Nektar++: Design and implementation of an implicit, spectral/hp element, compressible flow solver using a Jacobian-free Newton Krylov approach

Zhen-Guo Yan<sup>a,b,\*</sup>, Yu Pan<sup>b</sup>, Giacomo Castiglioni<sup>b</sup>, Koen Hillewaert<sup>c</sup>, Joaquim Peiró<sup>b</sup>, David Moxey<sup>d</sup>, Spencer J. Sherwin<sup>b</sup>

<sup>a</sup>State Key Laboratory of Aerodynamics, China Aerodynamics Research and Development Center, Mianyang, PRC

<sup>b</sup>Department of Aeronautics, Imperial College London, London, U.K.
<sup>c</sup>Université de Liège, Place du 20-Août, 7, B-4000 Liège, Belgique
<sup>d</sup>Department of Mechanical Engineering, University of Exeter, Exeter, U.K.

# Abstract

At high Reynolds numbers the use of explicit in time compressible flow simulations with spectral/hp element discretization can become significantly limited by time step. To alleviate this limitation we extend the capability of the spectral/hp element open-source software framework, Nektar++, to include an implicit discontinuous Galerkin compressible flow solver. The integration in time is carried out by a singly diagonally implicit Runge-Kutta method. The non-linear system arising from the implicit time integration is iteratively solved by the Jacobian-free Newton Krylov (JFNK) method. A favorable feature of the JFNK approach is its extensive use of the explicit operators available from the previous explicit in time implementation. The functionalities of different building blocks of the implicit solver are analyzed from the point of view of software design and placed in appropriate hierarchical levels in the C++ libraries. In the detailed implementation, the contributions of different parts of the solver to computational cost, memory consumption and programming complexity are also analyzed. A combination of analytical and numerical methods is adopted to simplify the programming complexity in forming the preconditioning matrix. The solver is verified and tested using cases such as manufactured compressible Poiseuille flow, Taylor-Green vor-

Email address: yanzhg@mail.ustc.edu.cn (Zhen-Guo Yan)

<sup>\*</sup>Corresponding author

tex, turbulent flow over a circular cylinder at Re = 3900 and shock wave boundary-layer interaction. The results show that the implicit solver can speed-up the simulations while maintaining good simulation accuracy.

Keywords: Nektar++, implicit time integration, Spectral/hp element, discontinuous Galerkin, Jacobian-free Newton Krylov.

#### 1. Introduction

Nektar++ is a C++ based cross-platform open-source framework with the purpose of making high-order spectral/hp element methods accessible to a wider range of researchers and engineers [6, 26]. Similar to many packages, the Nektar++ framework is composed of libraries and solvers. The six core libraries provide basic functions for different aspects of the implementation of the spectral/hp element method, including the elemental base functions, the extraction of geometric information, the numerical schemes for specific types of equations. All these libraries provide the basic building blocks for the numerical methods, which makes the construction of new solvers easier and more transparent. The Nektar++ framework supports various kinds of high-order curved meshes including hexahedral, prismatic, pyramidal, tetrahedral, quadrilateral and triangular meshes. Continuous Galerkin (CG) [20], discontinuous Galerkin (DG) [24] and flux reconstruction (FR) schemes (currently only support hexahedral and quadrilateral meshes) [25] are supported for spatial discretizations. Up to 12 built-in solvers have been developed to date providing the capability of multi-solver coupling [26]. These solvers make the Nektar++ framework applicable to a wide range of simulations |6, 26|.

In Nektar++, the development of the implicit solver for the compressible Navier-Stokes (NS) equations has been based on the explicit version of the solver. The explicit solver has potentially significant limitations on the time step because of the stability restrictions of high-order spectral/hp element schemes, see for instance reference [18]. These time step restrictions can significantly reduce the convective and diffusive time step near the boundaries, e.g. to resolve boundary layers with stretched cells at increasing Reynolds numbers [20] or near badly shaped cells in complex meshes. The unconditional stability of implicit solvers can not only speed-up the simulations by allowing much larger time steps but also reduce the risk of run time "blow-up" resulting from local instability in the nonlinear evolution of the flow

field.

However, the development of the implicit solver presents some challenges both in numerical algorithms and in software design. Firstly, the coupling of different equations in the compressible NS equations means the equations cannot be solved component by component as is the case for existing implicit solvers of decoupled equation systems in Nektar++. This not only makes the coupled implicit system much larger in size but also makes the implicit system stiffer due to the scattering of eigenvalues for different equations. Secondly, the problem naturally leads to a nonlinear system requiring highly optimized set up operators since the implicit system needs to be updated along with the simulations.

The development of implicit compressible flow solvers in high-order opensource software is still quite limited especially for unsteady simulations. Based on deal. II, Hartmann and Houston [16] developed a solver of the compressible NS equations using implicit time integration and the DG scheme. However, their solver is only tested on steady-state problems and only supports quadrilateral and hexahedral meshes. A fully implicit FR solver for compressible NS equations has been developed in the CoolFluid framework [39]. However, the implicit solver is only tested in steady state problems and only first and second order curved quadrilateral and hexahedral meshes are supported. Nek5000 provides an implicit solver based on Jacobian-free Newton Krylov method [9]. However, it only supports weakly compressible simulations through some modifications of the incompressible NS equations. An implicit solver based on MOOSE has been reported recently, but it currently only focuses on incompressible flow systems [33]. In summary, only a few implicit solvers with limited capabilities are available in the high-order open-source community.

We will discuss the development of an implicit solver based on the Nektar++ framework. The progress and capabilities of the solver are summarized and demonstrated using a few test cases. The long term goal is to improve the efficiency and robustness of the compressible flow simulations in Nektar++ therefore enabling large scale high fidelity simulations of unsteady problems. Meanwhile, it provides a good example for developing user defined solvers and offers a new pattern for implicit solvers of coupled nonlinear systems in Nektar++. The implicit solver also offers a good alternative for exploring different aspects of the implicit solvers and for large scale simulations in the open-source community.

The governing equations are presented in Section 2. Spatial and temporal

numerical schemes are described in Section 3. Section 4 gives a detailed description of the software implementation and a summary of the capabilities of the implicit solver. Results of various test cases are presented in Section 5 to verify the solver implementation and to illustrate its capabilities. Finally, conclusions are drawn in Section 6.

## 2. Governing equations

This section presents the Navier-Stokes governing equations of compressible viscous flow and discuss their non-dimensionalization.

## 2.1. Compressible Navier-Stokes equations

The three dimensional (3D) non-dimensional compressible NS equations can be expressed in abridged form as

$$\frac{\partial \mathbf{U}}{\partial t} = -\nabla \cdot \mathbf{H} = -\nabla \cdot (\mathbf{F} - \mathbf{G}). \tag{1}$$

Here, the summation convention is used for repeated indexes, the conservative variable vector  $\mathbf{U}$ , the advection flux vector  $\mathbf{F} = (\mathbf{F}_1, \mathbf{F}_2, \mathbf{F}_3)^T$  and the diffusion flux vector  $\mathbf{G} = (\mathbf{G}_1, \mathbf{G}_2, \mathbf{G}_3)^T$  are

$$\mathbf{U} = \begin{pmatrix} \rho \\ \rho u_1 \\ \rho u_2 \\ \rho u_3 \\ E \end{pmatrix}, \mathbf{F}_i = \begin{pmatrix} \rho u_i \\ \rho u_i u_1 + p \delta_{1i} \\ \rho u_i u_2 + p \delta_{2i} \\ \rho u_i u_3 + p \delta_{3i} \\ u_i (E + p) \end{pmatrix}, \mathbf{G}_i = \begin{pmatrix} 0 \\ \tau_{i1} \\ \tau_{i2} \\ \tau_{i3} \\ u_i \tau_{ij} - q_i \end{pmatrix}, \quad (2)$$

where  $\delta_{ij}$  is the Kronecker delta,  $\rho$  is the density, p is the pressure, T is the temperature,  $u_i$  is the ith velocity component,  $E = p/(\gamma - 1) + \rho u_k u_k/2$  is the total specific energy,  $\gamma = C_p/C_v = 1.4$  is the ratio of the specific heats with constant pressure  $(C_p)$  and constant volume  $(C_v)$  and the stress tensor is

$$\tau_{ij} = \mu \left( \frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right) - \frac{2}{3} \mu \frac{\partial u_k}{\partial x_k},\tag{3}$$

where  $x_i$  is the *i*th coordinate and  $\mu$  is the dynamic viscosity. Finally, the *i*th component of heat flux is given by

$$q_i = -\frac{\mu}{\operatorname{Ma}_{\infty}^2 \Pr(\gamma - 1)} \frac{\partial T}{\partial x_i},\tag{4}$$

where  $Ma_{\infty}$  is the Mach number based on the reference flow state, Pr is the Prandtl number. The Euler equations of inviscid compressible flow are recovered by ignoring the diffusion flux vector. Together with appropriate initial and boundary conditions, Eq. (1) is closed through the equation of state,

$$p = \frac{\rho T}{\gamma \text{Ma}_{\infty}^2},\tag{5}$$

and Sutherland's law,

$$\mu = \frac{T^{3/2}}{\text{Re}_{\infty}} \frac{1 + 110/T_{\infty}}{T + 110/T_{\infty}},\tag{6}$$

are provided. In the above equations,  $\text{Re}_{\infty}$  is the Reynolds number based on the reference flow state and  $T_{\infty}$  is the reference temperature which should be given in degrees Kelvin.

#### 3. Numerical methods

This section describes the discontinuous Galerkin method, the specific spectral/hp element method adopted, and the implicit temporal discretizations.

#### 3.1. Discontinuous Galerkin method

In the discontinuous Galerkin method, the computational domain  $(\Omega)$  is divided into  $N_e$  non-overlapping elements  $(\Omega_e)$ . The space of test function is defined as

$$V^{P} = \left\{ \phi : \phi |_{\Omega_{e}} \in \mathscr{P}^{P} \left( \Omega_{e} \right), e = 1, \cdots, N_{e} \right\}, \tag{7}$$

where  $\mathscr{P}^P(\Omega_e)$  is the polynomial space of degree P in  $\Omega_e$ . The weak form of Eq. (1) is obtained by multiplying by the test function  $\phi_p$  and performing integration by parts in  $\Omega_e$ ,

$$\int_{\Omega} \frac{\partial \mathbf{U}}{\partial t} \phi_p d\Omega_e = \int_{\Omega} \nabla \phi_p \cdot \mathbf{H} d\Omega_e - \int_{\Gamma} \phi_p \mathbf{H}^n d\Gamma_e, \tag{8}$$

where  $\Gamma_e$  is the element boundary,  $\mathbf{H}^n = \mathbf{F}^n - \mathbf{G}^n$ ,  $\mathbf{F}^n = \mathbf{F} \cdot \mathbf{n}$  and  $\mathbf{G}^n = \mathbf{G} \cdot \mathbf{n}$  are boundary fluxes on the elemental outward normal direction  $(\mathbf{n})$ .

The vector of conservative variables  $\mathbf{U}$  is approximated in the same polynomial space as the test functions and it is expressed as

$$\mathbf{U}\left(\mathbf{x},t\right) \simeq \sum_{q=1}^{N} \mathbf{u}_{q}\left(t\right) \phi_{q}\left(\mathbf{x}\right), \tag{9}$$

where  $\mathbf{u}_q(t)$  is the qth coefficient of the base (or trial) function  $\phi_q(\mathbf{x})$  and N is the total degrees of freedom (DoFs) in the element. The flow variable values on some quadrature points  $\mathbf{Q}$  are calculated using the backward transformation  $\mathbf{Q}_i = \sum_{q=1}^N \mathbf{B}_{iq} \mathbf{u}_q$  where  $\mathbf{B}$  is the backward transformation matrix with  $\mathbf{B}_{iq} = \phi_q(\mathbf{x}_i)$ . Here  $\mathbf{x}_i$  is coordinates of the *i*th quadrature point. The fluxes are calculated at these quadrature points and a quadrature rule with  $N_Q$  quadrature points is adopted to calculate the integration in the element and  $N_Q^\Gamma$  quadrature points on element boundaries. This leads to

$$\sum_{i=1}^{N_Q} \sum_{q=1}^{N} \phi_p(\mathbf{x}_i) w_i J_i \phi_q(\mathbf{x}_i) \frac{d\mathbf{u}_q}{dt} = \sum_{i=1}^{N_Q} w_i J_i \nabla \phi_p(\mathbf{x}_i) \cdot \mathbf{H}(\mathbf{Q}_i)$$

$$- \sum_{i=1}^{N_Q^{\Gamma}} \phi_p(\mathbf{x}_i^{\Gamma}) w_i^{\Gamma} J_i^{\Gamma} \hat{\mathbf{H}}_i^n,$$
(10)

where  $w_i$  and  $J_i$  are the quadrature weights and grid metric Jacobian on the ith quadrature point,  $\hat{\mathbf{H}}_i^n = \hat{\mathbf{F}}_i^n - \hat{\mathbf{G}}_i^n$ ,  $\hat{\mathbf{F}}_i^n$  and  $\hat{\mathbf{G}}_i^n$  are the numerical normal fluxes on the ith quadrature point  $\mathbf{x}_i^{\Gamma}$  of the element boundaries,  $w_i^{\Gamma}$  and  $J_i^{\Gamma}$  are the quadrature weights and grid Jacobian of the lower dimensional integrations on element boundaries, which are usually not equal to the  $w_i$  and  $J_i$  even if the element and element boundary quadrature points are at the same position.

The weak DG scheme for the advection term is complete as long as a Riemann numerical flux is used to calculate the normal flux,  $\hat{\mathbf{F}}^n(\mathbf{Q}^+, \mathbf{Q}^-, \mathbf{n})$ , in which  $\mathbf{Q}^+$  and  $\mathbf{Q}^-$  are variable values on the exterior and interior sides of the element boundaries, respectively.

#### 3.1.1. Interior penalty method

The interior penalty (IP) method is adopted to discretize the diffusion term. To simplify the complexity of the expressions for the multi-dimensional matrix operations in this section, all the indexes of matrices and vectors follow the summation convention. The diffusion flux in the IP method is expressed in the following equivalent form

$$\mathbf{G}_{ik} = \mathbf{K} \left( \mathbf{Q} \right)_{ijkl} \nabla_j \mathbf{Q}_l, \tag{11}$$

where i and j are indexes of different spatial directions, k and l are indexes of different flow variables. The expression of  $\mathbf{K}(\mathbf{Q})_{ijkl}$  can be found in [16]. For

clarity, only the diffusion terms are shown in the derivation. The IP method can be expressed in the following primal form

$$\int_{\Omega} \frac{\partial \mathbf{U}_{k}}{\partial t} \phi_{p} d\Omega = \sum_{e=1}^{N_{e}} \int_{\Gamma_{e}} \llbracket \phi_{p} \rrbracket_{i} \left\{ \mathbf{K} \left( \left\{ \mathbf{Q} \right\} \right)_{ijkl} \nabla_{j} \mathbf{Q}_{l} \right\} d\Gamma_{e} - \sum_{e=1}^{N_{e}} \int_{\Omega_{e}} \nabla_{i} \phi_{p} \mathbf{G}_{ik} d\Omega_{e} 
+ \beta_{s} \sum_{e=1}^{N_{e}} \int_{\Gamma_{e}} \llbracket \mathbf{Q}_{l} \rrbracket_{i} \left\{ \mathbf{K} \left( \left\{ \mathbf{Q} \right\} \right)_{ijkl} \nabla_{j} \phi_{p} \right\} d\Gamma_{e} 
+ \sum_{e=1}^{N_{e}} \beta_{p} \int_{\Gamma_{e}} \llbracket \phi_{p} \rrbracket_{i} \mathbf{K} \left( \left\{ \mathbf{Q} \right\} \right)_{ijkl} \llbracket \mathbf{Q}_{l} \rrbracket_{j} d\Gamma_{e},$$
(12)

with  $\beta_s = 1$ ,  $\beta_p = (P+1)^2/h$  for simulations with quadrilateral and hexahedral meshes,  $[\![w]\!]_i = w^1\mathbf{n}_i^1 + w^2\mathbf{n}_i^2$  and  $\{w\} = (w^1 + w^2)/2$ , where the superscripts 1 and 2 indicate values from the two sides of the element interface. Compared with the IP method proposed in reference [17],  $\mathbf{K}(\mathbf{Q})_{ijkl}$  is replaced by  $\mathbf{K}(\{\mathbf{Q}\})_{ijkl}$  for implementation purposes.  $\beta_s = 1$  makes the IP method symmetric and adjoint consistent. Values of  $\beta_p$  for other mesh types can be found in reference [19].

Using the relation between primal form and its flux form in reference [1], Eq. (12) can be expressed in its flux form. As a result, the whole discretization can be written in the following matrix form

$$\mathbf{M} \frac{d\mathbf{u}}{dt} = \sum_{j=1}^{d} \mathbf{B}^{T} \mathbf{D}_{j}^{T} \mathbf{\Lambda} (wJ) \mathbf{H}_{j} (\mathbf{Q})$$
$$- \left( \mathbf{B}^{\Gamma} \mathbf{M}_{c} \right)^{T} \mathbf{\Lambda} \left( w^{\Gamma} J^{\Gamma} \right) \hat{\mathbf{H}}^{n} ,$$
$$- \sum_{j=1}^{d} \mathbf{B}^{T} \mathbf{D}_{j}^{T} \mathbf{J}^{T} \mathbf{\Lambda} \left( w^{\Gamma} J^{\Gamma} \right) \hat{\mathbf{S}}_{j}^{n}$$
(13)

where d is the spatial dimension,  $\mathbf{M} = \mathbf{B}^T \mathbf{\Lambda} (wJ) \mathbf{B}$  is the mass matrix,  $\mathbf{\Lambda}$  represents a diagonal matrix,  $\mathbf{D}_j$  is the derivative matrix in the jth direction,  $\mathbf{B}^{\Gamma}$  is the backward transformation matrix of  $\phi^{\Gamma}$  and  $\mathbf{M}_{\mathbf{c}}$  is the mapping matrix between  $\phi^{\Gamma}$  and  $\phi$ ,  $\mathbf{J}$  is the interpolation matrix from quadrature points of a element to quadrature points of its element boundaries,  $\hat{\mathbf{H}}^n = \hat{\mathbf{F}}^n - \hat{\mathbf{G}}^n$ ,  $\hat{\mathbf{F}}^n$  is the Riemann flux,  $\hat{\mathbf{G}}_k^n = \mathbf{n}_i \left( \left\{ \mathbf{K} \left( \left\{ \mathbf{Q} \right\} \right)_{ijkl} \nabla_j \mathbf{Q}_l \right\} + \beta_p \mathbf{K} \left( \left\{ \mathbf{Q} \right\} \right)_{ijkl} \left[ \mathbf{Q}_l \right]_j \right)$ ,  $\hat{\mathbf{S}}_{j,k}^n = \beta_s [\![\mathbf{Q}_l]\!]_i \left\{ \mathbf{K} \left( \left\{ \mathbf{Q} \right\} \right)_{ijkl} \right\}$ . The  $\hat{\mathbf{G}}_k^n$ , which represents the flux integration

of the first and last terms on the right-hand side of Eq. (12), can be calculated together with the advection flux, while the symmetric flux  $\hat{\mathbf{S}}_{j,k}^n$  cannot. Some approximations of the preconditioning matrices are make based on these properties of  $\hat{\mathbf{G}}_k^n$  and  $\hat{\mathbf{S}}_{j,k}^n$ , as discussed in Section 4.2.2.

Both the IP and the local DG (LDG) method are available in Nektar++ for the compressible flow solver. Currently the preferred implicit solver is based on the IP method since it is easier to derive its Jacobian matrix. Moreover, the IP method leads to smaller stencil than the LDG in general, which is beneficial not only for lowering the memory consumption but also for simplifying the calculation of the Jacobian matrix needed when preconditioning (see Section 4.2.2).

## 3.1.2. Shock-capturing method

An artificial viscosity method is implemented to regularize discontinuous solutions. This effectively amounts to adding artificial viscosity,  $\mu_{av}$ , and artificial thermal conductivity,  $\kappa_{av}$ , terms to the physical  $\mu$  and  $\kappa$  terms. The expressions of the artificial  $\mu_{av}$  and  $\kappa_{av}$  terms are

$$\mu_{av} = \mu_0 \rho \frac{h}{P} \left( c + \sqrt{u_k u_k} \right) S, \tag{14}$$

$$\kappa_{av} = \mu_{av} \frac{C_p}{\Pr},\tag{15}$$

where h is the mesh size,  $\mu_0 = 0.25$  is a parameter that controls the magnitude of  $\mu_{av}$  and S is the shock sensor. The modal resolution-based shock sensor proposed in reference [32] and Ducros' sensor [12] are implemented in Nektar++.

#### 3.2. Time discretization

After the spatial discretization, the partial differential equations Eq. (8) become a set of ordinary differential equations (ODEs) of the form

$$\mathbf{M}\frac{d\mathbf{u}}{dt} = \mathcal{L}\left(\mathbf{u}, t\right),\tag{16}$$

where  $\mathcal{L}(\mathbf{u},t)$  is the discrete term representing the spatial discretization operator. Various time integration methods are implemented to discretize the temporal derivatives which are summarized in Section 4.3. Here we describe

the explicit multi-stage Runge-Kutta (ERK) and the singly diagonally implicit multi-stage Runge-Kutta (SDIRK) methods [21] as illustrative examples of implementation. The discrete equations of the Runge-Kutta methods can be expressed in general as

$$\hat{\mathbf{u}}^m = a^{mm} \mathbf{E}^m + \sum_{i=0}^{m-1} a^{mi} \mathbf{E}^i + \mathbf{u}^n, \ m = 0, 2, \cdots, M - 1;$$
 (17)

$$\mathbf{E}^{m} = \Delta t \mathbf{M}^{-1} \mathcal{L} \left( \hat{\mathbf{u}}^{m}, t^{n} + c^{m} \Delta t \right); c^{m} = \sum_{i=0}^{m} a^{mi}, \tag{18}$$

$$\mathbf{u}^{n+1} = \sum_{i=0}^{M-1} b^i \mathbf{E}^i + \mathbf{u}^n, \tag{19}$$

where  $\mathbf{u}^n$  is the flow solution vector evaluated at the *n*th time step, M is the total number of stages of the Runge-Kutta method,  $a^{mi}$  denotes the coefficients of the Runge-Kutta method, and  $\hat{\mathbf{u}}^m$  is the *m*th stage approximation of the solution. The coefficients of the Runge-Kutta schemes are given in Appendix A using Butcher tableaus.

For explicit methods  $(a^{mm}=0)$ , the calculation of Eq. (17) is trivial. However, Eq. (17) becomes a nonlinear system for implicit stages  $(a^{mm} \neq 0)$  that can be expressed as

$$\mathbf{N}\left(\hat{\mathbf{u}}^{m}\right) = \hat{\mathbf{u}}^{m} - \mathbf{S}_{m} - \triangle t a^{mm} \mathbf{M}^{-1} \mathcal{L}\left(\hat{\mathbf{u}}^{m}, t^{n} + c^{m} \triangle t\right) = \mathbf{0} \quad , \tag{20}$$

where  $\mathbf{S}_m = \sum_{i=0}^{m-1} a^{mi} \mathbf{E}^i + \mathbf{u}^n$  is a source term based on known approximations at stage m.

The Newton method [22] is adopted for solving the nonlinear system (20). The Newton iteration can be written as

$$\bar{\mathbf{u}}^0 = \mathbf{S}_m, \tag{21}$$

$$\frac{\partial \mathbf{N}}{\partial \mathbf{u}} \left( \bar{\mathbf{u}}^l \right) \left( \bar{\mathbf{u}}^{l+1} - \bar{\mathbf{u}}^l \right) = -\mathbf{N} \left( \bar{\mathbf{u}}^l \right), \ l = 0, 1, \cdots.$$
 (22)

When the  $L_2$  norm of the residual is sufficiently small, namely

$$\left\| \mathbf{N} \left( \bar{\mathbf{u}}^{l} \right) \right\|_{2} < \alpha \left\| \mathbf{N} \left( \bar{\mathbf{u}}^{0} \right) \right\|_{2}, \tag{23}$$

 $\hat{\mathbf{u}}^m = \bar{\mathbf{u}}^l$  is regarded as the converged solution. A reduction of the initial residual with  $\alpha = 10^{-3}$  will be enforced unless otherwise specified. In some

cases, a smaller  $\alpha$  is needed to maintain temporal accuracy. As discussed in Section 6.4, the solver will be slower using a smaller  $\alpha$ , but the change in efficiency will not be very large.

The restarted generalized minimal residual method (GMRES) [35] is used for solving the linear problem (22), which is restarted when the number of GMRES iterations exceeds 30. The GMRES iteration is terminated after the residual drops by  $5 \times 10^{-2}$ , which is chosen mainly based on efficiency considerations. An important aspect in GMRES is preconditioning, which will be discussed separately. Currently, a standard version of GMRES proposed in [35] is used. Other optimized linear system solvers, such as the SNESANDERSON and SNESNGMRES [5] in PETSc, can also be adopted to further improve the performance of the linear system solver.

The Jacobian matrix,  $\partial \mathbf{N}/\partial \mathbf{u}$ , in Eq. (22) is a large sparse matrix. The calculation and storage of  $\partial \mathbf{N}/\partial \mathbf{u}$  is usually excessively expensive compared with those of the explicit solver. Since  $\partial \mathbf{N}/\partial \mathbf{u}$  is only used to calculate its inner product with a vector  $(\partial \mathbf{N}/\partial \mathbf{u} \cdot \mathbf{q})$  in the GMRES, a Jacobian-free method is adopted to avoid explicitly calculating and storing the Jacobian matrix (see Section 4.2.2).

The use of good preconditioners in GMRES is very important for efficiently solving stiff linear systems. Instead of solving the system of Eq. (22) directly, one can get the same solution by solving the preconditioned linear system

$$\left(\frac{\partial \mathbf{N}}{\partial \mathbf{u}} \mathbf{P}^{-1}\right) \left(\mathbf{P} \triangle \bar{\mathbf{u}}^l\right) = -\mathbf{N} \left(\bar{\mathbf{u}}^l\right), \tag{24}$$

where  $\mathbf{P}$  is the preconditioning matrix. A good preconditioner should be able to effectively cluster the distribution of eigenvalues of the linear system. The preconditioner  $\mathbf{P}$  is usually an approximate matrix of the Jacobian matrix  $\partial \mathbf{N}/\partial \mathbf{u}$ . It should be relatively easy to invert since it has to be used repeatedly in every GMRES iteration. Usually it is not necessary to calculate  $\partial \mathbf{N}/\partial \mathbf{u} \, \mathbf{P}^{-1}$  explicitly. Only the inner product of the matrix  $\mathbf{P}^{-1}$  with some vectors is required in a practical implementation.

# 4. Implementation

This section presents the details of the software implementation of the implicit solver. The structure of the Nektar++ libraries and the general procedure followed in constructing a flow solver are introduced in Section 4.1.

Section 4.2 describes the software design of the implicit solver in detail. Finally Section 4.3 summarizes the capabilities of the solver and the techniques for code verification.

# 4.1. Introduction to Nektar++

Nektar++ mainly consists of the libraries and several built-in solvers. The Nektar++ libraries provide a structured hierarchy of C++ classes, which offers the major functions needed for spectral/hp element methods. The six core libraries are:

- LibUtilities: polynomial base functions, basic data classes such as array and matrix, linear algebra functions, time integration schemes, parallel communication and I/O.
- StdRegions: standard (or reference) elements, polynomial approximation of the solution using base functions, operations on standard elements such as derivation and integration operations.
- SpatialDomains: mapping from standard elements to physical elements, grid metrics and grid Jacobian.
- LocalRegions: extension of operations on physical elements, physical elements inheriting from StdRegions and SpatialDomains.
- MultiRegions: collection of physical elements comprising  $\Omega$ , adjacent relations between elements, elemental and global variable mapping.
- SolverUtils: high level building blocks of solvers, drivers controlling simulation procedures, general equation systems, specific numerical scheme for different types of equations.

In the latest version of Nektar++, four additional libraries have been included to enrich the libraries with functionalities encompassing efficient operator evaluation, quasi-3D simulations, high-order curved mesh generation and post-processing [26]. The Nektar++ libraries provide main functions related to data storage (such as matrix and array), linear/nonlinear algebra (such as GMRES), geometry related calculations and numerical methods (such as DG), but nothing related to specific equations.

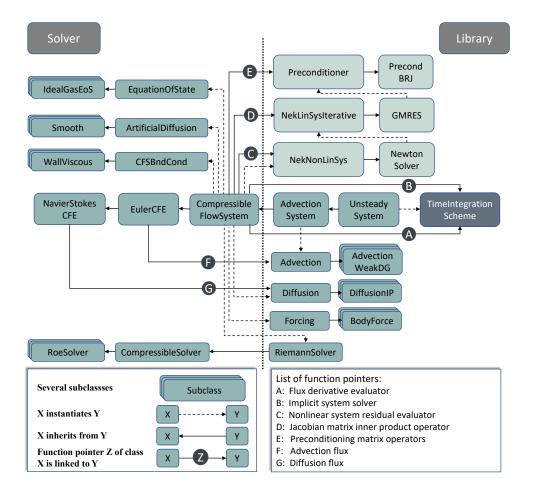


Figure 1: Class structure of the explicit and implicit solvers. The equation system classes (EulerCFE or NavierStokesCFE) contain access to the main functionalities of the solver, such as time integration, and the solution fields. They make use of numerical methods from the libraries, such as AdvectionWeakDG, and equations system related functions, such as the advection flux (F) and diffusion flux (G), to form the spatial discretization operator A. The spatial operator A and the time integration method form the explicit solver using the method of lines. For the implicit solver, additional classes like the Newton solver (NewtonSolver), GMRES solver (GMRES) and linear algebra solver of preconditioners (e.g. PrecondBRJ) are instantiated. Together with operators related to the nonlinear system (C), the Jacobian matrix (D) and preconditioning matrix (E), an implicit system solver (B) is constructed, which is linked to the implicit time integration scheme to form the implicit solver.

The main procedures for building a solver are briefly introduced using the explicit solver of the NS equations as an example. Besides the main solver class (CompressibleFlowSolver) which controls the solving procedure, an equation system class is needed. The equation system classes (EulerCFE or NavierStokesCFE) are instantiated and initialized dynamically based on user inputs using a factory method pattern [6], which is also extensively used for the dynamic object creation of classes in the solver. The equation system class inherits instantiations of classes related to geometry information, solution approximation, time integration and others from equations system classes in the libraries such as UnsteadySystem in SolverUtils. Thus the main functionality of the equation system class is to offer equation system related functions and to form spatial discretization operators using numerical methods from the libraries such as, for instance, AdvectionWeakDG. Fig. 1 illustrates the class structure of the explicit and implicit solvers. The equation of state (EquationOfState), boundary conditions (CFSBndCond), Riemann flux (RiemmanSolver), shock capturing method (ArtificialDiffusion), forcing term (Forcing), advection flux (function pointer F) and diffusion flux (function pointer G) are instantiated or implemented in the equation system class. Specific numerical schemes like the weak DG scheme for advection terms (AdvectionWeakDG) and the interior penalty method for diffusion terms (DiffusionIP) are instantiated to calculate  $\nabla \cdot \mathbf{F}$  and  $\nabla \cdot \mathbf{G}$  in Eq. (1) provided with all these equation system related functions (F and G). Finally the spatial discretization operator  $\mathcal{L}(\mathbf{u},t)$  is linked to the time integration class using the pointer-to-function A and the explicit solver is complete.

## 4.2. Implicit solver architecture and implementation

Section 4.2.1 presents the analysis of different building blocks of the implicit solver and their hierarchical structures. The implementation of the operations related to the NS equations is then discussed in Section 4.2.2.

#### 4.2.1. Analysis of the implicit solver architecture

We first analyze the different building blocks of the solver from the point of view of library and software design and place these blocks into the right hierarchical levels in the libraries and the solver. The Newton method, GM-RES method and linear algebra solvers of preconditioners, which are general mathematical methods independent of the type of equations, are placed in the LibUtilities as shown on the right upper corner of Fig. 1. Meanwhile operations that depend on the equations and the numerical schemes, i.e. the

implicit system solver (function pointer B in Fig. 1), nonlinear system residual evaluator (function pointer C in Fig. 1), Jacobian matrix operation (function pointer D in Fig. 1) and preconditioning matrix operations (function pointer E in Fig. 1), are coded at the solver level in the CompressibleFlowSystem class and they are linked to the corresponding classes from the libraries using pointers to these functions.

This structure provides a general framework for developing implicit solvers for different equation systems. Provided with equation system related functions, an implicit solver can also be constructed for other equations systems using the classes programmed in the libraries.

## 4.2.2. Implementation of the NS equations related operations

The top level class structure of the implicit solver is illustrated in Fig. 1. Even though only a nonlinear system solver is added to the Nektar++ libraries, the nonlinearity and coupling of the NS equations makes the design of the Jacobian matrix and preconditioning matrix related operations much more difficult than that of the existing implicit solvers of linear decoupled equation systems for the following reasons: a) the nonlinearity means the Jacobian and preconditioning matrices should be updated along with the simulation, which makes the computational speed of these operations vital for the overall efficiency; b) The coupling of eigenvalues of different equations increases the stiffness of the system; c) The coupling of different equations leads to a much larger Jacobian matrix, especially in large-scale 3D simulations. We will discuss the design of equation related operations in the following.

Nonlinear system residual. The nonlinear system residual,  $\mathbf{N}(\mathbf{u})$ , is easily available from an existing explicit solver. However, for the sake of efficiency it is important to optimize its evaluation since the implicit solver requires repeated calculations of the residual. Techniques such as sum-factorization [7, 28] and collections of similar linear algebra operations [26], have been studied and adopted in Nektar++ to speed-up the evaluation of  $\mathbf{N}(\mathbf{u})$ .

Jacobian matrix calculation. Explicitly calculating and storing the Jacobian matrix  $\partial \mathbf{N}/\partial \mathbf{u}$  accurately is expensive both in CPU cost and memory consumption, as discussed in Section 5. In the GMRES solver, we do not need to calculate the Jacobian matrix explicitly, but only need to evaluate the inner product of the Jacobian matrix with a vector  $(\partial \mathbf{N}/\partial \mathbf{u} \cdot \mathbf{q})$ , therefore

we adopt a Jacobian-free method which reads

$$\frac{\partial \mathbf{N}}{\partial \mathbf{u}}(\mathbf{u}) \cdot \mathbf{q} \simeq \frac{\mathbf{N}(\mathbf{u} + \epsilon \mathbf{q}) - \mathbf{N}(\mathbf{u})}{\epsilon},\tag{25}$$

$$\epsilon = \|\mathbf{u}\|_2 \sqrt{\hat{\epsilon}},\tag{26}$$

where  $\hat{\epsilon} = 2.22 \times 10^{-16}$  is the machine epsilon for our computation system. By storing  $\mathbf{N}(\mathbf{u})$ , only one evaluation of the nonlinear residual  $\mathbf{N}(\mathbf{u} + \epsilon \mathbf{q})$  is required for each  $\mathbf{q}$ . The Jacobian-free method offers a relatively efficient way of calculating the Jacobian matrix, while keeping enough accuracy to maintain high convergence speed [22]. The errors in numerically calculating the Jacobian matrix have been studied in references [38, 14]. Both references report similar convergence history using numerical Jacobian and analytical Jacobians. However, reference [45] finds that simulations with analytical Jacobian has better efficiency.

From the point of view of software design, this replaces a numerical method and a set of operations that are highly dependent of the type of equation solved by a more general finite difference scheme and an existing residual evaluator N(u). The Jacobian-free code will also be generally applicable to other equations systems if provided with the corresponding N(u).

Preconditioning matrix calculation. Constructing a good preconditioning matrix  $\mathbf{P}$  is highly dependent on appropriately approximating the Jacobian matrix  $\partial \mathbf{N}/\partial \mathbf{u}$ . Although some explorations on completely matrix-free preconditioners have been reported in the literature, the most efficient and most widely used preconditioners still require explicitly evaluating at least part of  $\partial \mathbf{N}/\partial \mathbf{u}$  for low to medium order simulations [15, 22]. Here a J step block relaxed Jacobi iterative preconditioner, BRJ(J), is developed instead of the widely adopted incomplete LU factorization method (ILU) mostly due to its lower memory requirements.

The whole Jacobian matrix  $\partial \mathbf{N}/\partial \mathbf{u}$  can be divided into  $N_e \times N_e$  small blocks, where  $N_e$  is the number of elements. The  $e_1$ th row and  $e_2$ th column block is  $\partial \mathbf{N}_{e_1}/\partial \mathbf{u}_{e_2}$  where  $\mathbf{N}_{e_1}$  is the nonlinear system residual of the  $e_1$ th element and  $\mathbf{u}_{e_2}$  is the independent variables of the  $e_2$ th element. The matrix  $\partial \mathbf{N}/\partial \mathbf{u}$  can be decoupled into

$$\frac{\partial \mathbf{N}}{\partial \mathbf{u}} = \hat{\mathbf{L}} + \hat{\mathbf{D}} + \hat{\mathbf{U}},\tag{27}$$

where  $\hat{\mathbf{L}}$ ,  $\hat{\mathbf{D}}$  and  $\hat{\mathbf{U}}$  correspond to the lower  $(e_1 < e_2)$ , diagonal  $(e_1 = e_2)$  and upper  $(e_1 > e_2)$  block part of  $\partial \mathbf{N}/\partial \mathbf{u}$ , respectively. Then the product  $\mathbf{P}^{-1}\mathbf{q}$  is implicitly calculated through the iteration,

$$\hat{\mathbf{q}}^{0} = \mathbf{0},$$

$$\hat{\mathbf{q}}^{j} = \omega \hat{\mathbf{D}}^{-1} \left[ \mathbf{q} - \left( \hat{\mathbf{L}} + \hat{\mathbf{U}} \right) \hat{\mathbf{q}}^{j-1} \right] + (1 - \omega) \hat{\mathbf{q}}^{j-1}, \ j = 1, 2, \cdots, J,$$
(28)

and  $\hat{\mathbf{q}}^J$  is the preconditioned vector. Currently the case with  $\omega=1$  is studied. To minimize the memory consumption, only the matrix  $\hat{\mathbf{D}}^{-1}$  is stored while the product  $(\hat{\mathbf{L}} + \hat{\mathbf{U}}) \hat{\mathbf{q}}^{j-1}$  is calculated on the fly. J=7 is used in this paper, if not specified.

The matrix  $\partial \mathbf{N}_{e_1}/\partial \mathbf{u}_{e_2}$  can be derived from Eqs (13) and (20) as

$$\frac{\partial \mathbf{N}_{e_1}}{\partial \mathbf{u}_{e_2}} = \mathbf{I} \delta_{e_1 e_2} - \Delta t a^{mm} \mathbf{M}^{-1} \frac{\partial \mathcal{L}_{e_1}}{\partial \mathbf{u}_{e_2}} \quad , \tag{29}$$

$$\frac{\partial \mathcal{L}_{e_{1}}}{\partial \mathbf{u}_{e_{2}}} = \sum_{j=1}^{d} \mathbf{B}_{e_{1}}^{T} \mathbf{D}_{j,e_{1}}^{T} \mathbf{\Lambda} \left[ w_{e_{1}} J_{e_{1}} \left( \frac{\partial \mathbf{H}_{j,e_{1}}}{\partial \mathbf{Q}_{e_{1}}} \right)_{\nabla \mathbf{Q}} \right] \mathbf{B}_{e_{1}} \delta_{e_{1}e_{2}} + \\
\sum_{i=1}^{d} \sum_{j=1}^{d} \mathbf{B}_{e_{1}}^{T} \mathbf{D}_{j,e_{1}}^{T} \mathbf{\Lambda} \left[ w_{e_{1}} J_{e_{1}} \left( \frac{\partial \mathbf{H}_{j,e_{1}}}{\partial \nabla_{i} \mathbf{Q}_{e_{1}}} \right)_{\mathbf{Q}} \right] \mathbf{D}_{i,e_{1}} \mathbf{B}_{e_{1}} \delta_{e_{1}e_{2}} + \\
\left( \mathbf{B}_{e_{1}}^{\Gamma} \mathbf{M}_{c} \right)^{T} \mathbf{\Lambda} \left[ w^{\Gamma} J^{\Gamma} \left( \frac{\partial \hat{\mathbf{H}}_{e_{1}}^{n}}{\partial \hat{\mathbf{Q}}^{+/-}} \right)_{\nabla \mathbf{Q}} \right] \frac{\partial \hat{\mathbf{Q}}^{+/-}}{\partial \mathbf{Q}_{e_{2}}} \mathbf{B}_{e_{2}} + \\
\sum_{j=1}^{d} \left( \mathbf{B}_{e_{1}}^{\Gamma} \mathbf{M}_{c} \right)^{T} \mathbf{\Lambda} \left[ w^{\Gamma} J^{\Gamma} \left( \frac{\partial \hat{\mathbf{H}}_{e_{1}}^{n}}{\partial \nabla_{j} \hat{\mathbf{Q}}^{+/-}} \right)_{\mathbf{Q}} \right] \frac{\partial \nabla_{j} \hat{\mathbf{Q}}^{+/-}}{\partial \nabla_{j} \mathbf{Q}_{e_{2}}} \mathbf{D}_{j,e_{2}} \mathbf{B}_{e_{2}} + \\
+ \sum_{j=1}^{d} \mathbf{B}^{T} \mathbf{D}_{j}^{T} \mathbf{J}^{T} \mathbf{\Lambda} \left[ w^{\Gamma} J^{\Gamma} \left( \frac{\partial \hat{\mathbf{S}}_{e_{1}}^{n}}{\partial \hat{\mathbf{Q}}^{+/-}} \right) \right] \frac{\partial \hat{\mathbf{Q}}^{+/-}}{\partial \mathbf{Q}_{e_{2}}} \mathbf{B}_{e_{2}},$$

$$(30)$$

where  $\partial \hat{\mathbf{Q}}^{+/-}/\partial \mathbf{Q}_{e_2}$  is the interpolation matrix from  $\mathbf{Q}_{e_2}$  to  $\hat{\mathbf{Q}}^{+/-}$  which is independent of  $\mathbf{Q}$  itself, superscript +/- is used to indicate that only one of them is a non-zero matrix for a specific  $e_2$ . The term  $\partial \nabla_j \hat{\mathbf{Q}}^{+/-}/\partial \nabla_j \mathbf{Q}_{e_2}$  is similarly an interpolation matrix.

The first two terms on the right-hand side of Eq. (30) are element integration evaluations, which are calculated using analytical methods. Following the ideas described in reference [19], single precision data type, continuous

memory operations and matrix padding are used to reduce the CPU cost and optimize cache performance in the implementation of these evaluations.

The remaining three terms are element boundary terms. Approximations are adopted to reduce the cost of preconditioner and simplify the coding. The last two terms on the right-hand side of Eq. (30) are neglected in all the simulations for the following two reasons. First, they are much more expensive than the third term on the right-hand side of Eq. (30), which has similar computation cost as the boundary integration of the advection flux only. Considering the diffusion term is much more expensive in 3D NS simulations, this simplification makes the computation cost of  $(\hat{\mathbf{L}} + \hat{\mathbf{U}}) \hat{\mathbf{q}}^{j-1}$  much smaller than one  $\mathbf{N}(\mathbf{u})$  evaluation, as discussed in Section 5.1. Second, numerical tests show that neglecting these two terms has very small impact on the preconditioning performance. For instance, neglecting the last two terms in the compressible Poiseuille flow test case in Section 6.1 with P=2 and  $20^2$  meshes only increases the total number of GMRES iterations in the first 10 time steps from 1136 to 1144.

To further simplify the coding complexity, the flux Jacobians  $\partial \hat{\mathbf{H}}_{e_1}^n/\partial \hat{\mathbf{Q}}^{\pm}$ , the only functions of  $\mathbf{Q}$  in the third term on the right-hand side of Eq. (30), are calculated using a finite difference approximation to avoid coding the Jacobian matrices of various Riemann fluxes and boundary conditions. The finite difference approximation adopted is

$$\left(\frac{\partial \hat{\mathbf{H}}_{e_1}^n}{\partial \hat{\mathbf{Q}}^{\pm}}\right)_{\nabla \mathbf{Q}} \simeq \frac{\hat{\mathbf{H}}_{e_1}^n \left(\hat{\mathbf{Q}}^{\pm} + \chi \mathbf{e}_j, \hat{\mathbf{Q}}^{\mp}, \nabla \hat{\mathbf{Q}}\right) - \hat{\mathbf{H}}_{e_1}^n \left(\hat{\mathbf{Q}}^{\pm}, \hat{\mathbf{Q}}^{\mp}, \nabla \hat{\mathbf{Q}}\right)}{\chi}, \quad (31)$$

where  $\mathbf{e}_j$  is the vector with one on the jth column but zeros on others. The parameter  $\chi$  is evaluated in a fashion similar to  $\epsilon$  in Eq. (26). However, the  $L_2$  norm  $\|\mathbf{u}\|_2$  in Eq. (26) is replaced by the  $L_2$  norm of each variable, which makes  $\chi$  different for each variable. Specific physical boundary conditions are imposed on physical boundaries in the  $\hat{\mathbf{H}}_{e_1}^n$  evaluation so that the Jacobian matrices of the boundary conditions are already included in Eq. (31). Because of the independence of  $\hat{\mathbf{H}}^n$  on each quadrature point, in total  $2N_{var}$  numerical flux evaluations on all the quadrature points of element boundaries and  $N_{var}$  boundary treatment operations are required to calculate  $\partial \hat{\mathbf{H}}_{e_1}^n/\partial \hat{\mathbf{Q}}^{\pm}$ . The matrices of  $\partial \hat{\mathbf{H}}_{e_1}^n/\partial \hat{\mathbf{Q}}^{\pm}$  are stored for the evaluation of  $(\hat{\mathbf{L}} + \hat{\mathbf{U}}) \hat{\mathbf{q}}^{l-1}$  in Eqs. (28). The computation cost and storage consumption of the matrix  $\partial \hat{\mathbf{H}}_{e_1}^n/\partial \hat{\mathbf{Q}}^{\pm}$  is small compared with the other parts of the implicit solver, as shown in Section 5.

The hybrid analytical and numerical method offers an efficient and accurate way of calculating the preconditioning matrices. In addition, it greatly reduces the complexity in programming. For simulations with shocks, the terms with  $\partial \mu_{av}/\partial \mathbf{u}$  and  $\partial \kappa_{av}/\partial \mathbf{u}$  are neglected in the preconditioning matrices calculations. As shown in Section 6.5, the implicit solver is able to stably and efficiently simulate discontinuous flows neglecting these two terms. The preconditioning matrices are frozen for 10 time steps in the following simulations, if not specified.

The main simulation flowchart of the implicit solver is given in Fig. 2, which also shows the flow between different classes.

Initialization	in CompressibleFlowSolver	
Time step loop $n$	in UnsteadySystem	
Runge-Kutta loop $m = 1, \dots, M$	in TimeIntegrationScheme	
Calculate source term $\mathbf{S}_m$	in TimeIntegrationScheme	
Netwon iteration l	in NewtonSolver	
Calculate residual $\mathbf{N}\left(ar{\mathbf{u}}^l ight)$	in CompressibleFlowSystem	
GMRES iteration $k$	in NekLinSysGMRES	
Calculate search vector $\mathbf{q}_k$	in NekLinSysGMRES	
BRJ iteration $j = 1, \dots, J$	in PrecondBRJ	
Calculate $\hat{\mathbf{q}}^j = \hat{\mathbf{D}}^{-1} \left( \mathbf{q}_k - \left( \hat{\mathbf{L}} + \hat{\mathbf{U}} \right) \hat{\mathbf{q}}^{j-1} \right)$	in CompressibleFlowSystem	
Calculate $\partial \mathbf{N}/\partial \mathbf{u} \cdot \hat{\mathbf{q}}^J$	in CompressibleFlowSystem	
Calculate linear system residual	in NekLinSysGMRES	
Calculate $\bar{\mathbf{u}}^{l+1}$ by linear combination of $\mathbf{q}_k$	in NekLinSysGMRES	
Calculate $\mathbf{u}^{n+1}$	in TimeIntegrationScheme	
Output and finalization	in CompressibleFlowSolver	

Figure 2: Nassi–Shneiderman diagram of the implicit solver with the corresponding class names. Brown: library class ; cyan: solver class.

## 4.3. Solver capabilities and code verification

The main capabilities of the explicit and implicit compressible flow solvers are summarized here. The explicit solver can use forward Euler, Adams Bashforth [40] and 2nd-4th order Runge-Kutta for temporal discretization, while the implicit solver can use backward Euler, 2nd order BDF (Backward differentiation formula) and 2nd-4th order SDIRK. The other aspects of the solvers are summarized as follows:

- Supported equation systems: Euler, Navier-Stokes;
- Advection discretization schemes: DG (with various nodal or modal bases functions), FR [25] (not supported by the implicit solver yet);
- Diffusion discretization schemes: IP [16] (can recover a specific version of direct DG [8]), LDG [10] (the implicit solver is partially supported);
- Riemann numerical flux: AUSM, Roe, HLL, HLLC, Steger-Warming, Lax-Friedrichs and the iterative exact Riemann solver in Section 4.9 of [36];
- Boundary conditions: periodic, time dependent Dirichlet, symmetry, slipping wall, isothermal non-slip wall, adiabatic non-slip wall, Riemann invariant, pressure inflow, pressure outflow, zero-order extrapolation [24];
- Shock capturing: modal resolution-based shock sensor [32], Ducros' sensor [12], physical and Laplace-based artificial viscosity;
- Mesh types: triangular, quadrilateral, hexahedral, prismatic, pyramidal, tetrahedral meshes;
- Parallel: MPI based parallelization for HPC computing.

All of these options are supported by the explicit solver, but are only partially supported by the implicit one. Options not supported or partially supported by the implicit solver have been marked above.

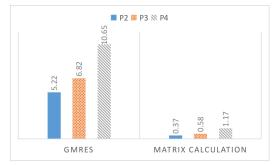
CTest facility of the CMake environment [6] are adopted to execute hundreds of test cases for testing different building blocks of the solver. These tests are automatically executed on a variety of platforms using a buildbot service to ensure the reliability of the solver [6].

# 5. Performance and memory consumption of the implicit solver

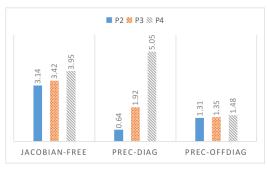
The implicit solver can significantly relax the constraint on the time step. However, it also requires much larger computational cost and memory consumption. The computational performance of different parts of the implicit solver is analyzed in Section 5.1, and the memory consumption is discussed in Section 5.2.

# 5.1. Performance analysis

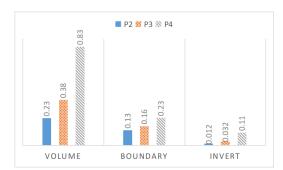
In order to provide a deeper insight of the implicit solver implementation, the computational costs of the different parts of the implicit solver are evaluated in the following. The evaluations are performed on a CentOS 7.7.1908 using standard -O3 compiler optimizations with the gcc 5.4.0 C++ compiler and version 2019.5.281 of the BLAS and Lapack implementation in the Intel MKL library. The evaluations are run on the HPC cluster of Imperial College London with one node and four nodes for the 2D and 3D cases, respectively. Each node is equipped with two Intel Xeon E5-2620 CPUs and 120 GB RAM. The computation CPU time is obtained using Oracle Developer Studio 12.6.



(a) The GMRES and preconditioning matrix calculations in the implicit solver



(b) Operations in GMRES



(c) Operations in preconditioning matrices calculation

Figure 3: Performance analysis of the implicit solver in 3D

#### 5.1.1. 3D case

The 3D cylinder test case described in Section 6.4 is adopted here for the performance analysis. The simulations are carried out at polynomial orders

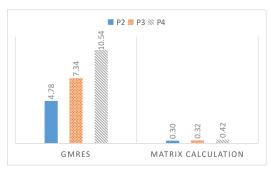
 $P=2,\,P=3$  and P=4, using seven preconditioning iterations (J=7) and a time step corresponding to CFL=40. The Newton iteration tolerance is  $\alpha=10^{-3}$ . In total 10 steps are simulated, with the preconditioning matrices calculated at the beginning of the simulations and frozen in the following 10 time steps. Fig. 4 shows the CPU time spent in different parts of the implicit solver, which is normalized by the total CPU time spent in evaluating  $\mathbf{N}(\mathbf{u})$  the same number of times as the total number of Newton iterations. Fig. 4a shows that the GMRES procedure, which calls the Jacobian-free operation and the preconditioners, makes up the majority of CPU time. The calculation of preconditioning matrices is much cheaper because of the freezing of the matrices. However, if the preconditioning matrices are updated every time step, their computational cost will be similar to that of the GMRES procedure and the implicit solver will be about twice as expensive as the current simulations.

Fig. 4b shows the three most expensive operations in the GMRES procedure: the JACOBIAN-FREE, PREC-OFFDIAG and PREC-DIAG. The JACOBIAN-FREE procedure corresponds to the operation of Eq. (25), in which the majority of CPU time is spent on one N(u) evaluation. As a result, the value of the normalized CPU time of the JACOBIAN-FREE correctly indicates the averaged number of GMRES iterations in each Newton iteration. The PREC-OFFDIAG corresponds to the  $(\hat{\mathbf{L}} + \hat{\mathbf{U}})$  vector inner product in Eq. (28), while the PREC-DIAG corresponds to the  $\hat{\mathbf{D}}^{-1}$  vector inner product in Eq. (28). Since the BRJ(7) is used with 0 initial guess, six PREC-OFFDIAG and seven PREC-DIAG are needed for each JACOBIAN-FREE operation. Although much more preconditioning operations are performed, their costs are smaller than that of the JACOBIAN-FREE operations for P=2 and P=3. Even though we have already adopted the single precision data type to reduce the operations of the matrix-vector inner product, the  $\hat{\mathbf{D}}^{-1}$  vector inner product still grows much faster than the other two parts when P increases. This is mainly because it uses a matrix-based implementation while most operations in the other parts use a sum-factorization implementation. This also indicates that matrix-based preconditioners will easily lose efficiency for high-order 3D simulations. Techniques proposed in reference [2] may be used to implement a matrix-free version of the BRJ preconditioner to avoid the fast growth of the  $\hat{\mathbf{D}}^{-1}$  vector inner product. As for the off-diagonal operations, the simplifications and implementations make sure they are relatively cheap compared with other parts.

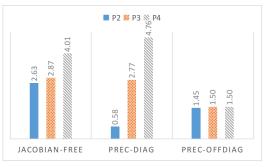
Fig. 4c shows the most expensive parts in the preconditioning matrices

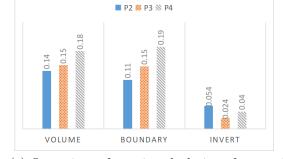
calculations. For the element volume part, techniques such as using single precision, maintaining memory continuous and padding the matrices as discussed in reference [19], have been adopted in our implementation to speed-up the calculation. However, this is still the most expensive part. In our implementation, the most expensive part of the volume term is implemented using BLAS::sger, but it still grows quite fast as P increases. Inverting  $\hat{\mathbf{D}}$ , which is implemented using Lapack::dgetrf and Lapack::dgetri, only makes up a small proportion in the test cases. However, it grows the fastest, which indicates it will become the dominant part when the polynomial order is high enough. In this case, techniques such as approximate matrix inversion [2, 30, 11] may be adopted.

Fig. 4 shows that a large proportion of computation costs in the implicit solver is from operations shared with the explicit solver. Therefore, the implicit solver will also greatly benefit from improvements in efficiency of the explicit solver.



(a) The GMRES and preconditioning matrix calculations in the implicit solver





- (b) Operations in GMRES
- (c) Operations of matrix calculation of preconditioner

Figure 4: Performance analysis of the implicit solver in 2D

# 5.1.2. 2D case

The 2D cylinder test cases of Section 6.2 are adopted for the performance analysis of the solver in 2D. The simulations are run at CFL=20 for 200 steps with preconditioning matrices updated every 10 time steps. The other parameters are the same as those of the 3D tests. Fig. 4 shows the CPU time of different parts of the implicit solver, which are normalized in the same way. Very similar trends as in the 3D case are observed. The  $\hat{\mathbf{D}}^{-1}$  vector inner product still grows much faster than the other parts, which means matrix-based preconditioners may also lose efficiency because of the large computational cost. However, the growth of CPU time in preconditioning matrices calculation is much slower than that in 3D case. The CPU time in inverting the  $\hat{\mathbf{D}}$  matrix is not increasing monotonously as that of the 3D case. A possible reason is that the functions in intel MKL library are sensitive to the rank of the matrix when the matrix size is small, as discussed in reference

# 5.2. Memory consumption analysis

Tab. 1 compares the memory requirements of the preconditioning matrices and Jacobian matrix  $\partial \mathbf{N}/\partial \mathbf{u}$  using the 3D cylinder case of Section 5.1.1. The memory consumption of the Jacobian-Sparse is that of the Jacobian matrix considering the sparsity patterns of the off diagonal blocks of  $\partial \mathbf{N}/\partial \mathbf{u}$ , which is calculated based on the estimation in reference [31]. The memory consumption of the Jacobian-Dense is used by storing the diagonal and off diagonal blocks as dense matrices. Note that the ILU preconditioner with zero fill-in should have the same memory consumption as the Jacobian-Dense since the matrices will generally lose sparsity after the factorization. The memory consumption of the Jacobian matrix becomes large in 3D. Only approximately 45 elements can be afforded per Gbyte of memory for simulations with P = 4 and hexahedral meshes if the Jacobian-Dense matrix is stored or ILU preconditioner is used. Moreover, frequently loading these large matrices into the CPU cache will seriously influence the computational speed since most modern high-performance computers (HPCs) are memory bandwidth limited [43]. The memory consumption can be further reduced if using other schemes like the CDG [31] by exploring the sparsity patterns, but it is still of the same order of magnitude as the Jacobian-Sparse.

Memory in operations BRJ-Diag and BRJ-Offdiag is required for storing  $\hat{\mathbf{D}}^{-1}$  and  $\partial \hat{\mathbf{H}}^n/\partial \hat{\mathbf{Q}}^{\pm}$ , both of which are stored in single precision to reduce memory and CPU costs. The storage of  $\partial \hat{\mathbf{H}}^n/\partial \hat{\mathbf{Q}}^{\pm}$  is small compared with that of  $\hat{\mathbf{D}}^{-1}$  especially for high order approximations. The total storage of BRJ is much smaller than that of the Jacobian matrix or the ILU preconditioner.

	BRJ-Diag	BRJ-Offdiag	BRJ-Total	Jacobian-Dense	Jacobian-Sparse
P=2	72 900	30 000	102 900	1 020 600	712 800
P=3	409 600	43 200	452 800	5 734 400	3 328 000
P=4	1 562 500	76 800	1 639 300	21 875 000	11 000 000

Table 1: Memory consumption of BRJ in bytes for each hexahedral mesh

## 6. Verification and applications

In the following, the implementations of spatial discretization methods are verified using accuracy tests described in Section 6.1. Section 6.2 discusses the temporal accuracy and efficiency of the solver in unsteady laminar flows over a circular cylinder. The temporal accuracy in turbulent simulations is studied in Section 6.3 using a Taylor-Green vortex case. Section 6.4 presents and discusses turbulent flow over a circular cylinder at Re = 3900 obtained with the implicit solver. Finally, a shock wave boundary layer interaction is studied to demonstrate the shock capturing ability of the compressible flow solver.

In this section, we use two different definitions of Courant-Friedrichs-Lewy (CFL) number which are important for the choices of time steps, are presented. The "standard" CFL number is defined as

$$CFL = \frac{c_{\lambda} \triangle t \left(c + \sqrt{u_k u_k}\right) P^2}{d_{RK} \triangle x},$$
(32)

where  $c_{\lambda} = 0.2$ ,  $d_{RK} = 2$  and c is the speed of sound. The time step of the explicit solver will be chosen based on Eq. (32), which ensures the maximum stable CFL is almost constant for different values of the polynomial order P [20]. Another CFL number, the convective CFL number, is defined as

$$CFL_c = \frac{\Delta t \sqrt{u_k u_k}}{\Delta x / P}.$$
 (33)

A value  $CFL_c = 1$  means the flow field can convect a distance of  $\Delta x/P$  in one time step, where the distance  $\Delta x/P$  in Eq. (33) is an estimate of the average length between two DoFs in the element. The time step of the implicit solver maintains a small  $CFL_c$ , which ensures the temporal error is similar to or smaller than the spatial error.

#### 6.1. Accuracy test

Testing the convergence order of accuracy (OoA) of the solver is a very effective way of code verification [34]. Here two different test cases are used to verify two different aspects of the solver. An isentropic vortex convection problem is used to verify the advection scheme and the treatment of the diffusion terms is verified by means of a manufactured compressible Poiseuille flow.

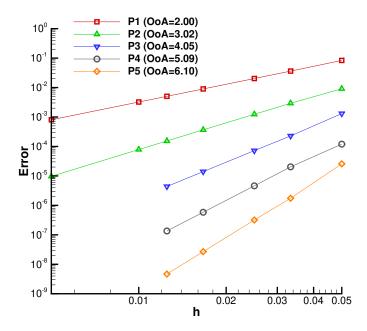


Figure 5: Accuracy test of the advection scheme

*Isentropic vortex convection*. In this case an inviscid isentropic vortex is convected down stream. In computational domain  $[0, 10]^2$ , the exact solution at time t is

$$\rho = \left(1 - \frac{\varphi^{2}(\gamma - 1)}{16\gamma\pi^{2}}e^{2(1-r^{2})}\right)^{\frac{1}{\gamma - 1}}$$

$$u = u_{0} + \frac{\varphi(y - y_{0})}{2\pi}e^{(1-r^{2})},$$

$$v = v_{0} + \frac{\varphi(x - x_{0})}{2\pi}e^{(1-r^{2})}$$

$$n = \rho^{\gamma}$$
(34)

where  $(u_0, v_0) = (1.0, 0.5)$  is the mean convection velocity field, the coordinates of the vortex center at time t are given by  $(x_0 + u_0t, y_0 + v_0t)$ , r is the distance from the vortex center, and  $\varphi = 5.0$  is a parameter that controls the strength of the vortex. Mesh convergence is analyzed through the use of progressively refined quadrilateral meshes. Periodic boundary conditions are applied to the boundaries in both directions. A fourth-order explicit Runge-Kutta with CFL=0.01 is adopted to guarantee the errors are dominated by

the spatial discretization. The error distributions are depicted in Fig. 5 with the order of accuracy (OoA) between the two finest mesh levels given in the legend. All simulations with polynomial orders P from 1 to 5 reach the designed order of accuracy.

Manufactured compressible Poiseuille flow. The method of manufactured solution (MMS) permit us to design specific test cases with analytical solutions for the study of different aspects of the solvers [34]. Using the MMS, a test case similar to Poiseuille flow [27] is designed to verify the DG methods for the diffusion terms. This problem is suitable for the verification of the IP method since it is dominated by the diffusion term. By design, the analytical solution is

$$\rho = 1$$

$$u = -\frac{1}{2\mu} \frac{dp}{dx} \left[ y (L - y) + \theta y^2 (L - y)^2 \right]$$

$$v = 0$$

$$p = \frac{1}{\gamma \text{Ma}^2} + \frac{dp}{dx} x$$
(35)

where L=1 is the height of the computational domain in the y direction. The corresponding forcing terms are given in Eqs. (36) in the Appendix and correspond to a free-stream flow with Ma = 0.1 and Re = 100 based on the maximum velocity, which is normalized to 1. The corresponding pressure gradient is  $dp/dx = -8\mu/L^2$ . Compared to incompressible Poiseuille flow, an extra term  $\theta y^2 (L-y)^2$  with  $\theta=0.01$  is added so that the analytical solution of the energy E, shown in Fig. 6, has a polynomial order of 8. This ensures the polynomial order of the solution is higher than that of the base functions and thus avoids the high-order coefficients becoming excessively small. A series of quadrilateral meshes with  $3^2$ ,  $4^2$ ,  $8^2$ ,  $12^2$ ,  $16^2$  and  $20^2$  elements are used with Dirichlet boundary conditions based on the analytical solution. The  $L_2$  and  $L_\infty$  norm error distributions are shown in Fig. 7. The results show that the designed order of accuracy is achieved for the IP method from P=2 to P=5. Simulations with triangular meshes can also achieve the designed order of accuracy but are not shown here.

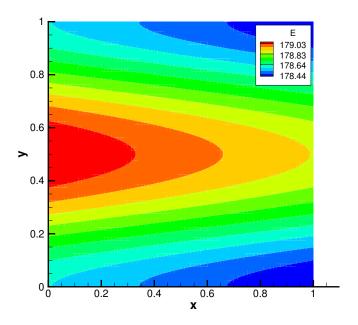


Figure 6: Total energy (E) distribution of the manufactured compressible Poiseuille flow

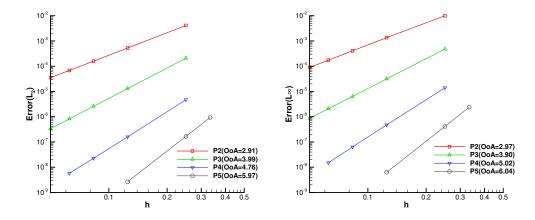


Figure 7: Convergence order of the IP methods in compressible Poiseuille flow

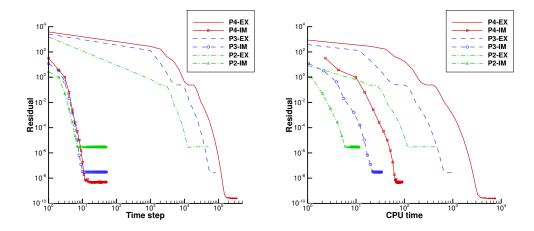


Figure 8: Residual distribution of implicit and explicit simulations in compressible Poiseuille flow

The efficiency of the two solvers is also compared in this steady problem. The initial condition is similar to that used in Eqs. (35) but with uniform velocity  $u = -(dp/dx)L^2/(8\mu)$ . The time integration used by the implicit solver is SDIRK2 whilst the explicit solver uses ERK2. The CFL number of the implicit solver grows from 500 to 5000 within the first 10 steps. Here the preconditioning matrices are updated in every other time step. The explicit solver uses its maximum stable CFL numbers of 0.08, 0.06 and 0.06 for P = 2, P = 3 and P = 4, respectively. Fig. 8 shows the convergence history of the residual norm  $\|\mathbf{M}^{-1}\mathcal{L}\|_2$ . All the implicit simulations reach steady state within 20 steps, and the explicit solver takes about 13000, 60000 and 170000 steps for P = 2, P = 3 and P = 4 to converge, respectively. Fig. 8 shows that even with no grid stretching the implicit solver is about 18.2, 33.1 and 53.4 times faster than the explicit solver for P = 2, P = 3 and P = 4, respectively.

Following references [27, 34], Dirichlet boundary conditions based on the analytical solution are applied weakly in the DGM in this subsonic problem. The difference between analytical solutions and numerical solutions will lead to large errors near boundaries, which prevents the simulations from converging to machine epsilon levels. However, the errors caused by boundary treatment does not invalidate the OoA tests since the difference between the numerical and analytical solutions on the boundaries tends to zero as the resolution increases.

# 6.2. Laminar flow over a circular cylinder

The temporal accuracy and efficiency of the implicit solver is tested next using unsteady simulations of flow over a circular cylinder. The parameters and meshes employed are similar to those presented in reference [3]. The flow conditions correspond to Ma = 0.3, Re =  $\rho uD/\mu$  = 1200 and we use 64 × 60 quadrilateral meshes. Starting from the same flow field, different time integration schemes and time steps are used to get the solutions after approximately 1.7 vortex shedding periods. The integrated lift and drag forces are the chosen metrics for the accuracy test with reference solutions obtained by the SDIRK4 in Tab. 16 of reference [21] with a very small time step  $\Delta t = 1 \times 10^{-5}$ . A smaller Newton iteration tolerance  $\alpha = 10^{-6}$  is used for SDIRK3 and SDIRK4. The convergence rates are presented in Fig. 9, which shows that the implicit solver achieves the desired order of accuracy.

The efficiency of the implicit solver is compared with that of the explicit solver in this test case. Starting from the same flow field, the solution is obtained after around 30 vortex shedding periods using both explicit and implicit solver. The details of the simulations are summarized in Tab. 2. The implicit solver can run with a time step 330 times larger than that of the explicit solver. The corresponding  $CFL_c$  is about 2.5 for this time step. The implicit solver is around 20.3 times quicker than the explicit solver. The implicit solver is also ran with a larger time step  $(CFL_c \simeq 5)$ , which achieves a speed-up of 27.9. By observing the evolution history of the flow field, it can be observed that the temporal errors mainly lead to a dispersion in the vortex shedding process. Thus the Strouhal number  $(Sr = \frac{Df}{u})$  with f the vortex shedding frequency) is a good metric to compare the temporal accuracy. Tab. 2 shows that all the simulations give almost the same prediction of Strouhal number (0.2419). The results show that the implicit solver can largely improve efficiency while keeping good temporal accuracy.

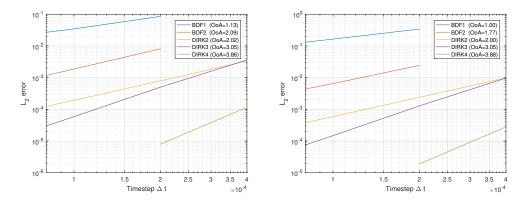


Figure 9: Convergence rate of different implicit time integration schemes

Method	ERK2	SDIRK2	
$\Delta t$	$4.5 \times 10^{-5}$	$1.5 \times 10^{-2}$	$2.0 \times 10^{-2}$
CFL	0.06	20	40
$\mathrm{CFL}_c$	$7.5 \times 10^{-3}$	2.5	5.0
Strouhal number	0.2419	0.2419	0.2419
Newton iterations	\	2.1	2.5
per RK stage			
GMRES iterations		3.5	5.0
per RK stage	,		
Speed-up	1.0	20.3	27.9

Table 2: Efficiency comparison between explicit and implicit time integration schemes

Polynomial order	P=2		P=3	
Method	ERK2	SDIRK2	ERK2	SDIRK2
$\Delta t$	$1.1 \times 10^{-3}$	0.1	$6.1 \times 10^{-4}$	0.046
CFL	0.08	7.5	0.10	7.5
$\mathrm{CFL}_c$	$2.6 \times 10^{-2}$	2.5	$1.8 \times 10^{-2}$	1.4
Newton iterations per RK stage	\	4	\	5
GMRES iterations per Newton iteration	\	3.7	\	4.6
Speed-up	1.0	2.76	1.0	1.2

Table 3: Efficiency comparison between explicit and implicit time integration schemes of TGV

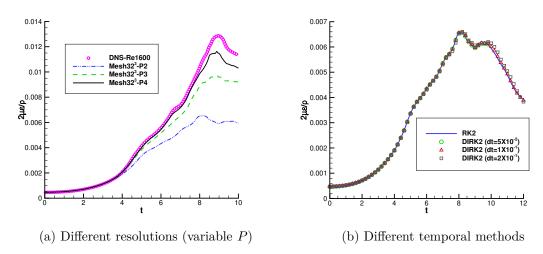


Figure 10: Evolution of the dissipation rate in Taylor-Green vortex

# 6.3. Taylor-Green vortex

The evolution of the Taylor-Green vortex (TGV) includes the break down of initially large vortices, transition to turbulence and decay of turbulence. With a very simple set up but a complex flow evolution, it is a very good test case for transitional and turbulent simulations. Initially, large vortices are

placed in a cubic domain  $[-\pi,\pi]^3$  with periodic boundary conditions. Flow conditions corresponding to Ma = 0.1 and Re = 1600 are used. The expression of the initial flow field can be found in reference [42]. The value  $2\mu\varepsilon/\rho$ , where  $\varepsilon$  is the enstrophy, is a good estimate of the energy dissipation rate at the incompressible limit [42]. Fig. 10a displays evolution curves of  $2\mu\varepsilon/\rho$  obtained with  $32^3$  meshes and using resolutions at different polynomial orders. Dealiasing via over-integration with 3(P+1)/2 quadrature points in each spatial direction [44] is used in the under-resolved simulations of TGV. The curves gradually tend to the DNS result [42] as the resolution increases, which shows the implicit solver can correctly resolve the main flow phenomena. All these simulations are run using SDIRK2 with  $CFL_c$  smaller than 2.5. The details of the implicit solver and their comparisons with ERK2 scheme are presented in Tab. 3. Although there is no grid stretching in this test case, the implicit solver can still be as efficient as, or more efficient than, the explicit solver. Fig. 10b compares results calculated using SDIRK2 with different time steps with the result obtained with ERK2, which is very accurate in time because of the very small time step used (CFL = 0.08). Three different time steps of 0.05, 0.1 and 0.2 are used for SDIRK2. Although the corresponding temporal errors of SDIRK2 should change by 16 times for different time steps, all these results coincide quite well with the results of ERK2. This means the temporal error is very small for the implicit solver. Given the large differences for different spatial resolutions, the implicit solver with the current parameter choice is accurate enough in the sense that the error of the solution is dominated by the spatial error. Finally, Fig. 11 shows instantaneous vorticity distributions obtained using ERK2 and SDIRK2. They are in good agreement with each other, which further verify the accuracy of the implicit solver.

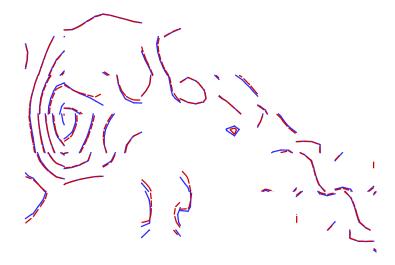


Figure 11: Comparison of contours of surface normal vorticity on  $y=\pi$  at t=8 in Taylor-Green vortex. Comparison between ERK2 (red dashed lines) and SDIRK2 with  $\Delta t=0.1$  (blue solid lines).

# 6.4. Turbulent flow over a circular cylinder at Re = 3900

The flow over a circular cylinder at Re = 3900 is another widely used test case for turbulence simulations. Various experimental and numerical results can be found for this test case, e.g. [29, 23]. Compared with the TGV, there are extra complexities because of grid stretching, the need for simulating the boundary layer, and the instability of the shear layer. The computational domain boundaries are located 40D away down stream and 20D away on other directions. Riemann invariant based boundary conditions are applied at the outer boundaries. The mesh near the cylinder is shown in Fig. 12. The mesh consists of 123 360 curved unstructured hexahedral elements in total. which corresponds to about 3.3M DoFs for P=2. The smallest mesh size is approximately 0.006D, which is located on the cylinder surface and is chosen based on the estimate advocated by reference [23]. As in the TGV test case, we use P=2 with over-integration. The computed vortex structure using P=2 is illustrated in Fig. 12, which shows the complex turbulent flow structures in the wake of the cylinder. The time-averaged velocity distributions in Fig. 13 show that the current result is in good agreement with the

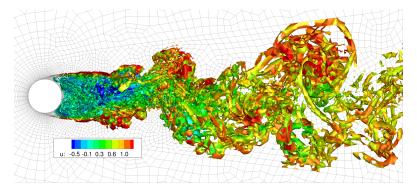


Figure 12: Q criteria iso-surface (Q=5) and the mesh of turbulent flow over a circular cylinder at Re = 3900

experimental results [29].

As shown in Tab. 4, a value CFL = 40 is employed in the simulations. The CFL<sub>c</sub> is around 2.5 at this CFL number, which ensures the temporal accuracy of the simulations. The corresponding time step is about 270 times larger than the explicit one with CFL = 0.15. On average, the implicit solver converges using 3 Newton iterations per implicit RK stage and 3.0 GMRES iterations per Newton iterations. The implicit solver is 14.7 times faster than the explicit one. The P=3 case is also run for efficiency comparison and, as shown in Tab. 4, a speed-up of 12.2 is achieved compared with the explicit solver. The seven stages third-order low-storage ERK scheme (ERK3-7) of Tab. A.15 of reference [37] is also compared, which runs at a much larger CFL number (CFL = 0.6). The implicit solver is 11.5 times faster than the ERK3-7 for the simulation with P=2.

The influence of the Newton iteration tolerance  $\alpha$  is also studied here. With  $\alpha = 10^{-4}$ ,  $\alpha = 10^{-5}$ ,  $\alpha = 10^{-6}$  and  $\alpha = 10^{-7}$ , the CPU time normalized by CPU time with  $\alpha = 10^{-3}$  is 1.0, 1.3, 1.6 and 1.7, respectively. Therefore decreasing  $\alpha$  will not seriously slow down the simulations.

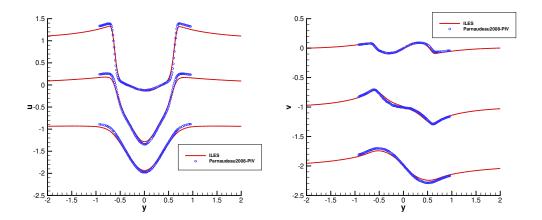


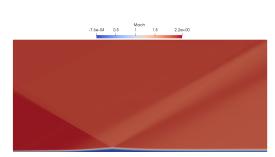
Figure 13: Time averaged velocity distributions of turbulent flow over a circular cylinder at  $\mathrm{Re} = 3900$ 

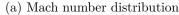
Polynomial order	P=2		P=3	
Method	ERK2	SDIRK2	ERK2	SDIRK2
$\Delta t$	$3.5 \times 10^{-5}$	$1.0 \times 10^{-2}$	$1.56 \times 10^{-5}$	$5.0 \times 10^{-3}$
CFL	0.15	40	0.15	48
$\mathrm{CFL}_c$	$9.0 \times 10^{-3}$	2.5	$6.0 \times 10^{-3}$	1.7
Newton iterations	\	3	\	3
per RK stage	\	3	\	3
GMRES iterations		3.0		3.7
per Newton iteration	\	5.0	\	5.7
Speed-up	1.0	14.7	1.0	12.2

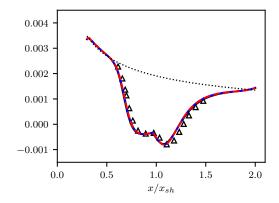
Table 4: Efficiency comparison between explicit and implicit time integration schemes of turbulent circular cylinder at  $\mathrm{Re}=3900$ 

# 6.5. Shock wave boundary-layer interaction

Finally, the shock capturing ability of the implicit solver is tested using the shock wave boundary-layer interaction (SWBLI) problem [4]. The computational domain and Mach number distribution are displayed in Fig. 14a. The viscous wall starts from x = 0 and the shock impingement position is  $x_{sh}$ . The computational domain starts from  $x = 0.3x_{sh}$  where the analytical compressible boundary layer solution [41] is imposed. Rankine-Hugoniot







(b) Skin friction coefficient distribution. Blue line: SDIRK-2; Red line: ERK2; triangles: simulation from [4]; dotted line: empirical solution by [13]

Figure 14: Mach number and skin friction coefficient distributions in shock wave boundary-layer interaction

relations are added to impose inflow boundary conditions consistent with the shock. Flow conditions corresponding to Ma = 2.15, Re $_{x_{sh}}$  =  $10^5$  and shock impingement angle of 30.8 are used. The simulations run with  $119 \times 39$  quadrilateral elements and P=3. The skin friction coefficient distributions are shown in Fig. 14b. The results of the explicit and implicit solvers coincide and are in good agreement with the simulation results presented in reference [4]. The efficiency of the solvers is compared in Tab. 5. Compared with the explicit solver, larger speed-up can be achieved when the CFL number of the implicit solver increases from 2.5 to 100. However, further increasing the CFL number to 1000, the implicit solver will have difficulty converging. The optimal numbers of preconditioning iterations are also given in Tab. 5. With a smaller time step, the system is less stiff and the optimal preconditioning number is smaller. Since the test case is regarded as a steady-state problem, a large CFL number can be used without influencing the accuracy. A speed-up of 22.8 is achieved with CFL=100.

Polynomial order	ERK2	SDIRK2		
CFL	0.08	2.5	5.0	100.0
$\Delta t$	$1.3 \times 10^{-4}$	$4.0 \times 10^{-3}$	$8.0 \times 10^{-3}$	$1.6 \times 10^{-2}$
J in BRJ	\	3	5	7
Newton iterations per RK stage	\	2.0	2.0	3.0
GMRES iterations per Newton iteration	\	3.0	3.2	8.9
Speed-up	1.0	2.7	4.2	22.8

Table 5: Efficiency comparison between explicit and implicit time integration schemes of shock wave boundary-layer interaction

#### 7. Discussion and future work

An implicit flow solver for the compressible Navier-Stokes equations has been developed in Nektar++ to improve the simulation efficiency using singly diagonally implicit Runge-Kutta and Jacobian-free Newton Krylov methods. To improve the efficiency, an efficient block relaxed Jacobi (BRJ) preconditioner is proposed, the computational cost and memory consumption of which are reduced by approximating the element boundary integrations, using single data precision and using a matrix-free implementation of the off diagonal terms. Performance analysis shows that the proposed BRJ preconditioner is low in computational cost at least for P=2 and P=3, which enables the solver to use more preconditioning iterations. The memory consumption of BRJ is significantly smaller than an ILU preconditioner. The software design and implementation details have been described with emphasis on providing a general pattern for the development of implicit solvers and for reducing the memory footprint. The solver capabilities were summarized and demonstrated using a variety of verification test cases. In the isentropic vortex convection, manufactured compressible Poiseuille flow and laminar flow over a circular cylinder test cases, the results achieved the designed convergent order of accuracy, thus verifying the implementations of the advection, diffusion and temporal discretization methods. The results of the TGV test case showed that the implicit solver achieves good temporal accuracy in the sense that the temporal error is much smaller than the spatial error. The implicit solver produced accurate predictions of turbulent flow over a circular cylinder at Re=3900 and shock wave boundary-layer interaction. The computational efficiency of the implicit solver is more than 10 times that of the explicit solver in both steady and unsteady subsonic simulations as well as steady supersonic simulations. The implicit solver is recommended for simulations with large grid stretching and/or low Mach number.

Further studies will focus on the development and implementation of efficient preconditioning methods. Moreover, the performance of the implicit solver for supersonic simulations has only been preliminarily investigated and will require to be studied thoroughly.

## Acknowledgements

The development the implicit solver in Nektar++ has been supported by EPSRC grant (EP/R029423/1) and UK Turbulence Consortium grant (EP/R029326/1). We would like to gratefully acknowledge access to the computing facilities provided by the Imperial College Research Computing Service (DOI: 10.14469/hpc/2232). Zhen-Guo Yan acknowledges support from the National Natural Science Foundation of China (Grant No. 11902344).

# Appendix A: Butcher tableaus of Runge-Kutta schemes

The coefficients of Runge-Kutta schemes are given using Butcher tableaus. The coefficients of the explicit second-order Runge-Kutta scheme (ERK2), the second-order singly diagonally implicit Runge-Kutta scheme (SDIRK2), the third-order singly diagonally implicit Runge-Kutta scheme (SDIRK3), the four-order singly diagonally implicit Runge-Kutta scheme (SDIRK4), are given in Tab. 6.

# Appendix B: Forcing terms of manufactured compressible Poiseuille flow

The analytical solution of the manufactured compressible Poiseuille flow has already been given in Eq. (35). The corresponding forcing terms are given by

$$f = \begin{pmatrix} 0 \\ \frac{dp}{dx}\theta \left(L^2 - 6yL + 6y^2\right) \\ 0 \\ f_4 \end{pmatrix}, \tag{36}$$

$$\begin{array}{c|ccccc}
0.0 & \lambda & \lambda \\
\hline
1.0 & 1.0 & 1.0 & 1.0 - \lambda & \lambda \\
\hline
0.5 & 0.5 & 1.0 - \lambda & \lambda \\
\end{array}$$
(a) ERK2 (b) SDIRK2 with  $\lambda = 1 - \sqrt{2}/2$ .

$$\begin{array}{c|ccccc}
\lambda & \lambda & \lambda \\
\hline
\frac{1+\lambda}{2} & \frac{1-\lambda}{2} & \lambda \\
\end{array}$$

Table 6: Butcher tableaus of Runge-Kutta schemes

where

$$f_4 = \frac{1}{2\mu} \frac{dp^2}{dx} (y(L-y) \left( 1 - \theta(L^2 - 6yL + 6y^2) - \gamma/(\gamma - 1) \right) - \frac{1}{4\mu} \frac{dp^2}{dx} (L - 2y)^2 (2\theta yL - 2\theta y^2 + 1)^2.$$
(37)

### References

- [1] Arnold, D., Brezzi, F., Cockburn, B., Marini, L.: Unified analysis of discontinuous Galerkin methods for elliptic problems. SIAM Journal on Numerical Analysis **39**(5), 1749–1779 (2002)
- [2] Bastian, P., Müller, E.H., Müthing, S., Piatkowski, M.: Matrix-free multigrid block-preconditioners for higher order discontinuous Galerkin discretisations. Journal of Computational Physics 394, 417–439 (2019)
- [3] Bijl, H., Carpenter, M.H., Vatsa, V.N., Kennedy, C.A.: Implicit time integration schemes for the unsteady compressible Navier–Stokes equations: Laminar flow. Journal of Computational Physics 179(1), 313–329 (2002)
- [4] Boin, J.P., Robinet, J.C., Corre, C., Deniau, H.: 3D steady and unsteady bifurcations in a shock-wave/laminar boundary layer interaction: A numerical study. Theoretical and Computational Fluid Dynamics 20(3), 163–180 (2006)
- [5] Brune, P.R., Knepley, M.G., Smith, B.F., Tu, X.: Composing scalable nonlinear algebraic solvers. SIAM Review **57**(4), 535–565 (2015)
- [6] Cantwell, C.D., Moxey, D., Comerford, A., Bolis, A., Rocco, G., Mengaldo, G., De Grazia, D., Yakovlev, S., Lombard, J.E., Ekelschot, D., Jordi, B., Xu, H., Mohamied, Y., Eskilsson, C., Nelson, B., Vos, P., Biotto, C., Kirby, R.M., Sherwin, S.J.: Nektar++: An open-source spectral/hp element framework. Computer Physics Communications 192, 205–219 (2015)
- [7] Cantwell, C.D., Sherwin, S.J., Kirby, R.M., Kelly, P.H.J.: From h to p efficiently: Strategy selection for operator evaluation on hexahedral and tetrahedral elements. Computers & Fluids 43(1), 23–28 (2011)

- [8] Cheng, J., Yang, X., Liu, X., Liu, T., Luo, H.: A direct discontinuous Galerkin method for the compressible Navier-Stokes equations on arbitrary grids. Journal of Computational Physics **327**, 484–502 (2016)
- [9] Chudanov, V., Aksenova, A., Goreinov, S., Makarevich, A., Pervichko, V.: Validation of a new method for solving of CFD problems in nuclear engineering using petascale HPC. pp. V004T10A009–V004T10A009. American Society of Mechanical Engineers (2014)
- [10] Cockburn, B., Shu, C.: The local discontinuous Galerkin method for time-dependent convection-diffusion systems. SIAM Journal on Numerical Analysis 35(6), 2440–2463 (1998)
- [11] Diosady, L.T., Murman, S.M.: Scalable tensor-product preconditioners for high-order finite-element methods: Scalar equations. Journal of Computational Physics 394, 759–776 (2019)
- [12] Ducros, F., Ferrand, V., Nicoud, F., Weber, C., Darracq, D., Gacherieu, C., Poinsot, T.: Large-eddy simulation of the shock/turbulence interaction. Journal of Computational Physics 152(2), 517–549 (1999)
- [13] Eckert, E.: Engineering relations for friction and heat transfer to surfaces in high velocity flow. Journal of the Aeronautical Sciences **22**(8), 585–587 (1955)
- [14] Ezertas, A., Eyi, S.: Performances of numerical and analytical jacobians in flow and sensitivity analysis. In: 19th AIAA Computational Fluid Dynamics. American Institute of Aeronautics and Astronautics, San Antonio, Texas (2009)
- [15] Franciolini, M., Crivellini, A., Nigro, A.: On the efficiency of a matrix-free linearly implicit time integration strategy for high-order discontinuous Galerkin solutions of incompressible turbulent flows. Computers and Fluids 159, 276–294 (2017)
- [16] Hartmann, R., Houston, P.: Symmetric interior penalty DG methods for the compressible Navier-Stokes equations I: Method formulation. International Journal of Numerical Analysis and Modeling **3**(1), 1–20 (2005)

- [17] Hartmann, R., Houston, P.: An optimal order interior penalty discontinuous Galerkin discretization of the compressible Navier–Stokes equations. Journal of Computational Physics **227**(22), 9670–9685 (2008)
- [18] Hesthaven, J.S., Warburton, T.: Nodal Discontinuous Galerkin Methods: Algorithms, Analysis, and Applications. Springer Science & Business Media (2007)
- [19] Hillewaert, K.: Development of the discontinuous Galerkin method for high-resolution, large scale CFD and acoustics in industrial geometries. Ph.D. thesis, Université Catholique de Louvain (2013)
- [20] Karniadakis, G., Sherwin, S.: Spectral/hp Element Methods for Computational Fluid Dynamics, second edn. Oxford Science Publications (2013)
- [21] Kennedy, C.A., Carpenter, M.H.: Diagonally implicit Runge-Kutta methods for ordinary differential equations. A review. NASA report TM-2016-219173 (2016)
- [22] Knoll, D.A., Keyes, D.E.: Jacobian-free Newton-Krylov methods: A survey of approaches and applications. Journal of Computational Physics **193**(2), 357–397 (2004)
- [23] Kravchenko, A.G., Moin, P.: Numerical studies of flow over a circular cylinder at Re<sub>D</sub>=3900. Physics of Fluids  $\mathbf{12}(2)$ , 403–417 (2000)
- [24] Mengaldo, G.: Discontinuous spectral/hp element methods: development, analysis and applications to compressible flows. Ph.D. thesis, Imperial College London (2015)
- [25] Mengaldo, G., De Grazia, D., Vincent, P.E., Sherwin, S.J.: On the connections between discontinuous Galerkin and flux reconstruction schemes: Extension to curvilinear meshes. Journal of Scientific Computing 67(3), 1272–1292 (2016)
- [26] Moxey, D., Cantwell, C.D., Bao, Y., Cassinelli, A., Castiglioni, G., Chun, S., Juda, E., Kazemi, E., Lackhove, K., Marcon, J., Mengaldo, G., Serson, D., Turner, M., Xu, H., Peiró, J., Kirby, R.M., Sherwin,

- S.J.: Nektar++: Enhancing the capability and application of high-fidelity spectral/hp element methods. Computer Physics Communications (2019)
- [27] Oliver, T.A.: Multigrid solution for high-order discontinuous Galerkin discretizations of the compressible Navier-Stokes equations. Thesis, Massachusetts Institute of Technology (2004)
- [28] Orszag, S.A.: Spectral methods for problems in complex geometries. Journal of Computational Physics **37**(1), 70–92 (1980)
- [29] Parnaudeau, P., Carlier, J., Heitz, D., Lamballais, E.: Experimental and numerical studies of the flow over a circular cylinder at Reynolds number 3900. Physics of Fluids **20**(8), 085101 (2008)
- [30] Pazner, W., Persson, P.O.: Approximate tensor-product preconditioners for very high order discontinuous Galerkin methods. Journal of Computational Physics 354, 344–369 (2018)
- [31] Peraire, J., Persson, P.: The compact discontinuous Galerkin (CDG) method for elliptic problems. SIAM Journal on Scientific Computing **30**(4), 1806–1824 (2008)
- [32] Persson, P.O., Peraire, J.: Sub-cell shock capturing for discontinuous Galerkin methods. In: 44th AIAA Aerospace Sciences Meeting and Exhibit. AIAA 2006-112, Reno, Nevada (2006)
- [33] Peterson, J.W., Lindsay, A.D., Kong, F.: Overview of the incompressible Navier–Stokes simulation capabilities in the MOOSE framework. Advances in Engineering Software 119, 68–92 (2018)
- [34] Roy, C.J.: Review of code and solution verification procedures for computational simulation. Journal of Computational Physics **205**(1), 131–156 (2005)
- [35] Saad, Y., Schultz, M.H.: GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems. SIAM Journal on Scientific and Statistical Computing 7(3), 856–869 (1986)
- [36] Toro, E.F.: Riemann Solvers and Numerical Methods for Fluid Dynamics, third edn. Springer (2009)

- [37] Toulorge, T., Desmet, W.: Optimal Runge-Kutta schemes for discontinuous Galerkin space discretizations applied to wave propagation problems. Journal of Computational Physics **231**(4), 2067–2091 (2012)
- [38] Vanden, K.J., Orkwis, P.D.: Comparison of numerical and analytical Jacobians. AIAA Journal **34**(6), 1125–1129 (1996)
- [39] Vandenhoeck, R., Lani, A.: Implicit high-order flux reconstruction solver for high-speed compressible flows. Computer Physics Communications **242**, 1–24 (2019)
- [40] Vos, P.E.J., Eskilsson, C., Bolis, A., Chun, S., Robert, M., Sherwin, S.J.: A generic framework for time-stepping partial differential equations (PDEs): General linear methods, object-oriented implementation and application to fluid problems. International Journal of Computational Fluid Dynamics 25(3), 107–125 (2011)
- [41] White, F.M.: Viscous fluid flow, second edn. McGraw-Hill, New York (1991)
- [42] Wiart, C.C.d., Hillewaert, K., Duponcheel, M., Winckelmans, G.: Assessment of a discontinuous Galerkin method for the simulation of vortical flows at high Reynolds number. International Journal for Numerical Methods in Fluids 74(7), 469–493 (2014)
- [43] Williams, S., Waterman, A., Patterson, D.: Roofline: An insightful visual performance model for multicore architectures. Commun. ACM **52**(4), 65–76 (2009)
- [44] Winters, A.R., Moura, R.C., Mengaldo, G., Gassner, G.J., Walch, S., Peiró, J., Sherwin, S.J.: A comparative study on polynomial dealiasing and split form discontinuous Galerkin schemes for under-resolved turbulence computations. Journal of Computational Physics 372, 1–21 (2018)
- [45] Xiaoquan, Y., Cheng, J., Luo, H., Zhao, Q.: Robust implicit direct discontinuous galerkin method for simulating the compressible turbulent flows. AIAA Journal **57**(3), 1113–1132 (2019)