{cases}$$

Here u and u_x are given at the two boundary points and $f(x, t)$ is chosen such that $u(x, t) = e^{-t} \sin(\pi x)/\pi^2$ is the exact solution of the problem above. The numerical results are shown in Table 1. They calibrate the fourth-order accuracy of the scheme.

mesh	9	Rate	17	Rate	33	Rate	65
$ e _h$	4.9254(-6)	4.02	3.0327(-7)	4.01	1.8884(-8)	4.00	1.1791(-9)
$ e_x _h$	3.2671(-6)	4.01	2.0308(-7)	4.00	1.2671(-8)	4.00	7.9157(-10)

TABLE 1. Compact scheme for $\partial_{xxt} u = \partial_x^4 u + f$ with exact solution: $u(x, t) = e^{-t} \sin(\pi x)/\pi^2$ on $x \in [0, 1], t > 0$. We display e and e_x , the l_2 errors for the u and for u_x . Here $\Delta t = h^2$ and $t = 0.5$.

7.2. **Example 2.** Consider the solution $u(x, t) = 16x^2(1-x)^2 \sin(1/((x-0.5)^2 + \varepsilon) \sin(2\pi t)$, with $\varepsilon = 0.05$ of the problem

$$(7.4) \quad \begin{cases} \partial_{xxt}u = \partial_x^4 u + f(x, t), & 0 < x < 1, \quad t > 0, \\ u(0, t) = 0, \quad u_x(0, t) = 0, & t > 0, \\ u(1, t) = 0, \quad u_x(1, t) = 0, & t > 0, \\ u(x, 0) = 0, & 0 \leq x \leq 1. \end{cases}$$

Here u and u_x are given at the two boundary points and $f(x, t)$ is chosen such that $16x^2(1-x)^2 \sin(1/((x-0.5)^2 + \varepsilon) \sin(2\pi t)$ is the exact solution of the problem above. Notice that the solution is very oscillatory near the center of the interval $[0, 1]$. The numerical results are shown in Table 2. They calibrate the fourth-order accuracy of the scheme.

mesh	33	Rate	65	Rate	129	Rate	257
$ e _h$	0.0523(0)	7.03	3.997(-4)	4.28	2.0579(-5)	4.07	1.2255(-6)
$ e_x _h$	1.4598(0)	4.41	0.0683(0)	4.03	0.0040(0)	4.04	2.4242(-4)

TABLE 2. Compact scheme for $\partial_{xxt}u = \partial_x^4 u + f$ with exact solution: $16x^2(1-x)^2 \sin(1/((x-0.5)^2 + \varepsilon) \sin(2\pi t)$ with $\varepsilon = 0.05$ on $x \in [0, 1], t > 0$. We display e and e_x , the l_2 errors for the u and for u_x . Here $\Delta t = h^2$ and $t = 0.25$.

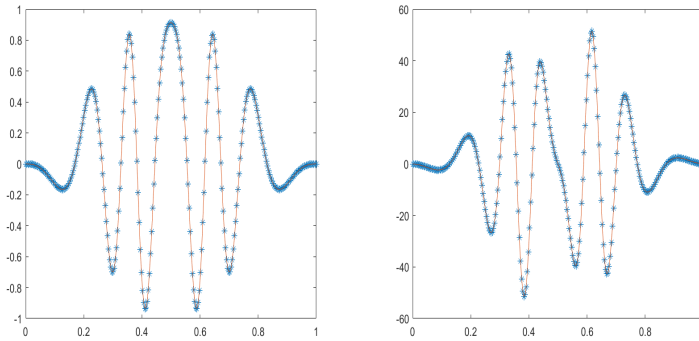


FIGURE 1. Example 2. Left: Approximation for $u(x, t)$. Right: Approximation for $\partial_x u(x, t)$

7.3. **Example 3.** Consider the solution $u(x, y, t) = (1-x^2)^2(1-y^2)^2 e^{-t}$ of the problem

$$(7.5) \quad \begin{cases} \partial_t(\Delta u) = \Delta^2 u + f(x, y, t), & -1 < x, y < 1, \quad t > 0, \\ u(-1, y, t) = u(1, y, t) = 0, \quad u_x(-1, y, t) = u_x(1, y, t), & t > 0, \\ u(x, -1, t) = u(x, 1, t) = 0, \quad u_y(x, -1, t) = u_y(x, 1, t) = 0, & t > 0, \\ u(x, 0) = (1-x^2)^2(1-y^2)^2, & -1 \leq x, y \leq 1. \end{cases}$$

Here u and u_x are given at the two boundary points and $f(x, t)$ is chosen such that $u(x, y, t) = (1-x^2)^2(1-y^2)^2 e^{-t}$ is the exact solution of the problem above. The numerical results are shown in Table 3. They assess the fourth-order accuracy of the scheme.

mesh	9×9	Rate	17×17	Rate	33×33	Rate	65×65
$ e _h$	1.2386(-4)	4.00	7.7408(-6)	4.00	4.8376(-7)	4.00	3.0235(-8)
$ e_x _h$	2.0259(-4)	3.99	1.2750(-5)	4.00	7.9731(-7)	4.00	4.9834(-8)

TABLE 3. Compact scheme for $\partial_t(\Delta u) = \Delta^2 u + f$ with exact solution: $u(x, y, t) = (1-x^2)^2(1-y^2)^2 e^{-t}$ on $(x, y) \in [-1, 1] \times [-1, 1], t > 0$. We display e and e_x , the l_2 errors for the u and for u_x . Here $\Delta t = h^2$ and $t = 0.25$.

7.4. **Example 4.** Consider the solution $u(x, y, t) = -0.5e^{-2t} \sin^2 x \cdot \sin^2 y$ of the problem

$$(7.6) \quad \begin{cases} \partial_t(\Delta u) = \Delta^2 u + f(x, y, t), & 0 < x, y < \pi \quad t > 0, \\ u(0, y, t) = u(\pi, y, t) = 0, & u_x(0, y, t) = u_x(\pi, y, t), \quad t > 0, \\ u(x, 0, t) = u(x, \pi, t) = 0, & u_y(x, 0, t) = u_y(x, \pi, t) = 0, \quad t > 0, \\ u(x, 0) = -0.5 \sin^2 x \cdot \sin^2 y, & 0 \leq x, y \leq \pi. \end{cases}$$

Here u and u_x are given at the two boundary points and $f(x, t)$ is chosen such that $u(x, y, t) = -0.5e^{-2t} \sin^2 x \cdot \sin^2 y$ is the exact solution of the problem above. The numerical results are shown in Table 4. They calibrate the fourth-order accuracy of the scheme.

mesh	9×9	Rate	17×17	Rate	33×33	Rate	65×65
$ e _h$	1.9508(-3)	3.96	1.2527(-4)	3.99	7.9061(-6)	4.00	4.9542(-7)
$ e_x _h$	2.6996(-3)	3.98	1.7134(-4)	3.99	1.0775(-5)	4.00	6.7459(-7)

TABLE 4. Compact scheme for $\partial_t(\Delta u) = \Delta^2 u + f$ with exact solution: $u(x, y, t) = -0.5e^{-2t} \sin^2 x \cdot \sin^2 y$ on $(x, y) \in [0, \pi] \times [0, \pi], t > 0$. We present e and e_x , the l_2 errors for the u and for u_x . Here $\Delta t = h^2$ and $t = 0.6168$.

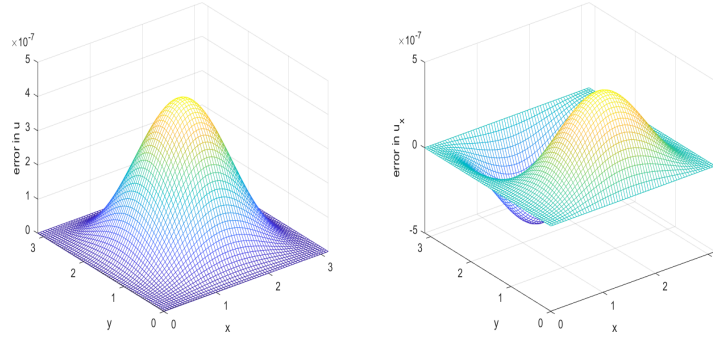


FIGURE 2. Example 4. Left: Error for $u(x, y, t)$. Right: Error for $\partial_x u(x, y, t)$

7.5. **Example 5.** Consider the solution $u(x, y, t) = (1 - x^2)^3(1 - y^2)^3 e^{-t}$ of the NS problem in the square $[-1, 1] \times [-1, 1]$.

$$(7.7) \quad \partial_t \Delta u + C(u) = \nu \Delta^2 u + f(x, y, t), \quad [-1, 1] \times [-1, 1], \quad t > 0,$$

where $C(u) = \nabla^\perp u \cdot \nabla \Delta u$ is the convective term of the NS equations. The boundary and initial data deduced from the exact solution. The viscosity is $\nu = 1$. The final time is $t_f = 1$. The numerical results are shown in Table 5.

mesh	9×9	Rate	17×17	Rate	33×33	Rate	65×65
$ e _h$	1.9373(-3)	4.00	1.2072(-4)	4.00	7.5424(-6)	4.00	4.7138(-7)
$ e_x _h$	1.9886(-3)	4.02	1.2255(-4)	4.00	7.6527(-6)	4.00	4.7827(-7)

TABLE 5. Compact scheme for $\partial_t(\Delta u) = -C(u) + \nu \Delta^2 u + f$ with exact solution: $u(x, y, t) = (1 - x^2)^3(1 - y^2)^3 e^{-t}$ in the square $[-1, 1] \times [-1, 1]$. We present e and e_x , the l_2 errors for the u and for u_x . Here $\Delta t = Ch^2$ and $t_f = 1$.

8. CONCLUSION

In [17], a $O(h^{3.5})$ error estimate was proved for (1.3) using energy methods. Here we establish that the error is in fact $O(h^4)$, for (1.3) and (1.1), despite the fact the the truncation error drops at interior points from $O(h^4)$ to $O(h)$ at near boundary points. Note also that in [3] we have proved second order convergence for a compact second order scheme for the linear equation close to (1.1) was established. In [4] we have proved that the error for the second order compact scheme for the full nonlinear Navier-Stokes equations is $O(h^{1.5})$. Our goal in a future research is to prove that fourth-order convergence can be extended to the Navier-Stokes system.

REFERENCES

- [1] S. Abarbanel, A. Ditzkowski and B. Gustafsson, "On error bounds of finite difference approximations to partial differential equations temporal behavior and rate of convergence", *J. Sci. Comput.*, **15** (2000), pp. 79-116.
- [2] M. Ben-Artzi, J.-P. Croisille and D. Fishelov, "*Navier-Stokes Equations in Planar Domains*", Imperial College Press, 2013.
- [3] M. Ben-Artzi, J.-P. Croisille and D. Fishelov and S. Trachtenberg, "A pure compact scheme for the streamfunction formulation of the Navier-Stokes equations", *J. Comp. Phys.*, **205** (2005), pp. 640-664.
- [4] M. Ben-Artzi, J.-P. Croisille and D. Fishelov, "Convergence of a compact scheme for the pure streamfunction formulation of the unsteady Navier-Stokes system", *SIAM J. Numer. Anal.*, **44** (2006), pp. 1997-2024.
- [5] M. Ben-Artzi, J.-P. Croisille and D. Fishelov, "A High Order Compact Scheme for the Pure-Streamfunction Formulation of the Navier-Stokes Equations", *J. Sci. Comput.*, **42** (2010), pp. 216-250.
- [6] M. Ben-Artzi, J.-P. Croisille and D. Fishelov, "Time evolution of discrete fourth-order elliptic operators", *Num. Meth. Part. Diff. Eqs.*, **35** (2019), pp. 1429-1457.
- [7] M. Ben-Artzi, J.-P. Croisille and D. Fishelov, "A fast direct solver for the biharmonic problem in a rectangular grid", *SIAM J. Sci. Comput.*, **31** (2008), pp. 303-333.
- [8] M. Ben-Artzi, J.-P. Croisille and D. Fishelov, "A Cartesian compact scheme for the Navier-Stokes equations in streamfunction formulation in irregular domains", *J. Sci. Comp.*, **81** (2019), pp. 1386-1408.
- [9] R. Bhatia, "*Matrix analysis*", GTM 169, Springer, 1997.
- [10] W. R. Briley, "A numerical study of laminar separation bubbles using the Navier-Stokes equations", *J. Flu. Mech.*, **47**, (1971). pp. 713-736.
- [11] G.F. Carey and W.F. Spitz, "Extension of High-Order Compact Schemes to Time Dependent Problems, *Numer. Meth. Part. Diff. Eqs.*, **17** (2001), pp. 657-672.
- [12] M. H. Carpenter and D. Gottlieb and S. Abarbanel, "The stability of numerical boundary treatments for compact high-order schemes finite difference schemes", *J. Comput. Phys.*, **108** (1993), pp. 272-295.
- [13] X. Chen and Y. Li and C. Drapaca and J. Cimbala, "A unified framework of continuous and discontinuous Galerkin methods for solving the incompressible Navier-Stokes equation", *J. Comp. Phys.*, **422**, (2020), doi: 10.1016/j.jcp.2020.109799
- [14] P. Drazin and N. Riley, "*The Navier-Stokes Equations - A classification of Flows and Exact Solutions*", London Math. Soc. Lecture Note Series 334, Cambridge Univ. Press, 2006
- [15] W. E and J.-G. Liu, "Essentially compact schemes for unsteady viscous incompressible flows", *J. Comput. Phys.*, **126(1)**, (1996), pp. 122-138.
- [16] D. Fishelov and M. Ben-Artzi and J.-P. Croisille, "Recent advances in the study of a fourth-order compact scheme for the one-dimensional biharmonic equation", *J. Sci. Comput.*, **53(1)**, (2012), pp. 55-79.
- [17] D. Fishelov, "Semi-discrete time-dependent fourth-order problems on an interval: error estimate", "Numerical Mathematics and Advanced Applications - ENUMATH 2013", in Assyr Abdulle, Simone Deparis, Daniel Kressner, Fabio Nobile, and Marco Picasso, *Lect. Notes in Comp. Sci. and Eng.*, Vol. 103, pp. 133-142, Springer, 2015.
- [18] M. M. Gupta and R.P.Manohar and J.W.Stephenson, "Single cell high order scheme for the convection-diffusion equation with variable coefficients", *Int. J. Numer. Meth. Fluids*, **4**, (1984), pp. 641-651.
- [19] B. Gustafsson, "The convergence rate for difference approximations to mixed initial boundary value problems", *Math. of Comput.*, **29** (1975), pp. 396-406.
- [20] B. Gustafsson, "The convergence rate for difference approximations to general mixed initial boundary value problems", *SIAM J. Numer. Anal.*, **18** (1981), pp. 179-190.
- [21] R. A. Horn and C. R. Johnson, *Topics in Matrix analysis*, vol. 1, Cambridge U. Press, (1985).
- [22] R. A. Horn and C. R. Johnson, *Topics in Matrix analysis*, vol. 2, Cambridge U. Press, (1991).
- [23] H. Johnston and R. Krasny, "Fourth-order finite difference simulation of a differentially heated cavity", *Int. J. Numer. Meth. Fluids*, **40**, (2002), pp 1031-1037.
- [24] J. C. Kalita and M. M. Gupta, "A streamfunction-velocity approach for 2D transient incompressible viscous flows", *Int. J. Numer. Meth. Fluids*, **62** (2010), pp. 237-266.
- [25] H. Lamb, "*Hydrodynamics*", 6th ed., Dover, 1932.
- [26] J. Lequeurre and A. Munnier, "Vorticity and Stream Function Formulations for the 2D Navier-Stokes Equations in a Bounded Domain", *J. Math. Fluid Mech.*, **22:15**, (2020). <https://doi.org/10.1007/s00021-019-0479-5>.
- [27] M. Li and T. Tang, "A compact fourth-order finite difference scheme for unsteady viscous incompressible flows", *J. Sci. Comput.*, **103 (1)** (2001), pp. 29-45.
- [28] J.-G. Liu, C. Wang, H. Johnston, "A fourth order scheme for incompressible Boussinesq equations", *Journal of Scientific Computing*, vol. 18 (2), 2003, 253-285.

- [29] G. Meurant, "A review on the inverse of symmetric tridiagonal and block tridiagonal matrices", *SIAM J. Matrix. Anal. App.*, **13(3)**, (1992), pp. 707-728.
- [30] S. Müller and F. Schweiger and E. Süli, "Optimal-Order Finite Difference Approximation of Generalized Solutions to the Biharmonic Equation in a Cube", *SIAM J. Numer. Anal.*, **58(1)**, (2020), pp. 298-329.
- [31] K. Mattsson and J. Nordström, "Summation by parts operators for finite difference approximations of second derivatives", *J. Comput. Phys.*, **199** (2004), pp. 503-540.
- [32] M. J. Naughton, "On numerical boundary conditions for the Navier-Stokes equations", PhD Thesis, Cal. Inst. Tech, 1986.
- [33] C. Pozrikidis, "Creeping flow in two-dimensional channels", *J. Fluid. Mech.*, **180** (1987), pp. 495-514.
- [34] A. Sifounakis and S. Lee and Y. You, "A conservative finite volume method for incompressible Navier–Stokes equations on locally refined nested Cartesian grids", *J. Comput. Phys.*, **326** (2016), pp. 845-861.
- [35] J. W. Stephenson, "Single cell discretizations of order two and four for biharmonic problems", *J. Comput. Phys.*, **55** (1984), pp. 65-80.
- [36] K. Svärd and J. Nordström, "On the convergence rates of energy-stable finite-difference schemes", *J. Comput. Phys.*, **397** (2019), 108819.
- [37] S. Wang and G. Kreiss, "Convergence of Summation-by-Parts Finite Difference Methods for the Wave Equation", *J. Sci. Comput.*, **71** (2017), pp. 219-245.
- [38] C. Wang, J.-G. Liu, "Analysis of finite difference schemes for unsteady Navier-Stokes equations in vorticity formulation", *Numerische Mathematik*, vol. 91, 2002, 543-576.
- [39] C. Wang, J.-G. Liu, "Fourth order convergence of compact difference solver for incompressible flow", *Communications in Applied Analysis*, vol. 7 (2), 2003, 171-191.
- [40] C. Wang, J.-G. Liu, H. Johnston, "Analysis of a fourth order finite difference method for incompressible Boussinesq equations", *Numerische Mathematik*, vol. 97 (3), 2004, 555-594.

DALIA FISHELOV: AFEKA - TEL-AVIV ACADEMIC COLLEGE OF ENGINEERING, 38 MIVZA KADESH ST., TEL-AVIV 69107, ISRAEL
Email address: daliaf@afeka.ac.il

JEAN-PIERRE CROISILLE: DEPARTMENT OF MATHEMATICS, IECL, UMR 7502, UNIV. DE LORRAINE, METZ 57045, FRANCE
Email address: jean-pierre.croisille@univ-lorraine.fr