

# An Explicit Jacobian for Newton's Method Applied to Numerical Approximations of Nonlinear Initial Boundary Value Problems\*

Fredrik Laurén<sup>†</sup>, Oskar Ålund<sup>†</sup>, and Jan Nordström<sup>†‡</sup>

**Abstract.** We derive an explicit form of the Jacobian for discrete approximations of nonlinear initial boundary value problems (IBVPs) on matrix-vector form. The technique is exemplified on the incompressible Navier-Stokes equations in two dimensions. The Jacobian facilitates the use of Newton's method to solve the corresponding nonlinear system of equations. Appropriate boundary conditions are weakly imposed and we show how to compute the Jacobian for those parts of the discretization as well. The convergence rate of the iterations is verified by using the method of manufactured solutions. The methodology in this paper that can be used on any numerical discretization of IBVPs on matrix-vector form.

**Key words.** Nonlinear initial boundary value problems, Jacobian, Newton's method, incompressible Navier-Stokes equations, summation-by-parts, weak boundary conditions.

**AMS subject classifications.** 65M06, 65M12

**1. Introduction.** Nonlinear systems of partial differential equations are common in computational science and engineering, and present multiple challenges. Stability is needed for reliability and high accuracy for fine solution details. For fast turnaround and timely result delivery, generic systems of nonlinear equations from the discretization of the form

$$(1.1) \quad \mathcal{F}(\phi) = 0$$

must be solved efficiently [15]. This is the topic of this paper. Several techniques exist to solve (1.1), for example dual-time stepping [7, 13], optimization algorithms [10] or iterative methods [15]. Among the classical iterative methods, Newton's method is an effective choice due to its quadratic convergence order. The obvious drawback with Newton's method is that the Jacobian must be known. Methods that bypass this requirement and instead approximate the Jacobian lead to lower convergence orders, a typical example is the Secant method [15]. Alternatively, by using Newton-Krylov methodologies [8], only the action of the Jacobian is required and that can be approximated by  $J_{\mathcal{F}}(\phi^k)\delta\mathbf{u} \approx (\mathcal{F}(\phi^k + \delta\mathbf{u}) - \mathcal{F}(\phi^k))/\delta$ , where  $\delta$  is small and  $\mathbf{u}$  depends on the subspaces in the Krylov iterations. The advantage of Newton-Krylov methods is that an explicit Jacobian is never required, but sophisticated preconditioners becomes necessary [2] instead.

The focus in this paper is to facilitate the use of Newton's method where the key component is an exact explicit form of the Jacobian of (1.1). To exemplify our technique, we will use finite-

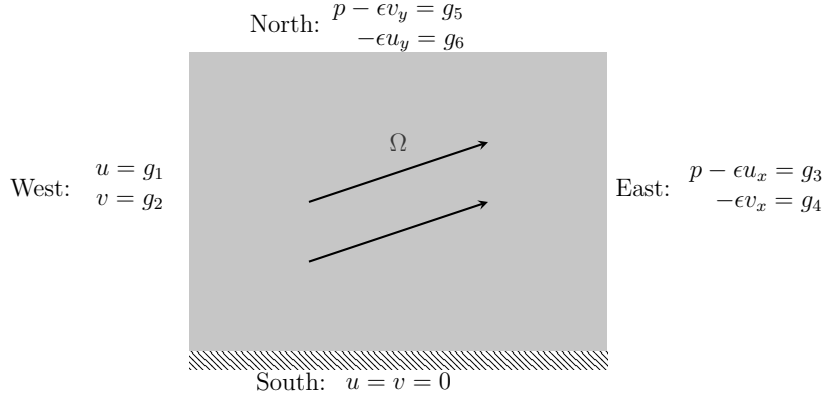
---

\*Submitted to the editors October 21, 2021.

**Funding:** This work was funded by the Swedish Research Council (Stokholm) under grant number 2018-05084.VR and SESSI.

<sup>†</sup>Department of Mathematics, Computational Mathematics, Linköping University, SE-581 83 Linköping, Sweden ([fredrik.lauren@liu.se](mailto:fredrik.lauren@liu.se)), ([oskar.alund@liu.se](mailto:oskar.alund@liu.se)), ([jan.nordstrom@liu.se](mailto:jan.nordstrom@liu.se)).

<sup>‡</sup>Department of Mathematics and Applied Mathematics, University of Johannesburg, P.O. Box 524, Auckland Park 2006, South Africa.



**Figure 1.** Illustration of the computational domain  $\Omega$  and the specific boundary conditions.

difference operators on summation-by-parts (SBP) form [18] to discretize the incompressible Navier-Stokes (INS) equations in space. The boundary conditions will be weakly imposed via the Simultaneous Approximation Term (SAT) technique [3]. In [12], such a discretization based on the SBP-SAT technique of the nonlinear INS equations was proven to be stable, which is the key prerequisite.

Based on the formulation in [12], we show how the Jacobian can be explicitly calculated. It is also shown that the Jacobian has a block structure, where several blocks are precomputed when forming  $\mathcal{F}$ , making the procedure very efficient.

To keep the paper focused on the derivation of the Jacobian, we follow [12] and consider a Cartesian grid. Exact Jacobians for numerical discretizations have recently been developed in [5] for so-called entropy stable numerical discretizations on SBP form in a periodic setting. Our new technique is not restricted to such specific discretizations, and we include the specific Jacobian related to the boundary conditions. The technique demonstrated in this paper can be used in a straightforward way on any numerical method for IBVPs that can be formulated on matrix-vector form. In addition, it can be readily extended to curvilinear grids [1], arbitrary dimensions, and other sets of linear and nonlinear equations. In principle all that is needed for the existence of the Jacobian is of course that  $\mathcal{F}$  is differentiable with respect to  $\phi$ . This covers the various nonlinearities that arise in discretizations of the Navier-Stokes equations (both compressible and incompressible). However, the feasibility of explicitly deriving the Jacobian is highly dependent on the way in which  $\mathcal{F}$  is presented. As we shall see, discretizations on SBP-SAT form are particularly simple to differentiate, making Newton's method an attractive solution method.

The rest of the paper proceeds as follows. We introduce the continuous problem in Section 2 and present the semi-discrete formulation in Section 3. The Jacobian of the discretization is derived in Section 4. Implicit time integration is discussed in Section 5 and numerical experiments are performed in Section 6. Finally, conclusions are drawn in Section 7.

**2. Problem formulation.** As an illustrative example of our technique, we consider the scenario illustrated in Figure 1. An incompressible fluid is moving from left to right. Hence, the left side is an inflow boundary, where Dirichlet conditions are imposed, and the right side

is an outflow boundary, where natural boundary conditions [14] are imposed. The lower part of the domain is a no-slip wall and the upper side is an outflow boundary, where again natural conditions are imposed. The initial-boundary value problem for the INS equations that we consider is

$$(2.1) \quad \begin{aligned} \tilde{I} \vec{w}_t + \mathcal{L}(\vec{w}) &= 0 & (x, y) \in \Omega, & \quad t > 0 \\ \mathcal{H} \vec{w} &= \vec{g} & (x, y) \in \partial\Omega, & \quad t > 0 \\ \tilde{I} \vec{w} &= \tilde{I} \vec{f} & (x, y) \in \Omega, & \quad t = 0. \end{aligned}$$

In (2.1),  $\vec{w} = (u, v, p)^\top$ , where  $u, v$  are the velocities in the  $x, y$  direction, respectively, and  $p$  is the pressure. Furthermore,  $\Omega = [0, 1]^2$  is the domain and  $\partial\Omega$  its boundary. The initial data  $\vec{f}$  and boundary data  $\vec{g}$  are sufficiently smooth and compatible functions and the spatial operator is given by [12]

$$(2.2) \quad \mathcal{L}(\vec{w}) = \frac{1}{2} [A \vec{w}_x + (A \vec{w})_x + B \vec{w}_y + (B \vec{w})_y] - \epsilon \tilde{I} [\vec{w}_{xx} + \vec{w}_{yy}].$$

The matrices in (2.1) and (2.2) are

$$A = \begin{pmatrix} u & 0 & 1 \\ 0 & u & 0 \\ 1 & 0 & 0 \end{pmatrix}, \quad B = \begin{pmatrix} v & 0 & 0 \\ 0 & v & 1 \\ 0 & 1 & 0 \end{pmatrix}, \quad \tilde{I} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

Lastly, the explicit form of the boundary conditions  $\mathcal{H} \vec{w} = \vec{g}$  reads

$$(2.3) \quad \begin{aligned} u &= g_1 & v &= g_2 & \text{at } x &= 0 & (\text{West}) \\ p - \epsilon u_x &= g_3 & -\epsilon v_x &= g_4 & \text{at } x &= 1 & (\text{East}) \\ u &= 0 & v &= 0 & \text{at } y &= 0 & (\text{South}) \\ p - \epsilon v_y &= g_5 & -\epsilon u_y &= g_6 & \text{at } y &= 1 & (\text{North}). \end{aligned}$$

**2.1. Boundedness.** We will for completeness show how to bound the solution. For simplicity, only the south side of the domain is discussed explicitly. Details of the upcoming analysis are found in [12].

For two vector functions  $\vec{\phi}, \vec{\psi}$  defined on  $\Omega$ , we introduce the inner product and norm

$$\langle \vec{\phi}, \vec{\psi} \rangle = \int_{\Omega} \vec{\phi}^\top \vec{\psi} d\Omega, \quad \|\vec{\phi}\|^2 = \langle \vec{\phi}, \vec{\phi} \rangle.$$

By multiplying (2.1) by  $2 \vec{w}^\top$  from the left and integrating over  $\Omega$ , we get

$$(2.4) \quad \frac{d}{dt} \|\vec{w}\|_I^2 + 2\epsilon \|\nabla \vec{w}\|_I^2 = BT,$$

where  $\nabla \vec{w} = (\nabla u, \nabla v, \nabla p)^\top$ ,  $\|\nabla \vec{w}\|_I^2$  is a dissipative volume term and

$$BT = \int_{\text{South}} \vec{w}^\top (B \vec{w} - 2\epsilon \tilde{I} \vec{w}_y) dx$$

contains the boundary terms evaluated at the south boundary. The other boundary terms are assumed dissipative and ignored. Imposing  $u = v = 0$  results in  $BT = 0$ . Integrating (2.4) in time (assuming homogeneous boundary conditions on all sides) leads to

$$(2.5) \quad \|\vec{w}\|_{\tilde{I}}^2(T) + 2\epsilon \int_0^T \|\nabla \vec{w}\|_{\tilde{I}}^2 dt \leq \|f\|_{\tilde{I}}^2,$$

which bounds the semi-norm of the solution ( $\|\vec{w}\|_{\tilde{I}}^2$ ) and its gradients ( $\|\nabla \vec{w}\|_{\tilde{I}}^2$ ) for any time.

**3. The semi-discrete scheme.** A brief introduction of the SBP-SAT technique is provided below and we recommend [6, 18] for extensive reviews.

We discretize the domain  $\Omega = [0, 1]^2$  with  $N + 1$  and  $M + 1$  grid points;  $x_i = i/N$ ,  $i = 0, \dots, N$  and  $y_j = j/M$ ,  $j = 0, \dots, M$  and let  $n = (N + 1)(M + 1)$  denote the total number of grid points. A scalar function  $q = q(x, y)$  defined on  $\Omega$  is thereby represented on the grid by  $\mathbf{q} = (q_{00}, \dots, q_{0M}, \dots, q_{N0}, \dots, q_{NM})^\top$  where  $q_{ij} = q(x_i, y_j)$ . For the vector-valued function  $\vec{w} = (u, v, p)^\top$ , the approximation is arranged as  $\vec{\mathbf{w}} = (\mathbf{u}^\top, \mathbf{v}^\top, \mathbf{p}^\top)^\top$ . Let  $\mathbf{D}_x = (P_x^{-1}Q_x) \otimes I_{M+1}$  and  $\mathbf{D}_y = I_{N+1} \otimes (P_y^{-1}Q_y)$ , where  $\otimes$  denotes the Kronecker product. Then the approximations of the spatial derivatives are given by

$$\mathbf{D}_x \mathbf{u} \approx \mathbf{u}_x, \quad \mathbf{D}_y \mathbf{u} \approx \mathbf{u}_y.$$

The matrices  $P_{x,y}$  are diagonal and positive definite, so that  $\mathbf{P} = P_x \otimes P_y$  forms a quadrature rule that defines the norm  $\|\vec{\mathbf{w}}\|_{I_3 \otimes \mathbf{P}}^2 = \vec{\mathbf{w}}^\top (I_3 \otimes \mathbf{P}) \vec{\mathbf{w}} \approx \iint_\Omega \vec{w}^\top \vec{w} d\Omega$ . We have also introduced  $I_k$ , which is the identity matrix of size  $k$ . Moreover, the matrices  $Q_{x,y}$  satisfy the SBP-property

$$(3.1) \quad Q_x + Q_x^\top = E_N - E_{0x} \quad Q_y + Q_y^\top = E_M - E_{0y},$$

where  $E_{0x,y} = \text{diag}(1, 0, 0, \dots, 0)$  and  $E_{N,M} = \text{diag}(0, 0, 0, \dots, 1)$  are matrices of appropriate sizes.

By using the notation above, the semi-discrete approximation of (2.1) becomes [12]

$$(3.2) \quad \tilde{\mathbf{I}} \vec{\mathbf{w}}_t + \mathcal{L}(\vec{\mathbf{w}}) = \mathcal{S}(\vec{\mathbf{w}}).$$

The discrete spatial operator is given by

$$\begin{aligned} \mathcal{L}(\vec{\mathbf{w}}) = & \frac{1}{2} [A(I_3 \otimes \mathbf{D}_x) \vec{\mathbf{w}} + (I_3 \otimes \mathbf{D}_x)A \vec{\mathbf{w}} + B(I_3 \otimes \mathbf{D}_y) \vec{\mathbf{w}} + (I_3 \otimes \mathbf{D}_y)B \vec{\mathbf{w}}] \\ & - \epsilon \tilde{\mathbf{I}} [(I_3 \otimes \mathbf{D}_x)^2 + (I_3 \otimes \mathbf{D}_y)^2] \vec{\mathbf{w}}, \end{aligned}$$

and the block matrices are

$$A = \begin{pmatrix} \mathbf{U} & \mathbf{0} & \mathbf{I} \\ \mathbf{0} & \mathbf{U} & \mathbf{0} \\ \mathbf{I} & \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad B = \begin{pmatrix} \mathbf{V} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{V} & \mathbf{I} \\ \mathbf{0} & \mathbf{I} & \mathbf{0} \end{pmatrix}, \quad \tilde{\mathbf{I}} = \begin{pmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix},$$

where  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times n}$  are diagonal matrices holding  $\mathbf{u}, \mathbf{v}$ , respectively. The matrices  $\mathbf{I}$  and  $\mathbf{0}$  are the identity and the zero matrix of size  $n \times n$ . Furthermore,  $\mathcal{S}(\vec{\mathbf{w}})$  contains penalty terms that enforce the boundary conditions.

116 The purpose of the SAT  $\mathcal{S}(\vec{w})$  is *i*) to enforce the boundary conditions in (2.3) and *ii*)  
 117 to stabilize the solution. One penalty term for each of the boundary conditions in (2.3) will  
 118 be constructed. Let  $k \in \{W, E, S, N\}$ . The SAT at boundary  $k$  that enforces the boundary  
 119 condition  $H^k \vec{w} = \vec{g}$  has the general form

$$120 \quad (3.3) \quad \mathcal{S}^k(\vec{w}) = (I_3 \otimes P^{-1}) \Sigma^k (I_3 \otimes P^k) (\mathcal{H}^k \vec{w} - \vec{g}).$$

121 In (3.3),  $\Sigma^k$  is the penalty matrix to be determined for stability at boundary  $k$ . The quadra-  
 122 tures are

$$123 \quad P^k = \begin{cases} E_{0x} \otimes P_y & \text{on the west boundary } (k = W) \\ E_N \otimes P_y & \text{on the east boundary } (k = E) \\ P_x \otimes E_{0y} & \text{on the south boundary } (k = S) \\ P_x \otimes E_M & \text{on the north boundary } (k = N). \end{cases}$$

124 For the boundary conditions listed in (2.3), the penalty terms are

$$\begin{aligned} \mathcal{S}^W(\vec{w}) &= (I_3 \otimes P^{-1}) \Sigma^W (I_3 \otimes P^W) \underbrace{\begin{pmatrix} u - g_1 \\ v - g_2 \\ u - g_1 \end{pmatrix}}_{\mathcal{H}^W \vec{w} - \vec{g}} \\ \mathcal{S}^E(\vec{w}) &= (I_3 \otimes P^{-1}) \Sigma^E (I_3 \otimes P^E) \underbrace{\begin{pmatrix} p - \epsilon D_x u - g_3 \\ -\epsilon D_x v - g_4 \\ 0 \end{pmatrix}}_{\mathcal{H}^E \vec{w} - \vec{g}} \\ \mathcal{S}^S(\vec{w}) &= (I_3 \otimes P^{-1}) \Sigma^S (I_3 \otimes P^S) \underbrace{\begin{pmatrix} u - 0 \\ v - 0 \\ v - 0 \end{pmatrix}}_{\mathcal{H}^S \vec{w} - 0} \\ \mathcal{S}^N(\vec{w}) &= (I_3 \otimes P^{-1}) \Sigma^N (I_3 \otimes P^N) \underbrace{\begin{pmatrix} -\epsilon D_y u - g_6 \\ p - \epsilon D_y v - g_5 \\ 0 \end{pmatrix}}_{\mathcal{H}^N \vec{w} - \vec{g}}, \end{aligned}$$

126 where

$$\begin{aligned} 127 \quad \Sigma^W &= \begin{pmatrix} -U/2 + \epsilon D_x^\top & 0 & 0 \\ 0 & -U/2 + \epsilon D_x^\top & 0 \\ 0 & 0 & -I \end{pmatrix}, & \Sigma^E &= (I_3 \otimes I) \\ 128 \quad \Sigma^S &= \begin{pmatrix} -V/2 + \epsilon D_y^\top & 0 & 0 \\ 0 & -V/2 + \epsilon D_y^\top & 0 \\ 0 & 0 & -I \end{pmatrix}, & \Sigma^N &= (I_3 \otimes I). \\ 129 \end{aligned}$$

As an example, the south penalty term can be written as

$$(3.5) \quad \mathcal{S}^S(\vec{w}) = (I_3 \otimes P^{-1}) \begin{pmatrix} -VP^S u/2 + \epsilon D_y^\top P^S u \\ -VP^S v/2 + \epsilon D_y^\top P^S v \\ -P^S v \end{pmatrix},$$

which is a more convenient notation for the derivation of the Jacobian in [Section 4](#).

We will show in the following section that this specific choice of penalty matrices leads to nonlinear stability. The total penalty term in [\(3.2\)](#) becomes

$$(3.6) \quad \mathcal{S}(\vec{w}) = \sum_{k \in \{W, E, S, N\}} \mathcal{S}(\vec{w})^k.$$

**3.1. Boundedness and Stability.** For completeness, we also show schematically how to obtain an energy estimate (again all details can be found in [\[12\]](#)). Similarly to the continuous analysis, we omit all boundaries except for the south one. By mimicking the continuous path [\[11\]](#), we multiply [\(3.2\)](#) by  $2\vec{w}^\top(I_3 \otimes P)$  from the left and use the SBP-property [\(3.1\)](#) to get

$$(3.7) \quad \frac{d}{dt} \|\vec{w}\|_{\tilde{I} \otimes P}^2 + 2\epsilon \|\nabla \vec{w}\|_{\tilde{I} \otimes P}^2 = BT,$$

where  $\|\nabla \vec{w}\|_{\tilde{I} \otimes P}^2 = (I_3 \otimes D_x \vec{w})^\top (I_3 \otimes P) \tilde{I} (I_3 \otimes D_x \vec{w}) + (I_3 \otimes D_y \vec{w})^\top (I_3 \otimes P) \tilde{I} (I_3 \otimes D_y \vec{w})$  is the dissipative volume term corresponding to the continuous one and

$$(3.8) \quad BT = \underbrace{\vec{w}^\top (I_3 \otimes P^S) B \vec{w} - 2\epsilon \vec{w}^\top (I_3 \otimes P^S) \tilde{I} (I_3 \otimes D_y) \vec{w}}_I \\ + \underbrace{2\vec{w}^\top (I_3 \otimes P) \mathcal{S}^S(\vec{w})}_{II}$$

contains all terms evaluated at the boundary.

The semi-norm of the solution ( $\|\vec{w}\|_{\tilde{I} \otimes P}^2$ ) is bounded if the right-hand side of [\(3.7\)](#) is non-positive. By expanding [\(3.8\)](#) and using the explicit form of  $\mathcal{S}^S(\vec{w})$  stated in [\(3.4\)](#), we find

$$BT = \underbrace{v^\top P^S (Uu + Vv + 2p - 2\epsilon D_y v) - 2\epsilon u^\top P^S D_y u}_I \\ - \underbrace{2v^\top P^S (Uu/2 + Vv/2 + p^\top v - \epsilon D_y v) + 2\epsilon u^\top P^S D_y u}_{II} = 0,$$

where term  $I$  is obtained from the governing equation and term  $II$  from the penalty term. As in the continuous setting, the boundary terms vanish. Integrating [\(3.7\)](#) in time (assuming homogeneous dissipative boundary conditions at all boundaries) leads to

$$(3.9) \quad \|\vec{w}\|_{\tilde{I} \otimes P}^2(T) + 2\epsilon \int_0^T \|\nabla \vec{w}\|_{\tilde{I} \otimes P}^2 dt \leq \|f\|_{\tilde{I} \otimes P}^2,$$

which is the semi-discrete version of the estimate in [\(2.5\)](#).

**4. Exact computation of the Jacobian.** In this section, we will explicitly compute the Jacobian of  $\mathcal{L}$  and  $\mathcal{S}$  in (3.2). Let  $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , where  $n = (N+1)(M+1)$  is the total number of grid points, be a differentiable vector function. For a given vector  $\mathbf{u} = (u_{00}, \dots, u_{NM})^\top \in \mathbb{R}^n$ ,  $\mathbf{h}$  outputs the vector  $\mathbf{h}(\mathbf{u}) = (h_{00}, \dots, h_{NM})^\top \in \mathbb{R}^n$ . The Jacobian matrix  $J_{\mathbf{h}} \in \mathbb{R}^{n \times n}$  of  $\mathbf{h}$  is given by

$$J_{\mathbf{h}} = \begin{pmatrix} \frac{\partial h_{00}}{\partial u_{00}} & \cdots & \frac{\partial h_{00}}{\partial u_{NM}} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_{NM}}{\partial u_{00}} & \cdots & \frac{\partial h_{NM}}{\partial u_{NM}} \end{pmatrix}.$$

We will first derive the Jacobian of the different terms in  $\mathcal{L}(\vec{w})$  and at the end, add the terms and state the complete result. To start, consider the vector function

$$\mathbf{h}(\mathbf{u}) = \begin{pmatrix} h_{00} \\ \vdots \\ h_{NM} \end{pmatrix} = \begin{pmatrix} u_{00} \\ \vdots \\ u_{NM} \end{pmatrix} = \mathbf{u}.$$

Since

$$\begin{array}{ccccc} \frac{\partial h_{00}}{\partial u_{00}} = 1 & \frac{\partial h_{00}}{\partial u_{01}} = 0 & \frac{\partial h_{00}}{\partial u_{02}} = 0 & \cdots & \frac{\partial h_{00}}{\partial u_{NM}} = 0 \\ \frac{\partial h_{01}}{\partial u_{00}} = 0 & \frac{\partial h_{01}}{\partial u_{01}} = 1 & \frac{\partial h_{01}}{\partial u_{02}} = 0 & \cdots & \frac{\partial h_{01}}{\partial u_{NM}} = 0 \\ & \vdots & & & \\ \frac{\partial h_{NM}}{\partial u_{00}} = 0 & \frac{\partial h_{NM}}{\partial u_{01}} = 0 & \frac{\partial h_{NM}}{\partial u_{02}} = 0 & \cdots & \frac{\partial h_{NM}}{\partial u_{NM}} = 1, \end{array}$$

the Jacobian of  $\mathbf{h}(\mathbf{u}) = \mathbf{u}$  becomes  $J_{\mathbf{u}} = \mathbf{I}$ . Now let

$$\begin{aligned} \mathbf{h}(\mathbf{u}) &= \begin{pmatrix} h_{00} \\ \vdots \\ h_{NM} \end{pmatrix} = \mathbf{D}_{\mathbf{x}} \mathbf{u} = \begin{pmatrix} D_{0,0} & \cdots & D_{0,NM} \\ \vdots & \ddots & \vdots \\ D_{NM,0} & \cdots & D_{NM,MN} \end{pmatrix} \begin{pmatrix} u_{00} \\ \vdots \\ u_{NM} \end{pmatrix} \\ &= \begin{pmatrix} D_{0,0}u_{00} + \cdots + D_{0,NM}u_{NM} \\ \vdots \\ D_{NM,0}u_{00} + \cdots + D_{NM,MN}u_{NM} \end{pmatrix}. \end{aligned}$$

Then, in the same way

$$\begin{array}{ccccc} \frac{\partial h_{00}}{\partial u_{00}} = D_{0,0} & \frac{\partial h_{00}}{\partial u_{01}} = D_{0,1} & \frac{\partial h_{00}}{\partial u_{02}} = D_{0,2} & \cdots & \frac{\partial h_{00}}{\partial u_{NM}} = D_{0,NM} \\ \frac{\partial h_{01}}{\partial u_{00}} = D_{1,0} & \frac{\partial h_{01}}{\partial u_{01}} = D_{1,1} & \frac{\partial h_{01}}{\partial u_{02}} = D_{1,2} & \cdots & \frac{\partial h_{01}}{\partial u_{NM}} = D_{1,NM} \\ & \vdots & & & \\ \frac{\partial h_{NM}}{\partial u_{00}} = D_{NM,0} & \frac{\partial h_{NM}}{\partial u_{01}} = D_{NM,1} & \frac{\partial h_{NM}}{\partial u_{02}} = D_{NM,2} & \cdots & \frac{\partial h_{NM}}{\partial u_{NM}} = D_{NM,NM}. \end{array}$$

Thus,  $J_{\mathbf{D}_{\mathbf{x}} \mathbf{u}} = \mathbf{D}_{\mathbf{x}}$ .

To derive the Jacobian of the nonlinear term,  $\mathbf{U}\mathbf{D}_x\mathbf{u}$ , we let

$$\mathbf{h}(\mathbf{u}) = \mathbf{U}\mathbf{D}_x\mathbf{u} = \begin{pmatrix} u_{00}[D_{0,0}u_{00} + \cdots + D_{0,NM}u_{NM}] \\ \vdots \\ u_{NM}[D_{NM,0}u_{00} + \cdots + D_{NM,MN}u_{NM}] \end{pmatrix} = \begin{pmatrix} u_{00}(\mathbf{D}_x\mathbf{u})_{00} \\ \vdots \\ u_{NM}(\mathbf{D}_x\mathbf{u})_{NM} \end{pmatrix}.$$

By using the product rule, we get that

$$\begin{aligned} \frac{\partial h_{00}}{\partial u_{00}} &= u_{00}D_{0,0} + (\mathbf{D}_x\mathbf{u})_{00} & \frac{\partial h_{00}}{\partial u_{01}} &= u_{00}D_{0,1} & \cdots & \frac{\partial h_{00}}{\partial u_{NM}} = u_{00}D_{0,NM} \\ \frac{\partial h_{01}}{\partial u_{00}} &= u_{01}D_{1,0} & \frac{\partial h_{01}}{\partial u_{01}} &= u_{01}D_{1,1} + (\mathbf{D}_x\mathbf{u})_{01} & \cdots & \frac{\partial h_{01}}{\partial u_{NM}} = u_{00}D_{0,NM} \\ & \vdots & & & & \\ \frac{\partial h_{NM}}{\partial u_{00}} &= u_{NM}D_{NM,0} & \frac{\partial h_{NM}}{\partial u_{01}} &= u_{NM}D_{NM,1} & \cdots & \frac{\partial h_{NM}}{\partial u_{NM}} = u_{NM}D_{NM,NM} + (\mathbf{D}_x\mathbf{u})_{NM}. \end{aligned}$$

Hence,

$$\begin{aligned} J_{\mathbf{U}\mathbf{D}_x\mathbf{u}} &= \begin{pmatrix} u_{00}D_{0,0} + (\mathbf{D}_x\mathbf{u})_{00} & u_{00}D_{0,1} & \cdots & u_{00}D_{0,NM} \\ u_{01}D_{1,0} & u_{01}D_{1,1} + (\mathbf{D}_x\mathbf{u})_{01} & \cdots & u_{01}D_{1,NM} \\ \vdots & \vdots & \ddots & \vdots \\ u_{NM}D_{NM,0} & u_{NM}D_{NM,1} & \cdots & u_{NM}D_{NM,NM} + (\mathbf{D}_x\mathbf{u})_{NM} \end{pmatrix} \\ &= \underbrace{\begin{pmatrix} u_{00}D_{0,0} & u_{00}D_{0,1} & \cdots & u_{00}D_{0,NM} \\ u_{01}D_{1,0} & u_{01}D_{1,1} & \cdots & u_{01}D_{1,NM} \\ \vdots & \vdots & \ddots & \vdots \\ u_{NM}D_{NM,0} & u_{NM}D_{NM,1} & \cdots & u_{NM}D_{NM,NM} \end{pmatrix}}_{\mathbf{U}\mathbf{D}_x} \\ &\quad + \underbrace{\begin{pmatrix} (\mathbf{D}_x\mathbf{u})_{00} & 0 & \cdots & 0 \\ 0 & (\mathbf{D}_x\mathbf{u})_{01} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & (\mathbf{D}_x\mathbf{u})_{NM} \end{pmatrix}}_{\underline{\mathbf{D}_x\mathbf{u}}} = \mathbf{U}\mathbf{D}_x + \underline{\mathbf{D}_x\mathbf{u}}, \end{aligned}$$

where  $\underline{\mathbf{D}_x\mathbf{u}} = \text{diag}(\mathbf{D}_x\mathbf{u})$ .

In a similar manner, for

$$\mathbf{h}(\mathbf{u}) = \mathbf{D}_x\mathbf{U}\mathbf{u} = \begin{pmatrix} D_{0,0}u_{00}^2 + \cdots + D_{0,NM}u_{NM}^2 \\ \vdots \\ D_{NM,0}u_{00}^2 + \cdots + D_{NM,MN}u_{NM}^2 \end{pmatrix}$$



we get that

$$\begin{aligned}
 \frac{\partial h_{00}}{\partial u_{00}} &= 2D_{0,0}u_{00} & \frac{\partial h_{00}}{\partial u_{01}} &= 2D_{0,1}u_{01} & \dots & \frac{\partial h_{00}}{\partial u_{NM}} &= 2D_{0,NM}u_{NM} \\
 \frac{\partial h_{01}}{\partial u_{00}} &= 2D_{1,0}u_{00} & \frac{\partial h_{01}}{\partial u_{01}} &= 2D_{1,1}u_{01} & \dots & \frac{\partial h_{01}}{\partial u_{NM}} &= 2D_{1,NM}u_{NM} \\
 &\vdots & & & & & \\
 \frac{\partial h_{NM}}{\partial u_{00}} &= 2D_{NM,0}u_{00} & \frac{\partial h_{NM}}{\partial u_{01}} &= 2D_{NM,1}u_{01} & \dots & \frac{\partial h_{NM}}{\partial u_{NM}} &= 2D_{NM,NM}u_{NM}.
 \end{aligned}$$

Hence,  $J_{D_x U} u = 2D_x U$ . To summarize, we have shown that

$$(4.1) \quad J_u = I, \quad J_{D_x u} = D_x, \quad J_{UD_x u} = UD_x + \underline{D_x u}, \quad J_{D_x U} u = 2D_x U.$$

**4.1. The Jacobian of the spatial operator.** Having established these building blocks, we next consider the terms in  $\mathcal{L}(\vec{w})$ . Since these terms have a block structure, their Jacobians will have that as well. Let  $\mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3 : \mathbb{R}^{3n} \rightarrow \mathbb{R}^n$  be differentiable functions of  $\vec{w}$  and define  $\tilde{\mathbf{h}} : \mathbb{R}^{3n} \rightarrow \mathbb{R}^{3n}$  given by

$$\tilde{\mathbf{h}}(\vec{w}) = \begin{pmatrix} \mathbf{h}^1(\vec{w}) \\ \mathbf{h}^2(\vec{w}) \\ \mathbf{h}^3(\vec{w}) \end{pmatrix}.$$

Since

$$\mathbf{h}^1 = \begin{pmatrix} h_{00}^1 \\ h_{01}^1 \\ \vdots \\ h_{NM}^1 \end{pmatrix} \quad \text{and} \quad \vec{w} = \begin{pmatrix} u \\ v \\ p \end{pmatrix}$$

it follows that

$$J_{h_{00}^1} = \frac{\partial h_{00}^1}{\partial \vec{w}} = \left( \frac{\partial h_{00}^1}{\partial u_{00}} \quad \dots \quad \frac{\partial h_{00}^1}{\partial u_{3NM}} \right) = \left( \frac{\partial h_{00}^1}{\partial u} \quad \frac{\partial h_{00}^1}{\partial v} \quad \frac{\partial h_{00}^1}{\partial p} \right) \in \mathbb{R}^{1 \times 3n}$$

and similarly for every element in  $\mathbf{h}^1$ . Therefore, the Jacobian of  $\mathbf{h}^1$  can be expressed as

$$J_{\mathbf{h}^1} = \frac{\partial \mathbf{h}^1}{\partial \vec{w}} = \begin{pmatrix} \frac{\partial h_{00}^1}{\partial u} & \frac{\partial h_{00}^1}{\partial v} & \frac{\partial h_{00}^1}{\partial p} \\ \frac{\partial h_{01}^1}{\partial u} & \frac{\partial h_{01}^1}{\partial v} & \frac{\partial h_{01}^1}{\partial p} \\ \vdots & \vdots & \vdots \\ \frac{\partial h_{NM}^1}{\partial u} & \frac{\partial h_{NM}^1}{\partial v} & \frac{\partial h_{NM}^1}{\partial p} \end{pmatrix} = \begin{pmatrix} \frac{\partial \mathbf{h}^1}{\partial u} & \frac{\partial \mathbf{h}^1}{\partial v} & \frac{\partial \mathbf{h}^1}{\partial p} \end{pmatrix} \in \mathbb{R}^{n \times 3n}.$$

The same holds for  $\mathbf{h}^2$  and  $\mathbf{h}^3$ . Thus, the Jacobian of  $\tilde{\mathbf{h}}$  is given by

$$J_{\tilde{\mathbf{h}}} = \begin{pmatrix} \frac{\partial \mathbf{h}^1}{\partial u} & \frac{\partial \mathbf{h}^1}{\partial v} & \frac{\partial \mathbf{h}^1}{\partial p} \\ \frac{\partial \mathbf{h}^2}{\partial u} & \frac{\partial \mathbf{h}^2}{\partial v} & \frac{\partial \mathbf{h}^2}{\partial p} \\ \frac{\partial \mathbf{h}^3}{\partial u} & \frac{\partial \mathbf{h}^3}{\partial v} & \frac{\partial \mathbf{h}^3}{\partial p} \end{pmatrix} \in \mathbb{R}^{3n \times 3n}$$

and each block in  $J_{\tilde{\mathbf{h}}}$  is of size  $n \times n$ .

For the first term in  $\mathcal{L}(\vec{w})$ ,  $A(I_3 \otimes D_x) \vec{w}$ , we get

$$\tilde{h}(\vec{w}) = \begin{pmatrix} h^1(\vec{w}) \\ h^2(\vec{w}) \\ h^3(\vec{w}) \end{pmatrix} = A(I_3 \otimes D_x) \vec{w} = \begin{pmatrix} UD_x u + D_x p \\ UD_x v \\ D_x u \end{pmatrix} = \begin{pmatrix} UD_x u + D_x p \\ D_x v u \\ D_x u \end{pmatrix}.$$

The last identities are useful when deriving  $J_{A(I_3 \otimes D_x) \vec{w}}$ . By using (4.1), we get that

$$\begin{aligned} \frac{\partial h^1}{\partial u} &= UD_x + \underline{D_x u} & \frac{\partial h^1}{\partial v} &= 0 & \frac{\partial h^1}{\partial p} &= D_x \\ \frac{\partial h^2}{\partial u} &= \underline{D_x v} & \frac{\partial h^2}{\partial v} &= UD_x & \frac{\partial h^2}{\partial p} &= 0 \\ \frac{\partial h^3}{\partial u} &= D_x & \frac{\partial h^3}{\partial v} &= 0 & \frac{\partial h^3}{\partial p} &= 0. \end{aligned}$$

Thus,

$$J_{A(I_3 \otimes D_x) \vec{w}} = \begin{pmatrix} UD_x + \underline{D_x u} & 0 & D_x \\ \underline{D_x v} & UD_x & 0 \\ D_x & 0 & 0 \end{pmatrix}.$$

Likewise for the second term in  $\mathcal{L}(\vec{w})$ ,  $(I_3 \otimes D_x) A \vec{w}$ , note that

$$(I_3 \otimes D_x) A \vec{w} = \begin{pmatrix} D_x U u + D_x p \\ D_x U v \\ D_x u \end{pmatrix} = \begin{pmatrix} D_x U u + D_x p \\ D_x V u \\ D_x u \end{pmatrix},$$

where we have used that  $Vu = Uv$ . Hence,

$$J_{(I_3 \otimes D_x) A \vec{w}} = \begin{pmatrix} 2D_x U & 0 & D_x \\ D_x V & D_x U & 0 \\ D_x & 0 & 0 \end{pmatrix}.$$

The next two terms in  $\mathcal{L}(\vec{w})$  are treated in a similar manner and we get that

$$J_{B(I_3 \otimes D_y) \vec{w}} = \begin{pmatrix} VD_y & \underline{D_y u} & 0 \\ 0 & VD_y + \underline{D_y v} & D_y \\ 0 & D_y & 0 \end{pmatrix}$$

$$J_{(I_3 \otimes D_y) B \vec{w}} = \begin{pmatrix} D_y V & D_y U & 0 \\ 0 & 2D_y V & D_y \\ 0 & D_y & 0 \end{pmatrix}.$$

Finally, the contribution to the Jacobian of the linear viscous terms simply becomes

$$J_{-\epsilon \tilde{I}[(I_3 \otimes D_x)^2 + (I_3 \otimes D_y)^2] \vec{w}} = - \begin{pmatrix} \epsilon(D_x^2 + D_y^2) & 0 & 0 \\ 0 & \epsilon(D_x^2 + D_y^2) & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

Adding all terms proves the following proposition, which is the first of the two main results of this paper.

**Proposition 4.1.** *The Jacobian  $J_{\mathcal{L}}$  of the discrete operator  $\mathcal{L}$  in (3.2) is*

$$(4.2) \quad J_{\mathcal{L}} = \begin{pmatrix} J_{11} & \frac{1}{2}(\underline{D_y u} + D_y U) & D_x \\ \frac{1}{2}(\underline{D_x v} + D_x V) & J_{22} & D_y \\ D_x & D_y & 0 \end{pmatrix}$$

where

$$J_{11} = \frac{1}{2} (UD_x + \underline{D_x u} + 2D_x U + VD_y + D_y V) - \epsilon(D_x^2 + D_y^2)$$

$$J_{22} = \frac{1}{2} (VD_y + \underline{D_y v} + 2D_y V + UD_x + D_x U) - \epsilon(D_x^2 + D_y^2).$$

**4.2. The Jacobian of the penalty terms.** By following the procedure presented above, we next derive the Jacobian for  $\mathcal{S}(\vec{w})$ . To start, we rewrite  $\mathcal{S}^S(\vec{w})$  as

$$(4.3) \quad \mathcal{S}^S(\vec{w}) = \begin{pmatrix} \mathcal{S}_1^S \\ \mathcal{S}_2^S \\ \mathcal{S}_3^S \end{pmatrix} = (I_3 \otimes P^{-1}) \begin{pmatrix} -VP^S u/2 + \epsilon D_y^\top P^S u \\ -VP^S v/2 + \epsilon D_y^\top P^S v \\ -P^S v \end{pmatrix} \in \mathbb{R}^{3n}.$$

The Jacobian of  $\mathcal{S}^S(\vec{w})$  is

$$J_{\mathcal{S}^S} = \begin{pmatrix} \frac{\partial \mathcal{S}_1^S}{\partial u} & \frac{\partial \mathcal{S}_1^S}{\partial v} & \frac{\partial \mathcal{S}_1^S}{\partial p} \\ \frac{\partial \mathcal{S}_2^S}{\partial u} & \frac{\partial \mathcal{S}_2^S}{\partial v} & \frac{\partial \mathcal{S}_2^S}{\partial p} \\ \frac{\partial \mathcal{S}_3^S}{\partial u} & \frac{\partial \mathcal{S}_3^S}{\partial v} & \frac{\partial \mathcal{S}_3^S}{\partial p} \end{pmatrix} \in \mathbb{R}^{3n \times 3n}.$$

The first block in  $J_{\mathcal{S}^S}$  becomes

$$\frac{\partial \mathcal{S}_1^S}{\partial u} = \left( -\underbrace{\frac{\partial}{\partial u} P^{-1} V P^S u/2}_{=P^{-1} V P^S/2} + \underbrace{\frac{\partial}{\partial u} \epsilon P^{-1} D_y^\top P^S u}_{=\epsilon P^{-1} D_y^\top P^S} \right) = P^{-1} (-V/2 + \epsilon D_y^\top) P^S.$$

Since  $P^S$  is diagonal, we have  $VP^S u = UP^S v$  and the second block is

$$\frac{\partial \mathcal{S}_2^S}{\partial v} = \left( -\underbrace{\frac{\partial}{\partial v} P^{-1} U P^S v/2}_{=P^{-1} U P^S/2} + \underbrace{\frac{\partial}{\partial v} \epsilon P^{-1} D_x^\top P^S u}_{=0} \right) = -P^{-1} U P^S/2.$$

Note that  $\mathcal{S}^S$  does not depend on  $p$  and also that  $\mathcal{S}_2^S$  and  $\mathcal{S}_3^S$  are both independent of  $u$ .

Hence, the remaining non-zero blocks of  $J_{\mathcal{S}^S}$  are

$$\frac{\partial \mathcal{S}_2^S}{\partial v} = P^{-1} (-V + \epsilon D_y^\top) P^S, \quad \frac{\partial \mathcal{S}_3^S}{\partial v} = -P^{-1} P^S,$$

where we have used that  $VP^S v = P^S V v$ . Therefore,

$$J_{\mathcal{S}^S} = (I_3 \otimes P^{-1}) \begin{pmatrix} -V/2 + \epsilon D_y^\top & -U/2 & 0 \\ 0 & -V + \epsilon D_y^\top & 0 \\ 0 & -I & 0 \end{pmatrix} (I_3 \otimes P^S).$$

For non-homogeneous boundary conditions, the boundary data  $\mathbf{g}$  will affect the Jacobian if the SAT:s are nonlinear with respect to  $\vec{\mathbf{w}}$ . We illustrate this by considering  $\mathcal{S}_1^W$  (the first block in  $\mathcal{S}^W$ ), which we rewrite in a similar manner as we did for  $\mathcal{S}_1^S$  in (4.3) and get

$$\begin{aligned}\mathcal{S}_1^W(\vec{\mathbf{w}}) &= \mathbf{P}^{-1}(-\mathbf{U}/2 + \epsilon \mathbf{D}_x^\top) \mathbf{P}^W (\mathbf{u} - \mathbf{g}_1) \\ &= -\mathbf{P}^{-1} \mathbf{U} \mathbf{P}^W (\mathbf{u} - \mathbf{g}_1)/2 + \epsilon \mathbf{P}^{-1} \mathbf{D}_x^\top \mathbf{P}^W (\mathbf{u} - \mathbf{g}_1).\end{aligned}$$

Note that the terms  $-\mathbf{P}^{-1} \mathbf{U} \mathbf{P}^W (\mathbf{u} - \mathbf{g}_1)/2$  and  $\epsilon \mathbf{P}^{-1} \mathbf{D}_x^\top \mathbf{P}^W (\mathbf{u} - \mathbf{g}_1)$  are nonlinear and linear with respect to  $\vec{\mathbf{w}}$  (via  $\mathbf{u}$ ), respectively. The Jacobian to the linear term simply becomes

$$\frac{\partial}{\partial \mathbf{u}} \left( \epsilon \mathbf{P}^{-1} \mathbf{D}_x^\top \mathbf{P}^W (\mathbf{u} - \mathbf{g}_1) \right) = \underbrace{\frac{\partial}{\partial \mathbf{u}} \left( \epsilon \mathbf{P}^{-1} \mathbf{D}_x^\top \mathbf{P}^W \mathbf{u} \right)}_{=\epsilon \mathbf{P}^{-1} \mathbf{D}_x^\top \mathbf{P}^W} - \underbrace{\frac{\partial}{\partial \mathbf{u}} \left( \epsilon \mathbf{P}^{-1} \mathbf{D}_x^\top \mathbf{P}^W \mathbf{g}_1 \right)}_{=0}.$$

For the nonlinear term, we use that  $\mathbf{U} \mathbf{P}^W = \mathbf{P}^W \mathbf{U}$  and  $\mathbf{U} \mathbf{g}_1 = \underline{\mathbf{g}}_1 \mathbf{u}$ , which yield

$$\begin{aligned}-\frac{\partial}{\partial \mathbf{u}} \left( \mathbf{P}^{-1} \mathbf{P}^W \mathbf{U} (\mathbf{u} - \mathbf{g}_1)/2 \right) &= -\frac{\partial}{\partial \mathbf{u}} \left( \mathbf{P}^{-1} \mathbf{P}^W \mathbf{U} \mathbf{u} \right)/2 + \frac{\partial}{\partial \mathbf{u}} \left( \mathbf{P}^{-1} \mathbf{P}^W \underline{\mathbf{g}}_1 \mathbf{u} \right)/2 \\ &= \underbrace{-\mathbf{P}^{-1} \mathbf{P}^W \mathbf{U}}_{=-\mathbf{P}^{-1} \mathbf{P}^W \mathbf{U}} + \underbrace{\mathbf{P}^{-1} \mathbf{P}^W \underline{\mathbf{g}}_1/2}_{=\mathbf{P}^{-1} \mathbf{P}^W \underline{\mathbf{g}}_1/2} \\ &= \mathbf{P}^{-1} (-\mathbf{U} + \underline{\mathbf{g}}_1/2) \mathbf{P}^W\end{aligned}$$

Since  $\mathcal{S}_1^W$  is independent of both  $\mathbf{v}$  and  $\mathbf{p}$ , its Jacobian becomes

$$J_{\mathcal{S}_1^W}(\vec{\mathbf{w}}) = \begin{pmatrix} \mathbf{P}^{-1}(-(\mathbf{U} - \underline{\mathbf{g}}_1/2) + \epsilon \mathbf{D}_x^\top) \mathbf{P}^W & \mathbf{0} & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{n \times 3n}$$

The Jacobian of the other penalty terms are derived in a similar manner and we have therefore proved the second main result of this paper.

**Proposition 4.2.** *The Jacobian of the total penalty term (3.6) is*

$$(4.4) \quad J_{\mathcal{S}}(\vec{\mathbf{w}}) = \sum_{k \in \{W, E, S, N\}} J_{\mathcal{S}^k}(\vec{\mathbf{w}}),$$

where

$$\begin{aligned}J_{\mathcal{S}^W}(\vec{\mathbf{w}}) &= (\mathbf{I}_3 \otimes \mathbf{P}^{-1}) \begin{pmatrix} -(\mathbf{U} - \underline{\mathbf{g}}_1/2) + \epsilon \mathbf{D}_x^\top & \mathbf{0} & \mathbf{0} \\ -(\mathbf{V} - \underline{\mathbf{g}}_2)/2 & -\mathbf{U}/2 + \epsilon \mathbf{D}_x^\top & \mathbf{0} \\ -\mathbf{I} & \mathbf{0} & \mathbf{0} \end{pmatrix} (\mathbf{I}_3 \otimes \mathbf{P}^W) \\ J_{\mathcal{S}^E}(\vec{\mathbf{w}}) &= (\mathbf{I}_3 \otimes \mathbf{P}^{-1} \mathbf{P}^E) \begin{pmatrix} -\epsilon \mathbf{D}_x & \mathbf{0} & \mathbf{I} \\ \mathbf{0} & -\epsilon \mathbf{D}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \\ J_{\mathcal{S}^S}(\vec{\mathbf{w}}) &= (\mathbf{I}_3 \otimes \mathbf{P}^{-1} \mathbf{P}^S) \begin{pmatrix} -\mathbf{V}/2 + \epsilon \mathbf{D}_y^\top & -\mathbf{U}/2 & \mathbf{0} \\ \mathbf{0} & -\mathbf{V} + \epsilon \mathbf{D}_y^\top & \mathbf{0} \\ \mathbf{0} & -\mathbf{I} & \mathbf{0} \end{pmatrix} (\mathbf{I}_3 \otimes \mathbf{P}^S) \\ J_{\mathcal{S}^N}(\vec{\mathbf{w}}) &= (\mathbf{I}_3 \otimes \mathbf{P}^{-1} \mathbf{P}^N) \begin{pmatrix} -\epsilon \mathbf{D}_y & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & -\epsilon \mathbf{D}_y & \mathbf{I} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}.\end{aligned}$$

*Remark 4.3.* We see from [Proposition 4.1](#) and [Proposition 4.2](#) that parts of the blocks in the Jacobian of both  $J_{\mathcal{L}}$  and  $J_{\mathcal{S}}$  are obtained directly from the construction of  $\mathcal{L}$ . The few remaining parts are obtained by *i*) matrix multiplications between a diagonal matrix and a non-diagonal one, for example  $U\mathbf{D}_x$  and *ii*) matrix additions. This leads to few new additional operations and hence efficiency.

**5. The fully discrete scheme.** To evolve the system [\(3.2\)](#) in time, we will for simplicity and ease of explanation use the implicit Backward Euler method. More accurate and efficient methods could be used in the same manner in practice. For an ordinary differential system of equations of the form

$$\mathcal{M}\phi_t + \mathcal{H}(\phi) = 0,$$

where  $\phi$  is a function defined on the grid and  $\mathcal{M}$  is a constant matrix, the backward Euler schemes becomes

$$(5.1) \quad \frac{\mathcal{M}(\phi^{i+1} - \phi^i)}{\Delta t} + \mathcal{H}(\phi^{i+1}) = 0.$$

In [\(5.1\)](#),  $\Delta t$  is the size of the time step and the superindices  $i$  and  $i + 1$  are the solution at time level  $i$  and  $i + 1$ , respectively.

In order to obtain  $\phi^{i+1}$ , the system of nonlinear equations in [\(5.1\)](#) must be solved. One strategy is to first form the function in [\(1.1\)](#), which results in

$$(5.2) \quad \mathcal{F}(\phi^{i+1}) = \frac{\mathcal{M}(\phi^{i+1} - \phi^i)}{\Delta t} + \mathcal{H}(\phi^{i+1}).$$

If we find a vector  $\phi^*$  such that  $\mathcal{F}(\phi^*) = 0$ , then  $\phi^{i+1} = \phi^*$ . To solve [\(5.2\)](#), we employ Newton's method [\[15\]](#), which is described in [Algorithm 5.1](#). This allows us to solve a sequence of linear systems of equations and arrive at an approximation of  $\phi^{i+1}$ .

---

**Algorithm 5.1** Newton's method

---

Input:  $\phi^0$  and tolerance  $tol$

Output: An approximation of  $\phi^*$ , where  $\mathcal{F}(\phi^*) = 0$

**for**  $j = 0, 1, 2, \dots$  **do**

    solve  $J_{\mathcal{F}}(\phi^j)\mathbf{h}^j = -\mathcal{F}(\phi^j)$

    set  $\phi^{j+1} = \phi^j + \mathbf{h}^j$

**if**  $\|\mathcal{F}(\phi^{j+1})\| < tol$  **then**

**return**  $\phi^{j+1}$

**end if**

**end for**

---

For the INS equations,  $\phi = \vec{w}$ ,  $\mathcal{H}(\phi) = \mathcal{L}(\phi) - \mathcal{S}(\phi)$ , and  $\mathcal{M} = \tilde{\mathbf{I}}$ . Hence, [\(5.2\)](#) becomes

$$(5.3) \quad \mathcal{F}(\vec{w}^{i+1}) = \frac{1}{\Delta t} \left( \begin{bmatrix} \mathbf{u}^{i+1} \\ \mathbf{v}^{i+1} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{u}^i \\ \mathbf{v}^i \\ \mathbf{0} \end{bmatrix} \right) + \mathcal{L}(\vec{w}^{i+1}) - \mathcal{S}(\vec{w}^{i+1}).$$

Furthermore,  $J_{\mathcal{H}}(\vec{w}) = J_{\mathcal{L}}(\vec{w}) - J_{\mathcal{S}}(\vec{w})$ , which yields

$$(5.4) \quad J_{\mathcal{F}}(\vec{w}) = \frac{1}{\Delta t} \tilde{\mathbf{I}} + J_{\mathcal{L}}(\vec{w}) - J_{\mathcal{S}}(\vec{w}),$$

to be used in the Newton iterations. In (5.4),  $J_{\mathcal{L}}(\vec{w})$  and  $J_{\mathcal{S}}(\vec{w})$  are given in Proposition 4.1 and Proposition 4.2, respectively.

**6. Numerical Experiments.** A simple finite-difference approximation of the Jacobian is given by [8]

$$(6.1) \quad J_{i,j} \approx \hat{J}_{i,j} = \frac{\mathcal{F}_i(\vec{w} + \delta_j \mathbf{e}_j) - \mathcal{F}_i(\vec{w})}{\delta_j}.$$

The approximation in (6.1) was used during the implementation of the analytical expression of  $J_{\mathcal{F}}$  since we expected  $\|J - \hat{J}\|_{\infty}$  to be small. This allowed us to write unit tests ensuring that the Jacobian has been correctly implemented, by comparing it to the approximation. In (6.1), a small  $\delta$  leads to a good approximation. However, note that if  $\delta$  is chosen too small, the approximation will be contaminated by floating-point roundoff errors, which limits the practically achievable accuracy of  $J$  [8].

Computing difference approximations of the Jacobian also allowed us to compare the efficiency of Newton's method using approximate versus analytical Jacobians. Note that computing the approximation (6.1) requires  $n$  evaluations of  $\mathcal{F}$ , resulting in  $O(n^2)$  complexity, compared to the  $O(n)$  complexity of evaluating the exact Jacobian. Table 1 shows the execution times for evaluating the analytical Jacobian versus computing the approximation (6.1) at increasing resolutions.

| Resolution     | Exact  | FD approximation |
|----------------|--------|------------------|
| $5 \times 5$   | 0.005s | 0.858s           |
| $10 \times 10$ | 0.006s | 13.85s           |
| $15 \times 15$ | 0.006s | 76.49s           |
| $20 \times 20$ | 0.007s | 261.6s           |

Table 1

Execution times for computing the exact Jacobian of  $\mathcal{F}$  versus the finite difference approximation (6.1). Even at low resolutions, using difference approximations of the Jacobian is clearly unrealistic.

As expected, due to the large number of evaluations of  $\mathcal{F}$  needed to compute the approximation, such a strategy quickly becomes infeasible.

It is readily seen that the number of floating point operations needed to evaluate the discrete spatial operator  $\mathcal{L}$  grows linearly with the degrees of freedom  $n$ . Consider for example the term  $\mathbf{A}(I_3 \otimes \mathbf{D}_{\mathbf{x}}) \vec{w}$ . The first product,  $(I_3 \otimes \mathbf{D}_{\mathbf{x}}) \vec{w}$ , are finite difference approximations at each point in the grid, resulting in  $Cn$  operations, where  $C$  depends on the width of the difference stencil. The matrix  $\mathbf{A}$  is a 3-by-3 block matrix with diagonal blocks, and so results in another  $O(n)$  number of operations. Analogously, the remaining terms in  $\mathcal{L}$  each contribute  $O(n)$  operations. The arithmetic complexity of evaluating a penalty term  $\mathcal{S}$  is  $O(\sqrt{n})$  (assuming equal resolution in the horizontal and vertical directions), since  $\mathcal{S}$  acts only on the grid boundary. Hence, the arithmetic complexity of evaluating  $\mathcal{F}$  is  $O(n)$ .

Let us study the arithmetic complexity of evaluating the Jacobian  $J_{\mathcal{F}}$  of  $\mathcal{F}$ . Inspecting the form of the Jacobian  $J_{\mathcal{L}}$  in Proposition 4.1 we see a number of terms that need to be evaluated. The partial derivatives  $\underline{D}_x u$ ,  $\underline{D}_y v$ , etc, have already been computed as part of the evaluation of  $\mathcal{L}$ , and hence can be disregarded. Similarly, terms that do not depend on the solution, such as  $\underline{D}_x$ ,  $\underline{D}_x^2$ , etc, can be disregarded since they remain constant throughout the simulation. Finally we have terms of the type  $\underline{U} \underline{D}_x$ ,  $\underline{D}_y \underline{V}$ , etc. These are all products of a diagonal matrix and a banded difference stencil matrix, and each contribute with  $O(n)$  operations. Summing the terms uses  $O(n)$  operations. Therefore, the arithmetic complexity of evaluating  $J_{\mathcal{L}}$  is  $O(n)$ . In fact, the number of operations needed to evaluate products like  $\underline{U} \underline{D}_x$  or  $\underline{D}_y \underline{V}$  do not exceed the number of operations needed to compute the discrete partial derivatives involved in  $\mathcal{L}$ . Hence, the cost ratio of evaluating  $J_{\mathcal{L}}$  and evaluating  $\mathcal{L}$  is less than 1 (i.e. the additional cost of evaluating  $J_{\mathcal{L}}$  is small). As before, the arithmetic complexity of evaluating the Jacobian with respect to a boundary penalty  $\mathcal{S}$  is  $O(\sqrt{n})$  since it acts only on the boundary of the grid. Thus, the total arithmetic complexity of evaluating  $J_{\mathcal{F}}$  is less than the cost of evaluating  $\mathcal{F}$ .

**6.1. The order of accuracy.** The method of manufactured solution [16] is used to verify the implementation. In all computations in this subsection, the initial guess is the solution from the previous time step and the tolerance  $tol$  in Algorithm 5.1 is set to  $10^{-12}$ . For the SBP-operators SBP21 and SBP42, the expected orders of accuracy for the system (3.2) are 2 and 3, respectively [17, 19]. The manufactured solution we have used is

$$\begin{aligned} u &= 1 + 0.1 \sin(3\pi x - 0.01t) \sin(3\pi y - 0.01t) \\ v &= \sin(3\pi x - 0.01t) \sin(3\pi y - 0.01t) \\ p &= \cos(3\pi x - 0.01t) \cos(3\pi y - 0.01t). \end{aligned} \quad (6.2)$$

Inserting (6.2) into (2.1) leads to a non-zero right-hand side  $\vec{k}(t, x, y)$ , which is evaluated on the grid and added to the right-hand side of (3.2) by the vector  $\vec{k}(t)$ . Since  $\vec{k}$  is independent of  $\vec{w}$ , it does not affect the Jacobian. The initial and boundary data are also taken from (6.2). The step size is chosen to be  $\Delta t = 10^{-5}$  and the computations are terminated at  $t = 1$ . Next, we compute the pointwise error vector  $\vec{e}$  and its  $L_2$ -norm  $\|\vec{e}\|_{I_3 \otimes P}$ . The spatial convergence rate for the SBP operators is given by  $r = \log(\|\vec{e}\|_i / \|\vec{e}\|_j) / \log((j-1)/(i-1))$ , where  $i$  and  $j$  refer to the number of grid points in both spatial dimensions. The order of accuracy in space are presented in Table 2 and agree well with theory.

Next, we consider the steady-state problem of (2.1) and (3.2), which means that the goal is to find  $\vec{w}^*$  such that

$$\mathcal{L}(\vec{w}^*) = \mathcal{S}(\vec{w}^*). \quad (6.3)$$

As before, we want to find an approximation to the vector  $\vec{w}^*$  which satisfies

$$\mathcal{F}(\vec{w}^*) = \mathcal{L}(\vec{w}^*) - \mathcal{S}(\vec{w}^*) = 0. \quad (6.4)$$

The Jacobian of  $\mathcal{F}$  is  $J_{\mathcal{F}}(\vec{w}) = J_{\mathcal{L}}(\vec{w}) - J_{\mathcal{S}}(\vec{w})$ . When the iterate  $\vec{w}^k$  is far away from  $\vec{w}^*$ , Newton's method may not converge and other techniques must initially be applied. We

| operator    | SBP21    |      | SBP42    |      |
|-------------|----------|------|----------|------|
| N           | $\ e\ $  | $r$  | $\ e\ $  | $r$  |
| 21          | 4.13e-02 | –    | 1.90e-02 | –    |
| 41          | 9.73e-03 | 2.16 | 2.19e-03 | 3.23 |
| 61          | 4.17e-03 | 2.13 | 6.34e-04 | 3.12 |
| 81          | 2.28e-03 | 2.12 | 2.70e-04 | 3.01 |
| Theoretical | 2        |      | 3        |      |

Table 2

Error and convergence rate.

371 choose the SOR method [15] until  $\|F(\vec{w}^k)\|_\infty$  is sufficiently small. For SOR, the next iterate  
 372 is given by  $\vec{w}^{k+1} = \vec{w}^k(1 - \alpha) + (\vec{w}^k - \mathbf{h}^k)\alpha$ , where  $\mathbf{h}^k$  is the Newton step from Algorithm 5.1  
 373 and  $\alpha \in (0, 1]$ .

374 To verify our procedure, we choose the steady manufactured solution to be [9]

$$\begin{aligned}
 (6.5) \quad u &= 1 - e^{\lambda x} \cos(2\pi y) & v &= \frac{1}{2\pi} \lambda e^{\lambda x} \sin(2\pi y) \\
 p &= \frac{1}{2} \left(1 - e^{2\lambda x}\right) & \lambda &= \frac{1}{2\epsilon} - \sqrt{\frac{1}{4\epsilon^2} + 4\pi^2}
 \end{aligned}$$

376 and the computational domain is changed to  $\Omega = [-0.5, 1] \times [-1, 1]$  for  $\epsilon = 1/20$ . Inserting  
 377 (6.5) into the time-independent version of (2.1) leads to  $k(t, x, y) = 0$ . The initial guess is  
 378  $\vec{w}^0 = (1, 1, \dots, 1)^\top$  and the tolerance  $tol$  in Algorithm 5.1 is again set to  $10^{-12}$ . Table 3  
 379 shows the error and convergence rates, which again agrees well with theory. In Figure 2,  
 380 the streamlines and the velocity field is illustrated for the converged solution on the grid  
 381 containing  $100 \times 100$  points. They agree well with previous results [9].

382 **6.2. The convergence rate of the Newton iteration.** Next, we will test the main develop-  
 383 ment in this paper. For  $\vec{w}^k$  sufficiently close to  $\vec{w}^*$ , Newton's method converges quadratically  
 384 in any norm [15], which means that  $e_{k+1} = C e_k^2$ , where  $C$  varies marginally between iterations  
 385 and  $e_k = \|\vec{w}^k - \vec{w}^*\|$ . To verify that, we consider a grid of size  $100 \times 100$  with the SBP42  
 386 operator. The exact solution  $\vec{w}^*$  is approximated by the last iterate. By the assumption that  
 387  $C$  is constant, the relation

$$\frac{e_{k+1}}{e_k} \approx \left( \frac{e_k}{e_{k-1}} \right)^p$$

389 is obtained for a general convergence rate  $p$ , which yields

$$p \approx \frac{\log(e_{k+1}/e_k)}{\log(e_k/e_{k-1})}.$$

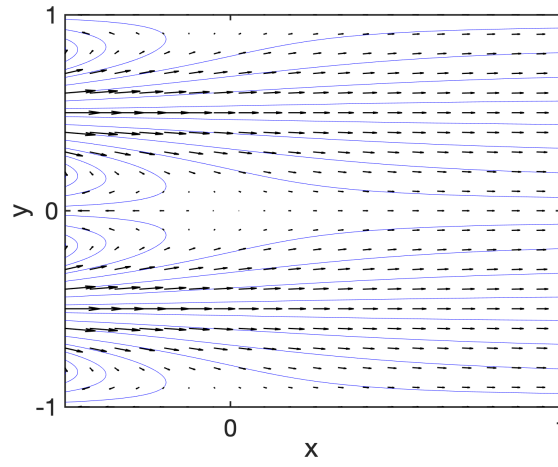
391 The error  $e_k = \|\vec{w}^k - \vec{w}^*\|_\infty$  is presented in Table 4 together with the estimations of  $p$ .  
 392 The convergence rate agrees well with the expected theoretical one, which verifies that the  
 393 Jacobian of  $\mathcal{F}$  is correct.



| operator    | SBP21    |      | SBP42    |      |
|-------------|----------|------|----------|------|
| N           | $\ e\ $  | $r$  | $\ e\ $  | $r$  |
| 21          | 2.04e-01 | –    | 4.95e-02 | –    |
| 41          | 4.56e-02 | 2.16 | 6.86e-03 | 2.85 |
| 61          | 2.04e-02 | 1.98 | 2.20e-03 | 2.80 |
| 81          | 1.16e-02 | 1.97 | 9.76e-04 | 2.83 |
| 101         | 7.46e-03 | 1.97 | 5.16e-04 | 2.85 |
| Theoretical | 2        |      | 3        |      |

**Table 3**

Error and (accuracy) convergence rate of (6.5).

**Figure 2.** Streamlines and the velocity field of (6.5).**Table 4**

Errors and the estimated (iterative) convergence rates of (6.5).

| k           | $\ e_k\ _\infty$ | $p$   |
|-------------|------------------|-------|
| 1           | 3.56e+00         | –     |
| 2           | 1.85e+01         | –     |
| 3           | 1.89e+00         | -1.38 |
| 4           | 1.21e+00         | 0.20  |
| 5           | 5.53e-01         | 1.74  |
| 6           | 1.10e-01         | 2.07  |
| 7           | 3.21e-03         | 2.19  |
| 8           | 3.31e-06         | 1.94  |
| 9           | 4.14e-12         | 1.98  |
| Theoretical | 2                |       |

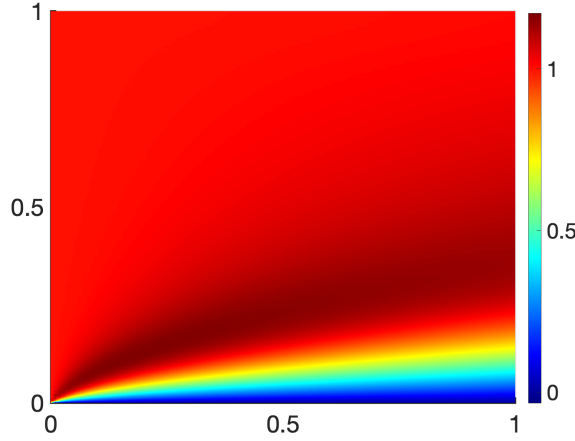


Figure 3. Flow over a solid surface.

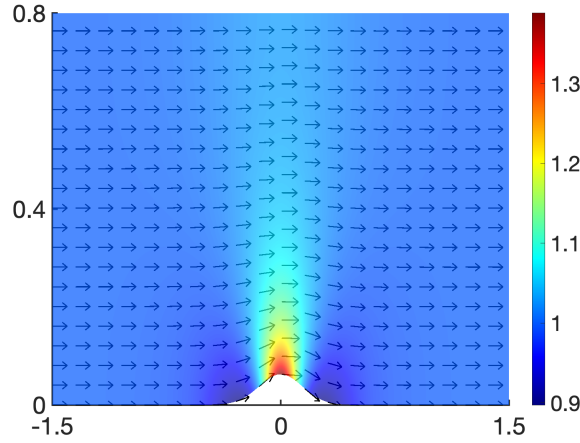
Table 5

Errors and the estimated (iterative) convergence rates for the flow over a solid surface.

| k           | $\ e_k\ _\infty$ | $p$  |
|-------------|------------------|------|
| 1           | 1.68e+00         | –    |
| 2           | 6.17e-01         | –    |
| 3           | 1.15e-01         | 1.67 |
| 4           | 4.37e-03         | 1.95 |
| 5           | 1.02e-05         | 1.85 |
| 6           | 5.51e-11         | 2.00 |
| Theoretical |                  | 2    |

Next, we move on to a more realistic case where the boundary data is set to  $g_1 = 1$ ,  $g_2 = g_3 = g_4 = g_5 = g_6 = 0$  and  $\epsilon = 0.01$ , which will lead to a boundary layer. The computations are performed on  $\Omega = [0, 1]^2$  with  $200 \times 200$  grid points with the SBP42 operator. Figure 3 illustrates  $\mathbf{u}$  for the converged solution and the iterative convergence order,  $p$ , is presented in Table 5. The estimated iterative convergence order agrees well with what is theoretically expected.

In the last experiment, we consider a curved grid [1] for the incompressible Euler equation (i.e.  $\epsilon = 0$ ). Both the south and north sides are solid surfaces, where the normal velocity is zero. The west side is an inflow boundary where  $u = 1$  and  $v = 0$  are specified and at the east side,  $p = 0$  is imposed. We change the domain to  $\Omega = [-1.5, 1.5] \times [0, 0.8]$  and include a smooth bump at the south boundary given by  $y(x) = 0.0625e^{-25x^2}$  [4]. In Figure 4, the converged solution is illustrated and the estimated iterative convergence rate  $p$  is presented in Table 6 for the initial guess  $(\mathbf{u}^0; \mathbf{v}^0; p^0) = (1, \dots, 1; 0, \dots, 0; 1, \dots, 1)$ . Again, the results agree well with the theoretical value.



**Figure 4.** Flow over a smooth bump. The plot illustrates the velocity field (arrows) and  $u$  (color figure) at the converged solution.

**Table 6**

Errors and the estimated (iterative) convergence rates for the bump.

| k           | $\ e_k\ _\infty$ | $p$  |
|-------------|------------------|------|
| 1           | 3.56e+00         | –    |
| 2           | 4.91e-01         | –    |
| 3           | 1.74e-02         | 1.69 |
| 4           | 3.33e-05         | 1.87 |
| 5           | 1.55e-10         | 1.96 |
| Theoretical |                  | 2    |

**7. Conclusions.** We derived an explicit expression for the Jacobian of a finite-difference discretization of the incompressible Navier-Stokes equations. Both the Jacobian of the system of equations and the Jacobian of the related boundary condition was computed exactly. By using the block-structure of the discretization, we showed that the Jacobian had a block structure as well, which lead to a compact and clean expression. We also showed that large parts of the Jacobian were computed by evaluating the discretization. We showed that the Jacobian could be used both in steady-state and time-dependent simulations. The numerical discretization was verified by manufactured solutions and the spatial convergence rates agreed well with the theoretical expectations. Furthermore, the computed estimates of the iterative convergence rates for Newton's method was two, and verified that the Jacobian was correctly computed. The methodology used in this paper is general and can be used in a straightforward manner for any numerical discretization of initial boundary value problems that can be written in matrix-vector form.

## REFERENCES

- [1] O. ÅLUND AND J. NORDSTRÖM, *Encapsulated high order difference operators on curvilinear non-conforming grids*, Journal of Computational Physics, (2019), pp. 209–224.
- [2] P. N. BROWN AND Y. SAAD, *Hybrid krylov methods for nonlinear systems of equations*, SIAM Journal on Scientific and Statistical Computing, 11 (1990), pp. 450–481.
- [3] M. H. CARPENTER, D. GOTTLIEB, AND S. ABARBANEL, *Time-stable boundary conditions for finite-difference schemes solving hyperbolic systems: methodology and application to high-order compact schemes*, Journal of Computational Physics, 111 (1994), pp. 220–236.
- [4] CENAERO, *VI2 Smooth Gaussian bump*, 2020 (accessed December 9, 2020). <https://how5.cenaero.be/content/vi2-smooth-gaussian-bump>.
- [5] J. CHAN AND C. TAYLOR, *Explicit Jacobian matrix formulas for entropy stable summation-by-parts schemes*, arXiv preprint arXiv:2006.07504, (2020).
- [6] D. C. D. R. FERNÁNDEZ, J. E. HICKEN, AND D. W. ZINGG, *Review of summation-by-parts operators with simultaneous approximation terms for the numerical solution of partial differential equations*, Computers & Fluids, 95 (2014), pp. 171–196.
- [7] A. JAMESON, *Time dependent calculations using multigrid, with applications to unsteady flows past airfoils and wings*, in 10th Computational Fluid Dynamics Conference, 1991, p. 1596.
- [8] D. A. KNOLL AND D. E. KEYES, *Jacobian-free Newton–Krylov methods: a survey of approaches and applications*, Journal of Computational Physics, 193 (2004), pp. 357–397.
- [9] L. KOVASZNY, *Laminar flow behind a two-dimensional grid*, Mathematical Proceedings of the Cambridge Philosophical Society, 44 (1948), pp. 58–62.
- [10] J. NOCEDAL AND S. WRIGHT, *Numerical optimization*, Springer Science & Business Media, 2006.
- [11] J. NORDSTRÖM, *A roadmap to well posed and stable problems in computational physics*, Journal of Scientific Computing, 71 (2017), pp. 365–385.
- [12] J. NORDSTRÖM AND C. LA COGNATA, *Energy stable boundary conditions for the nonlinear incompressible Navier–Stokes equations*, Mathematics of Computation, 88 (2019), pp. 665–690.
- [13] J. NORDSTRÖM AND A. A. RUGGIU, *Dual time-stepping using second derivatives*, Journal of Scientific Computing, 81 (2019), pp. 1050–1071.
- [14] T. C. PAPANASTASIOU, N. MALAMATARIS, AND K. ELLWOOD, *A new outflow boundary condition*, International journal for numerical methods in fluids, 14 (1992), pp. 587–608.
- [15] A. QUARTERONI, R. SACCO, AND F. SALERI, *Numerical mathematics*, vol. 37, Springer Science & Business Media, 2010.
- [16] P. J. ROACHE, *Code verification by the method of manufactured solutions*, J. Fluids Eng., 124 (2002), pp. 4–10.
- [17] M. SVÄRD AND J. NORDSTRÖM, *On the order of accuracy for difference approximations of initial-boundary value problems*, Journal of Computational Physics, 218 (2006), pp. 333–352.
- [18] M. SVÄRD AND J. NORDSTRÖM, *Review of summation-by-parts schemes for initial-boundary-value problems*, Journal of Computational Physics, 268 (2014), pp. 17–38.
- [19] M. SVÄRD AND J. NORDSTRÖM, *On the convergence rates of energy-stable finite-difference schemes*, Journal of Computational Physics, 397 (2019), p. 108819.