

FINITE ELEMENT DECOMPOSITION AND MINIMAL EXTENSION FOR FLOW EQUATIONS*

R. ALTMANN¹ AND J. HEILAND²

Abstract. In the simulation of flows, the correct treatment of the pressure variable is the key to stable time-integration schemes. This paper contributes a new approach based on the theory of differential-algebraic equations. Motivated by the index reduction technique of minimal extension, a remodelling of the flow equations is proposed. It is shown how this reformulation can be realized for standard finite elements *via* a decomposition of the discrete spaces and that it ensures stable and accurate approximations. The presented decomposition preserves sparsity and does not call on variable transformations which might change the meaning of the variables. Since the method is eventually an index reduction, high index effects leading to instabilities are eliminated.

Mathematics Subject Classification. 76M10, 65L80, 65J10.

Received April 11, 2014. Revised February 23, 2015.
Published online September 9, 2015.

1. INTRODUCTION

A semi-discretization in space of the Navier-Stokes equations (NSE) leads to differential-algebraic equations (DAE) of differentiation index 2, *cf.* [42], that take the form

$$M\dot{u} + K(u) - B^T p = f, \quad u(0) = a, \quad (1.1a)$$

$$Bu = 0. \quad (1.1b)$$

For time integration one has to take care of the differential-algebraic structure that requires implicit schemes and that can cause a reduction of the convergence order up to possible divergence [16, 24]. To avoid divergence for low-order schemes, a general approach is the remodelling of the equations as an equivalent or arbitrarily close system of index 1 [16, 24, 42].

Keywords and phrases. Navier–Stokes equations, time integration schemes, finite element method, index reduction, operator DAE.

* *The work of the first author was supported by the ERC Advanced Grant “Modeling, Simulation and Control of Multi-Physics Systems” MODSIMCONMP and the Berlin Mathematical School BMS. The work of J. Heiland was funded by a scholarship by “Studienstiftung des deutschen Volkes”.*

¹ Institut für Mathematik MA4-5, Technische Universität Berlin, Straße des 17. Juni 136, 10623 Berlin, Germany.
raltmann@math.tu-berlin.de

² Max Planck Institute for Dynamics of Complex Technical Systems, Sandtorstraße 1, 39106 Magdeburg, Germany.

For the semi-discrete NSE a variety of methods has been developed that are successfully applied in the numerical approximation of unsteady flow. To illustrate the basic ideas and properties, we roughly classify the most common approaches into penalty methods [20, 35], pressure correction or projection methods [16], and divergence-free methods.

In penalization methods, one adds a term $\lambda^{-1}p$, $\lambda \gg 1$, to the left-hand side of (1.1b) and obtains an ODE for u via $M\dot{u} + K(u) - \lambda B^T B u = f$. In projection methods, also referred to as operator splitting or pressure correction methods, one uses a guess for the pressure to compute an approximate velocity update \tilde{u} via (1.1a) in every time step. Then, one computes the components $\tilde{u} = u_0 + u_\perp$, with u_0 satisfying $Bu_0 = 0$ and u_\perp being in the span of $M^{-1}B^T$.

The two mentioned approaches decouple the pressure and velocity computation. This is computationally beneficial since (1.1) decomposes into two smaller systems. As elaborated in [42], this decoupling is incomplete and depends strongly on the heuristic penalization parameter, cf. [35], or the time step. Also, as we will reason in Chapter 3.1, the computation of the pressure from the velocities is ill-conditioned. Another common difficulty of projection schemes is the need for boundary conditions for the pressure which are unphysical [16].

A complete decoupling is obtained in divergence-free formulations that reduce (1.1a) to an ODE for the divergence-free components of u . The presence of a divergence-free basis for u is optimal for the approximation of the velocity since the system is reduced to a subspace of the velocity space and the constraints (1.1b) are fulfilled *a priori*. An overview of divergence-free elements is provided in ([16], Chap. 3.13.7). However, these elements are rarely used in simulations because of their difficult implementation ([16], Chap. 3.12.2).

As an alternative, one may resolve the algebraic constraints numerically, *e.g.* via a QR-decomposition of B . This approach is not taken in practice, because the variable transformation $u \leftarrow Q\tilde{u}$ is computationally unfeasible already for moderately sized systems. And again, and this holds also for divergence-free elements, the associated equation for the pressure is ill-conditioned, as demonstrated in Section 3.1 below. Recent approaches [25] for the numerical construction of sparse divergence-free bases, *i.e.*, a null space of B , only tackle the problem of infeasibility.

Note that the use of quasi divergence-free elements, see *e.g.* [29], reduces the system's size but leaves the DAE structure unchanged.

There are plenty of other sophisticated methods in computational fluid dynamics (CFD) that can cope with the mentioned above difficulties in the stable and consistent approximation of flows, see, *e.g.*, the review article [33] or the DAE-based approach proposed in [26].

We present a new modelling approach which is appealing because of several reasons. In contrast to established methods it neither depends on a heuristic parameter nor does it require boundary conditions for the pressure. It is consistent, *i.e.* the solution set remains unaltered, and it has a valid representation on the operator level. Additionally, this technique allows for more general constraints of the form $Bu = g \neq 0$. Such a constraint appears because of the incorporation of boundary conditions, or, *e.g.*, in the dual equations of optimal control problems where the pressure is included in the cost functional [21].

Basically, we propose a variant of minimal extension [23, 24] tailored to finite element discretizations of flow equations. As for the analysis of the abstract setting, we adapt the ideas of [1] where problems from elastodynamics were considered. In practice, the crucial part in minimal extensions is the right choice of what has to be added to the equations. For mechanical systems, the needed dummy variables are easily determined [23]. For flow equations, the extension is readily determined only in theory. Considering particular but popular discretization schemes, we provide algorithms that make the assembling of a minimally extended system feasible.

This new approach can be seen as resolving the algebraic constraints while at the same time avoiding the difficulties mentioned above. The variables are transformed only *via* a permutation and thus, keep their physical meaning. The sole application of a permutation preserves sparsity and is well conditioned. Since the so-called hidden constraint $B\dot{u} = 0$ is added to the system instead of implicitly eliminated, instabilities are reduced. In particular, we will show that the method is robust with respect to an error due to the approximate solution of the algebraic equations.

Furthermore, the pressure p remains a physically valid part of the system, rather than being eliminated or functioning as a velocity correction. The increase of the system size may be compensated by the direct applicability of efficient time stepping schemes.

This paper is organized as follows. In Section 2, the unsteady NSE is formulated as a constrained operator differential equation. Following the ideas of minimal extension, we reformulate the so-called operator DAE such that a spatial discretization leads to a DAE of differentiation index 1. This property requires certain assumptions on the finite element spaces which are presented in Section 3. In particular, a splitting of the velocity ansatz space is necessary. We show the advantages of this method and give examples for such splittings for standard discretization schemes such as Crouzeix–Raviart [13] and Taylor–Hood [37]. In Section 4, we present the benefits of the presented approach for numerical time integration for a non-viscous two-dimensional internal flow and for the two-dimensional cylinder wake.

2. OPERATOR FORMULATION

We consider the unsteady NSE on a domain $\Omega \subset \mathbb{R}^n$, $n \in \{2, 3\}$ with twice differentiable boundary $\partial\Omega$ in a time interval $(0, T)$,

$$\dot{u} + (u \cdot \nabla)u - \frac{1}{\text{Re}} \Delta u + \nabla p = \beta \quad \text{in } \Omega \times (0, T), \quad (2.1a)$$

$$\text{div } u = 0 \quad \text{in } \Omega \times (0, T), \quad (2.1b)$$

$$u = 0 \quad \text{on } \partial\Omega \times (0, T), \quad (2.1c)$$

$$u(\cdot, 0) = a. \quad (2.1d)$$

This system describes the evolution of the velocity field $u(t) \in (\Omega \rightarrow \mathbb{R}^n)$ and the pressure $p(t) \in (\Omega \rightarrow \mathbb{R})$ for a given parameter $\text{Re} > 0$, an initial value $a \in (\Omega \rightarrow \mathbb{R}^n)$ and a volume force $\beta(t) \in (\Omega \rightarrow \mathbb{R}^n)$.

2.1. Preliminaries

For the basic definition of Sobolev spaces on a domain Ω , as the space $H^1(\Omega)$ of square integrable functions $L^2(\Omega)$ that possess a weak derivative in $L^2(\Omega)$, its subspace $H_0^1(\Omega)$ of functions that are weakly differentiable and vanish on the boundary $\partial\Omega$, and of Bochner spaces like $L^2(0, T; L^2(\Omega))$ or $H^1(0, T; L^2(\Omega))$, we refer the reader to [36].

To shorten notation, we define the spaces

$$\mathcal{V} := [H_0^1(\Omega)]^n, \quad \mathcal{H} := [L^2(\Omega)]^n, \quad \text{and } \mathcal{Q} := L^2(\Omega)/\mathbb{R}.$$

The space \mathcal{V} is densely and continuously embedded in \mathcal{H} and thus, the identification of \mathcal{H} with its dual \mathcal{H}' via the Riesz isomorphism gives the evolution triple $\mathcal{V} \subset \mathcal{H} \subset \mathcal{V}'$. Let $\mathcal{W}^{1,2}(0, T)$ be the space of functions $u \in L^2(0, T; \mathcal{V})$ with weak time derivative $\dot{u} \in L^1(0, T; \mathcal{V}')$.

We consider a weak formulation of (2.1): given right-hand sides $\mathcal{F} \in L^2(0, T; \mathcal{V}')$, $\mathcal{G} \in L^2(0, T; \mathcal{Q}')$, and an initial condition $a \in \mathcal{H}$, we seek for $(u, p) \in \mathcal{W}^{1,2}(0, T) \times L^2(0, T; \mathcal{Q})$ satisfying

$$\dot{u}(t) + \mathcal{K}(u(t)) - \mathcal{B}'p(t) = \mathcal{F}(t) \quad \text{in } \mathcal{V}', \quad (2.2a)$$

$$\mathcal{B}u(t) = \mathcal{G}(t) \quad \text{in } \mathcal{Q}', \quad (2.2b)$$

$$u(0) = a \quad \text{in } \mathcal{H}, \quad (2.2c)$$

a.e. on $(0, T)$. Because of the differential-algebraic structure in an abstract setting, we call (2.2) an *operator DAE*. Therein, the operators $\mathcal{K}: \mathcal{V} \rightarrow \mathcal{V}'$ and $\mathcal{B}: \mathcal{V} \rightarrow \mathcal{Q}'$ are defined via

$$\langle \mathcal{K}(u), v \rangle = \int_{\Omega} (u \cdot \nabla)u \cdot v \, dx + \frac{1}{\text{Re}} \int_{\Omega} \nabla u \cdot \nabla v \, dx \quad (2.3)$$

and

$$\langle \mathcal{B}u, q \rangle = \int_{\Omega} (\operatorname{div} u)q \, dx = \langle u, \mathcal{B}'q \rangle, \tag{2.4}$$

respectively, given $u \in \mathcal{V}$ and for all $v \in \mathcal{V}$ and $q \in \mathcal{Q}$. Note that system (2.2) not only covers the NSE but also more general flow equations since we have introduced an inhomogeneity \mathcal{G} .

Since \mathcal{B} is bounded, the ansatz space \mathcal{V} can be decomposed into the divergence-free space \mathcal{V}_{df} and its orthogonal complement $\mathcal{V}_{\text{df}}^{\perp}$ with respect to the inner product of \mathcal{V} , *i.e.*,

$$\mathcal{V}_{\text{df}} := \ker \mathcal{B} = \{u \in \mathcal{V} \mid \operatorname{div} u = 0\}, \quad \mathcal{V} = \mathcal{V}_{\text{df}} \oplus \mathcal{V}_{\text{df}}^{\perp}. \tag{2.5}$$

This implies a unique decomposition of $u \in \mathcal{V}$ into $u = u_1 + u_2$ with $u_1 \in \mathcal{V}_{\text{df}}$ and $u_2 \in \mathcal{V}_{\text{df}}^{\perp}$.

2.2. Existence of solutions

Classical existence results consider (2.2) with $\mathcal{G}(t) = 0$ formulated on the subspace of divergence-free functions $\mathcal{V}_{\text{df}} \subset \mathcal{V}$. The problem then turns to: find $u_1 \in L^2(0, T; \mathcal{V}_{\text{df}})$ with $\dot{u}_1 \in L^1(0, T; \mathcal{V}'_{\text{df}})$ satisfying

$$\dot{u}_1(t) + \mathcal{K}_1(u_1(t)) = \mathcal{F}_1(t) \quad \text{in } \mathcal{V}'_{\text{df}}, \tag{2.6a}$$

$$u_1(0) = a_1 \quad \text{in } \mathcal{H}, \tag{2.6b}$$

a.e. on $(0, T)$. Therein, let $\mathcal{F}_1 \in L^2(0, T; \mathcal{V}'_{\text{df}})$, a_1 be in the closure of \mathcal{V}_{df} w.r.t. the norm of \mathcal{H} , and $\mathcal{K}_1: \mathcal{V}_{\text{df}} \rightarrow \mathcal{V}'_{\text{df}}$ be defined as in (2.3). The formulation *via* divergence-free functions particularly eliminates the pressure from the equations.

There exists a solution $u_1 \in L^2(0, T; \mathcal{V}_{\text{df}})$ satisfying (2.6), see ([38], Thm. III.3.1), which is unique in the two-dimensional case ([38], Thm. III.3.2). Given u , one can generally establish a corresponding pressure p as a distribution on $(0, T) \times \Omega$, (*cf.* [38], p. 307). However, the pair (u_1, p) only solves (2.2) under additional regularity conditions, *cf.* [32], and if $a = a_1$. In particular, if the values in (2.2a) are in \mathcal{H} rather than in \mathcal{V}' , then system (2.2) can be split *via* the Helmholtz decomposition ([15], Cor. I.3.4) into a part defining $u_1 \in L^2(0, T; \mathcal{V}_{\text{df}})$ and a part that uniquely defines $p \in L^2(0, T; Q)$. This additional regularity is given globally in 2D and locally in time in 3D, if $\mathcal{F} \in L^2(0, T; \mathcal{H})$ and $a \in \mathcal{V}_{\text{df}}$, (*cf.* [36], Lems. 25.1, 25.2). Since a solution u_1 to (2.2) always solves (2.6), in 2D, it is unique.

An inhomogeneity \mathcal{G} in the constraint (2.2b) is likely to appear in discretized schemes and for more general boundary conditions. For maximal generality, we will consider it present, as it imposes restrictions on the solvability of the equations.

Since $\mathcal{B}: \mathcal{V} \rightarrow \mathcal{Q}'$ has a right-inverse \mathcal{B}^- , see ([15], Lem. I.4.1), the complement to u_1 is eventually given by $u_2 = \mathcal{B}^- \mathcal{G}$. Plugging this relation into (2.2), we obtain the remainder system

$$\dot{u}_1(t) + \mathcal{K}(u_1(t) + \mathcal{B}^- \mathcal{G}(t)) - \mathcal{B}'p(t) = \mathcal{F}(t) - \mathcal{B}^- \dot{\mathcal{G}}(t) \quad \text{in } \mathcal{V}', \tag{2.7a}$$

$$\mathcal{B}u_1(t) = 0 \quad \text{in } \mathcal{Q}', \tag{2.7b}$$

$$u_1(0) = a - \mathcal{B}^- \mathcal{G}(0) \quad \text{in } \mathcal{H}, \tag{2.7c}$$

which is well-posed, only if $\dot{\mathcal{G}}$ is at least in $L^1(0, T; \mathcal{Q}')$. Then, solvability for (2.7) can be established analogously to solvability for (2.2).

2.3. Index reduction

The spatial discretization of (2.2) leads to a DAE of index 2. Here we use the concept of the perturbation index [17], that for semi-explicit systems as (2.2) and (2.8) coincides with the differentiation index [11]. Thus, a system is of index d if d is the smallest integer such that a perturbation of the right hand side δ causes a deviation in the solution that can be bounded *via* the first $d - 1$ time derivatives of δ , (*cf.* [17], Def. 1.1).

In numerical simulations, the occurrence of derivatives of perturbations appears as divisions by powers of the small discretization parameters ([17], p. 1).

Thus, it may be preferable to use equivalent formulations of lower index. The semi-explicit structure of the NSE allows for minimal extension [23], which reduces the index without transforming the variables. This is done by adding the time derivatives of the constraints, which leads to an overdetermined system, and then introducing a minimal number of variables to make the system square again.

Following this idea, we reformulate the operator DAE (2.2) to index-1 form. By this we mean that a certain discretization in space leads to a DAE of index 1. Using (2.5), we seek in system (2.2) for u_1 and u_2 instead of u . The corresponding ansatz spaces read

$$u_1 \in \mathcal{W}^{1,2}(0, T) \cap L^2(0, T; \mathcal{V}_{\text{df}}), \quad u_2 \in \mathcal{W}^{1,2}(0, T) \cap L^2(0, T; \mathcal{V}_{\text{df}}^\perp).$$

Assuming sufficient regularity, we add the derivative of the constraint, the so-called *hidden constraint*. Since the operator \mathcal{B} is independent of time, the hidden constraint reads

$$\mathcal{B}\dot{u}_2(t) = \mathcal{B}[\dot{u}_1(t) + \dot{u}_2(t)] = \dot{\mathcal{G}}(t).$$

As a second step, we introduce a new variable $\tilde{u}_2 := \dot{u}_2$. The reformulated and extended problem then reads: find $u_1 \in \mathcal{W}^{1,2} \cap L^2(0, T; \mathcal{V}_{\text{df}})$, $u_2, \tilde{u}_2 \in L^2(0, T; \mathcal{V}_{\text{df}}^\perp)$, and $p \in L^2(0, T; \mathcal{Q})$ such that

$$\dot{u}_1(t) + \tilde{u}_2(t) + \mathcal{K}(u_1(t) + u_2(t)) - \mathcal{B}'p(t) = \mathcal{F}(t) \quad \text{in } \mathcal{V}', \quad (2.8a)$$

$$\mathcal{B}u_2(t) = \mathcal{G}(t) \quad \text{in } \mathcal{Q}', \quad (2.8b)$$

$$\mathcal{B}\tilde{u}_2(t) = \dot{\mathcal{G}}(t) \quad \text{in } \mathcal{Q}', \quad (2.8c)$$

$$u_1(0) = a_1 \quad \text{in } \mathcal{H}. \quad (2.8d)$$

Remark 2.1. Since \mathcal{B} is constant in time, equations (2.8b) and (2.8c) imply that u_2 is in $H^1(0, T; \mathcal{V}_{\text{df}}^\perp)$ provided that $\mathcal{G} \in H^1(0, T; \mathcal{Q}')$.

Remark 2.2. For the spatial discretization in Section 3, we add to (2.8b) and to (2.8c) the vanishing terms $\mathcal{B}u_1(t)$ and $\mathcal{B}\dot{u}_1(t)$, respectively. This is necessary since we will deal with nonconforming finite elements, where the discrete version of u_1 does not vanish under the action of \mathcal{B} .

The derivation of system (2.8) gives rise to the following theorem.

Theorem 2.3. *If $\mathcal{G} \in H^1(0, T; \mathcal{Q}')$ and $a = a_1 + \mathcal{B}^- \mathcal{G}(0)$, then the operator DAEs (2.2) and (2.8) have the same solution set.*

For the proof that the extended operator DAE (2.8) leads to an index-1 DAE, we refer to Section 3.2.

Remark 2.4. For completeness, we want to mention that a straight forward index reduction can be obtained by adding $\lambda \mathcal{B}\dot{u} - \lambda \dot{\mathcal{G}}$ to (2.2b), which is known as *Baumgarte stabilization*. We will not consider this method here because of its strong dependence on the parameter λ , cf. [5, 30].

3. DISCRETE FORMULATION

This section is devoted to the spatial and temporal discretization of the NSE in two space dimensions. We show that the introduced splitting of the ansatz spaces provides an efficient simulation procedure. Furthermore, we comment on possible extensions to the three-dimensional case.

We consider spatial discretizations by finite elements, *i.e.*, we construct finite dimensional subspaces V_h and Q_h of \mathcal{V} and \mathcal{Q} , respectively, based on a triangulation \mathcal{T} of the polygonal Lipschitz domain Ω . The triangulation is assumed to be regular in the sense of Ciarlet [12]. Furthermore, we take for granted that the

triangulation is shape regular ([8], Chap. II.5). In the sequel, \mathcal{N} denotes the set of vertices of \mathcal{T} and \mathcal{E} the set of edges. The latter consists of interior and boundary edges, namely \mathcal{E}_{int} and \mathcal{E}_{ext} . We focus on triangular meshes but will comment on quadrilateral elements in Section 3.6.

The finite dimensional approximation of the velocity $u(t)$ is given by the coefficient vector $q(t) \in \mathbb{R}^n$, which corresponds to a function in V_h . The discrete representative of the pressure is again denoted by $p(t) \in \mathbb{R}^m$. The semi-discretized version of system (2.2) reads

$$M\dot{q}(t) + K(q(t)) - B^T p(t) = f(t), \tag{3.1a}$$

$$Bq(t) = g(t). \tag{3.1b}$$

Therein, for a given basis $\{\Psi_j\}$ of V_h and $\{\varphi_i\}$ of Q_h ,

$$M = [m_{jk}] \in \mathbb{R}^{n \times n}, \quad m_{jk} := \int_{\Omega} \Psi_j \cdot \Psi_k \, dx \tag{3.2a}$$

denotes the positive definite mass matrix and the nonlinear function K is the discrete version of the operator \mathcal{K} ,

$$K(q(t)) = [K_j(q(t))], \quad K_j(q(t)) := \int_{\Omega} (q(t) \cdot \nabla)q(t) \cdot \Psi_j \, dx + \frac{1}{\text{Re}} \int_{\Omega} \nabla q(t) \cdot \nabla \Psi_j \, dx, \tag{3.2b}$$

where we have assigned $q(t)$ with its function representation in V_h . The matrix $B = [b_{ij}] \in \mathbb{R}^{m \times n}$ is defined via

$$b_{ij} = \int_{\Omega} \varphi_i \, \text{div} \Psi_j \, dx. \tag{3.2c}$$

In the next subsection, we recall solution strategies of solving system (3.1) with the help of the QR algorithm and divergence-free finite elements. Afterwards, we propose a different ansatz which is based on the index-1 formulation which arises from the discretization of the operator DAE (2.8). This includes a decomposition of the finite element space V_h .

3.1. QR Decomposition and divergence-free elements

For completeness, we address the case of eliminating all algebraic constraints and reducing the system to the so-called *inherent* or *underlying ODE*. In other words, we consider here the index-0 formulation of the NSE.

Numerically, this can be achieved by a QR decomposition of $B = [0 \ R]Q$, with R invertible and Q unitarian. With the transformation $q =: Q^T \tilde{q}$ and the splitting $\tilde{q} = [\tilde{q}_1^T \ \tilde{q}_2^T]^T$, the divergence constraint (3.1b) becomes $\tilde{q}_2 = R^{-1}g$. Then, a scaling of the momentum equation (3.1a) by Q^T gives the decoupled system

$$\tilde{M}_{11} \dot{\tilde{q}}_1 + \tilde{K}_{11}(\tilde{q}_1) = \tilde{f}_1, \tag{3.3a}$$

$$-R^T p = -\tilde{M}_{21} \dot{\tilde{q}}_1 - \tilde{K}_{21}(\tilde{q}_1) + \tilde{f}_2. \tag{3.3b}$$

The subscripts refer to the block structure corresponding to the splitting of \tilde{q} and the tilde denotes the coefficients after the transformation of the system and the substitution of \tilde{q}_2 by $R^{-1}g$.

Since \tilde{M}_{11} is invertible, (3.3a) is equivalent to a standard ODE for \tilde{q}_1 . Thus, one can expect stable approximations of $q = Q^T \tilde{q}$. However, the pressure p as defined by (3.3b) requires $\dot{\tilde{q}}_1$ and $\tilde{K}_{21}(\tilde{q}_1)$, *i.e.*, discrete time and space derivatives. In a numerical realization, this amplifies a non-smooth error in \tilde{q}_1 by τ^{-1} or h^{-2} , where τ and h are length scales of time and space discretization, respectively.

Any such decomposition would suffer from these instabilities in the pressure approximation for τ and h tending to zero. In particular, divergence-free elements that directly provide a basis for \tilde{q}_1 , thus, come with the same difficulties for the pressure reconstruction.

3.2. Index-1 formulation

As announced above, in this subsection, we show that a proper semi-discretization in space of system (2.8) leads to a DAE of index 1. Since \mathcal{V} was decomposed in Section 2.3 into $\mathcal{V}_{\text{df}} \oplus \mathcal{V}_{\text{df}}^\perp$, we also decompose the finite dimensional space V_h . We denote the approximation space of \mathcal{V}_{df} by $V_{h,1}$, its complement $\mathcal{V}_{\text{df}}^\perp$ is discretized by $V_{h,2}$. Furthermore, we assume that the direct sum of $V_{h,1}$ and $V_{h,2}$ is again V_h . Note that we do not assume the discretization to be conform, *i.e.*, we allow for $V_{h,1} \not\subset \mathcal{V}_{\text{df}}$ and $V_{h,2} \not\subset \mathcal{V}_{\text{df}}^\perp$, even if $V_h \subset \mathcal{V}$.

With q_1 , q_2 , and \tilde{q}_2 denoting the semi-discrete approximations of u_1 , u_2 , and \tilde{u}_2 , respectively, we obtain the spatial discretization of system (2.8) which reads

$$M \begin{bmatrix} \dot{q}_1(t) \\ \tilde{q}_2(t) \end{bmatrix} + K \begin{bmatrix} q_1(t) \\ q_2(t) \end{bmatrix} - B^T p(t) = f(t), \quad (3.4a)$$

$$B \begin{bmatrix} q_1(t) \\ q_2(t) \end{bmatrix} = g(t), \quad (3.4b)$$

$$B \begin{bmatrix} \dot{q}_1(t) \\ \tilde{q}_2(t) \end{bmatrix} = \dot{g}(t) \quad (3.4c)$$

with M , K , and B as defined in (3.2) and with the basis of V_h ordered according to its decomposition into $V_{h,1}$ and $V_{h,2}$. In the sequel, we analyse for which discretizations the DAE (3.4) is of index 1.

We require the standard stability condition for the spatial discretization ([15], Chap. II), *i.e.*, there exists a positive constant $c \in \mathbb{R}$ such that

$$\inf_{q_h \in Q_h} \sup_{v_h \in V_h} \frac{\langle Bv_h, q_h \rangle}{\|v_h\|_{\mathcal{V}} \|q_h\|_{\mathcal{Q}}} \geq c > 0. \quad (3.5)$$

From (3.5) we infer that B has full row rank and that there is a decomposition $V_h = V_{h,1} \oplus V_{h,2}$ such that the submatrix of the columns accounting for $V_{h,2}$ is invertible. Formally, we put this into the following assumption.

Assumption 3.1. The finite element spaces $V_{h,1}$, $V_{h,2}$, and Q_h satisfy that the matrix representation B as defined in (3.2c) has the block structure $B = [B_1 \ B_2]$ with a nonsingular square matrix B_2 that contains the columns corresponding to $V_{h,2}$.

As a direct consequence of Assumption 3.1, we have that $\dim V_{h,2} = \dim Q_h$. In Section 3.5 we give examples how to decompose V_h for certain finite element spaces used in CFD to meet Assumption 3.1.

Theorem 3.2. *Every finite element discretization of (2.8) with spaces $V_{h,1}$, $V_{h,2}$, and Q_h satisfying Assumption 3.1 leads to a DAE of index 1.*

Proof. We follow the proof of ([24], Thm. 6.12) and show that under Assumption 3.1 the DAE (3.4) has index 1. A multiplication of (3.4a) by BM^{-1} from the left and the relation (3.4c) give

$$-BM^{-1}B^T p = BM^{-1}f - BM^{-1}K \begin{bmatrix} q_1 \\ q_2 \end{bmatrix} - \dot{g}. \quad (3.6)$$

Since M is positive definite and B is of full rank by Assumption 3.1, $BM^{-1}B^T$ is invertible. Thus, we can express the pressure p in terms of f , \dot{g} , q_1 , and q_2 . By Assumption 3.1 and (3.4b), we have

$$q_2 = B_2^{-1}g - B_2^{-1}B_1q_1. \quad (3.7)$$

Finally, if we insert (3.6) and (3.7) into equation (3.4a), we obtain

$$M \begin{bmatrix} \dot{q}_1 \\ \tilde{q}_2 \end{bmatrix} = f + B^T p - K \begin{bmatrix} q_1 \\ q_2 \end{bmatrix} =: f^*(f, g, \dot{g}, q_1).$$

Since M is invertible, this provides an ODE in q_1 . Thus, we can solve for q_1 , q_2 (by (3.7)), \tilde{q}_2 (by (3.4c)), and p (by (3.6)), *i.e.* the DAE (3.4) is of index 1. \square

Remark 3.3. The *inf-sup condition* (3.5) ensures a bound on the inverse of $BM^{-1}B^T$ from equation (3.6) independent of the discretization parameter h and thus, stability in the spatial approximation of the pressure.

Remark 3.4. The reordering of the basis of V_h , that ensures Assumption 3.1, always exists – if (3.5) holds – and is basically a permutation of the velocity variables. This time-independent transformation is applied and inverted in exact arithmetics and preserves sparsity of the coefficient matrices.

3.3. Time integration

The spatially discretized NSE (3.1) represents a semi-explicit index-2 DAE. For these systems, implicit time-stepping schemes such as the *Radau Ila* or *backward differencing* methods provide stable approximations of arbitrary order, provided that the inhomogeneities are sufficiently smooth [18]. These methods, however, require the solution of the full coupled nonlinear system at every stage. A compromise of the stability of implicit and the low computational load of explicit schemes is given by half-explicit schemes, that are explicit in the dynamic equations and implicit in the algebraic part.

Half-explicit Runge–Kutta (RK) methods were investigated for index-1 DAEs in [4, 27]. Methods for the index-2 case are provided *e.g.* in [3, 18]. Generally speaking, the application to index-1 problems is straightforward while index-2 problems require specific treatments and possibly additional stages in the RK method.

We illustrate the different behavior with respect to inaccuracies of the index-1 and index-2 formulation of the NSE, using an explicit Euler method ([16], Chap. 3.16.1) for the dynamical part. Superscripts $+$, c , and $-$ denote the next, current, and previous iterates, respectively. For the index-2 equation (3.1), the update to (q^+, p^c) from the current iterate (q^c, p^-) *via* a time step of length τ is obtained *via*

$$\begin{bmatrix} \frac{1}{\tau}M & -B^T \\ B & 0 \end{bmatrix} \begin{bmatrix} q^+ \\ p^c \end{bmatrix} = \begin{bmatrix} \frac{1}{\tau}Mq^c + f^c - K(q^c) \\ g^+ \end{bmatrix}. \tag{3.8}$$

For the update of the index-1 formulation (3.4), we propose the solution of

$$\begin{bmatrix} \frac{1}{\tau}M_{11} & M_{12} & -B_1^T & 0 \\ \frac{1}{\tau}M_{21} & M_{22} & -B_2^T & 0 \\ \frac{1}{\tau}B_1 & B_2 & 0 & 0 \\ B_1 & 0 & 0 & B_2 \end{bmatrix} \begin{bmatrix} q_1^+ \\ \tilde{q}_2^c \\ p^c \\ q_2^+ \end{bmatrix} = \begin{bmatrix} \frac{1}{\tau}M_{11}q_1^c + f_1^c - K_1(q_1^c, q_2^c) \\ \frac{1}{\tau}M_{21}q_1^c + f_2^c - K_2(q_1^c, q_2^c) \\ \frac{1}{\tau}B_1q_1^c + \dot{g}^c \\ g^+ \end{bmatrix}. \tag{3.9}$$

The different stability properties become evident, if one examines the inherent equation for the pressure update p^c , derived *via* premultiplying the upper part of the equations by BM^{-1} . In the index-2 case (3.8), this leads to

$$-BM^{-1}B^T p^c = \frac{Bq^c - Bq^+}{\tau} + BM^{-1}[f^c - K(q^c)]. \tag{3.10}$$

The index-1 formulation yields for the pressure

$$-BM^{-1}B^T p^c = \frac{1}{\tau}B \begin{bmatrix} q_1^c - q_1^+ \\ -\tau\tilde{q}_2^c \end{bmatrix} + BM^{-1}[f^c - K(q_1^c, q_2^c)]. \tag{3.11}$$

If the equations are solved up to a residual of size ε^c , the dominating difference in the pressure definition in the index-1 and index-2 formulation lies in the terms

$$\frac{Bq^c - Bq^+}{\tau} = \frac{g^c - g^+}{\tau} + \frac{\varepsilon^c - \varepsilon^+}{\tau} \quad \text{and} \quad \frac{1}{\tau}B \begin{bmatrix} q_1^c - q_1^+ \\ -\tau\tilde{q}_2^c \end{bmatrix} = -\dot{g}^+ + \varepsilon^c. \tag{3.12}$$

Thus, unlike for the index-1 case, in the index-2 formulation an error in the algebraic constraints is amplified by $1/\tau$. This instability is observed in the numerical example in Section 4.

3.4. Stable discretization schemes

In this section we summarize the most common finite element schemes used in CFD. All mentioned schemes satisfy the *inf-sup* (also called Ladyzhenskaya–Babuška–Brezzi) condition (3.5) which is necessary to ensure stability of the pressure variable ([10], Chap. VI.3). Additional stable schemes are addressed in ([15], Chap. II) as well as in ([16], Chap. 3).

Using standard notation, we denote by $\mathcal{P}_k(\mathcal{T})$ the space of piecewise polynomials of degree k . The space of piecewise polynomials which are globally continuous is denoted by

$$\mathcal{S}_k(\mathcal{T}) := \mathcal{P}_k(\mathcal{T}) \cap H^1(\Omega).$$

With zero boundary conditions, we write $\mathcal{S}_{k,0}(\mathcal{T})$. For the pressure variable, we introduce the space

$$\mathcal{P}_0^0(\mathcal{T}) := \mathcal{P}_0(\mathcal{T})/\mathbb{R} = \mathcal{P}_0(\mathcal{T} \setminus \{T_0\}) \quad (3.13)$$

for some triangle $T_0 \in \mathcal{T}$, *i.e.*, we *fix* the pressure by setting it to zero at one triangle T_0 . Similarly, we define $\mathcal{S}_1^0(\mathcal{T}) := \mathcal{S}_1(\mathcal{T})/\mathbb{R}$. The discontinuous Crouzeix–Raviart finite element space [13] with zero boundary conditions is given by

$$\text{CR}_0(\mathcal{T}) := \mathcal{P}_1(\mathcal{T}) \cap C(\{\text{mid}(E) \mid E \in \mathcal{E}\}) \cap \{v \mid v(\text{mid}(\mathcal{E}_{\text{ext}})) = 0\}.$$

This space contains piecewise affine functions which are continuous at the midpoints of interior edges and vanish at the midpoints of boundary edges. It is well-known that the $[\mathcal{S}_{1,0}(\mathcal{T})]^2 - \mathcal{P}_0^0(\mathcal{T})$ scheme is not stable ([10], Ex. VI.3.1). An alternative low-order scheme was introduced in [13] and is given by

$$V_h = [\text{CR}_0(\mathcal{T})]^2, \quad Q_h = \mathcal{P}_0^0(\mathcal{T}). \quad (3.14)$$

In [22] yet another finite element space is introduced with less degrees of freedom by the mixture of continuous and discontinuous velocity components. An alternative approach is to enrich the discrete velocity space $\mathcal{S}_{1,0}(\mathcal{T})$ by bubble functions [41]. Bernardi and Raugel [7] use edge-bubble functions multiplied by the outer normal vector of the corresponding edge. Such an edge-bubble function is defined as scaled product of the two nodal hat functions corresponding to the two nodes of an edge. Thus, its support is locally bounded by the two adjacent triangles. In this way, the fluxes through interior edges provide additional degrees of freedom. This ansatz is analysed in more detail in Section 3.5.2.

Quite popular are approaches of Taylor–Hood type [37]. Therein, the velocities are approximated by polynomials of one degree higher than the pressure. The Taylor–Hood element of lowest order is defined by

$$V_h = [\mathcal{S}_{2,0}(\mathcal{T})]^2, \quad Q_h = \mathcal{S}_1^0(\mathcal{T}). \quad (3.15)$$

Note that the ansatz for the pressure is continuous which yields a more natural model.

3.5. Decompositions of \mathbf{V}_h

In this subsection, we derive decompositions of the finite elements schemes mentioned above such that Assumption 3.1 is satisfied. Since we do not deal with divergence-free elements, all resulting discretizations schemes for system (3.4) will be of nonconforming nature. Nonconforming finite element methods, for which the discrete space is no subspace of the continuous ansatz space, are analysed in ([9], Chap. 10).

We show the construction of $V_{h,1}$ and $V_{h,2}$ by means of three examples.

3.5.1. Discontinuous velocity

As first example, we consider the discontinuous Crouzeix–Raviart ansatz, introduced in (3.14). This ansatz is often used since it provides an efficient tool for CFD [6]. A proof of the *inf-sup* condition (3.5) is given in [6, 13].

Let $T_0 \in \mathcal{T}$ denote the triangle on which the pressure is fixed. The following algorithm defines a one-to-one mapping $\iota : \mathcal{T} \setminus \{T_0\} \rightarrow \mathcal{E}_{\text{int}}$ which allows to define the discrete space $V_{h,2}$.

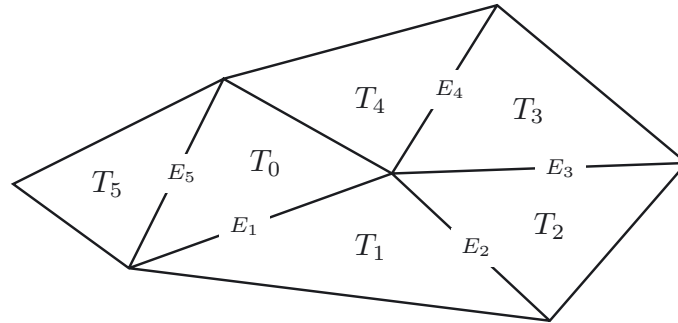


FIGURE 1. Illustration of Algorithm 3.5, $\iota(T_i) = E_i$ for $i = 1, \dots, 5$. Step 2b of the algorithm is applied once to reset $T := T_0$.

Algorithm 3.5 (Mapping ι).

Step 1. Choose any $T \in \mathcal{T} \setminus \{T_0\}$ which shares an edge with T_0 and denote this edge by $E := T_0 \cap T \in \mathcal{E}_{\text{int}}$. Then, define $\iota(T) := E$ and $\mathcal{T}_R := \mathcal{T} \setminus \{T_0, T\}$. If $\mathcal{T}_R = \emptyset$, then stop.

Step 2. If T from the previous step has an edge-neighbour in \mathcal{T}_R , then continue with Step 2a. Otherwise, go to Step 2b.

Step 2a. Select such a neighbouring triangle $S \in \mathcal{T}_R$ and set $E := T \cap S \in \mathcal{E}_{\text{int}}$. Furthermore, set $\iota(S) := E$ and $\mathcal{T}_R := \mathcal{T}_R \setminus \{S\}$. If $\mathcal{T}_R = \emptyset$, then stop. Otherwise, return to Step 2 with $T := S$.

Step 2b: Reset $T \in \mathcal{T} \setminus \mathcal{T}_R$ such that there exists an edge-neighbour in \mathcal{T}_R and return to Step 2.

An illustration of the algorithm is shown in Figure 1.

Remark 3.6. Step 2b of Algorithm 3.5 is realizable since $\mathcal{T}_R \neq \emptyset$ and Ω is assumed to be connected with Lipschitz boundary. Furthermore, the algorithm terminates in finite time since the number of triangles is finite and Step 2a reduces the set of triangles \mathcal{T}_R by one in at least every second iteration.

Remark 3.7. In the sequel, we will benefit of an order of the triangles $\mathcal{T} \setminus \{T_0\}$, given by their first appearance in Algorithm 3.5, namely $\{T_j\}_{j=1, \dots, |\mathcal{T}|-1}$.

Let ϕ_E denote the Crouzeix–Raviart basis function for an edge $E \in \text{Range}(\iota) \subset \mathcal{E}_{\text{int}}$, i.e., ϕ_E is piecewise linear with the value 1 at the midpoint of E and 0 at the midpoint of any other edge. The corresponding triangle $T = \iota^{-1}(E)$ lies in the support of ϕ_E and thus, $\phi_E|_T$ cannot be constant. As a consequence, the divergence of either

$$\begin{bmatrix} \phi_E \\ 0 \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} 0 \\ \phi_E \end{bmatrix}$$

has to be nonzero. Let Φ_E denote one of these basis functions with $\text{div}(\Phi_E|_T) \neq 0$. In the same manner, we select a basis function for every edge in the range of ι and obtain the ansatz space

$$V_{h,2} := \text{span}\{\Phi_E \mid E \in \text{Range}(\iota)\}. \tag{3.16}$$

All remaining basis functions span the discrete space $V_{h,1}$. With the given decomposition of V_h , we obtain the following result.

Lemma 3.8 (Decomposition for Crouzeix–Raviart). *The discretization scheme $V_h - Q_h$ from (3.14) with the decomposition $V_h = V_{h,1} \oplus V_{h,2}$ defined in (3.16) satisfies Assumption 3.1.*

Proof. The matrix B_2 from Assumption 3.1 corresponds to the discrete space $V_{h,2}$ and is defined by

$$B_{2,ij} = \int_{\Omega} \chi_i \operatorname{div}(\Phi_j) \, dx = \int_{T_i} \operatorname{div}(\Phi_j) \, dx. \quad (3.17)$$

Therein, $\{\Phi_j\}$ denote the basis functions of $V_{h,2}$ and $\{\chi_i\}$ the basis functions of Q_h , *i.e.*, $\chi_i = 1$ on the triangle T_i and 0 elsewhere. Note that χ_0 , where the pressure is fixed to be 0, is not a basis function of Q_h and, thus, excluded from the considerations. Also, note that we assume an order of the basis functions according to Remark 3.7. Since $\operatorname{div}(\Phi_i) \neq 0$ on T_i by construction, the diagonal entries of B_2 are nonzero. Furthermore, every column can only have two entries because of the support of edge-bubble functions. By the construction of Algorithm 3.5, the second entry can only be above the diagonal and thus, B_2 is upper triangular and nonsingular. \square

Remark 3.9 (Outflow boundary conditions). For flow problems that have an outflow with *homogeneous Neumann* or *do nothing* conditions, the pressure must not be fixed, *cf.* (3.13). In this case, Algorithm 3.5 defines $V_{h,2}$ if one starts with a T_0 that shares an edge E_0 with the outflow boundary. Then, the inclusion of χ_0 in the definition (3.17) leads to a $B_2 \in \mathbb{R}^{m,m-1}$ that is *Hessenberg* with the last column missing and with nonzero entries on the subdiagonal. Adding the Φ_{E_0} to $V_{h,2}$ for which $\operatorname{div}(\Phi_{E_0}|_{T_0}) \neq 0$, we add a column that is zero except from the first row's entry and that makes B_2 square and invertible.

Remark 3.10 (Condition number). The condition number of the matrix B_2 obtained by Algorithm 3.5 and Lemma 3.8 scales as h^{-1} where h denotes the mesh-size. For a uniform mesh of the unit square where Algorithm 3.5 runs without reset, *i.e.*, without entering step 2b, the matrix B_2 has the structure

$$B_2 = \begin{bmatrix} h & h & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & h \\ & & & & h \end{bmatrix}, \quad B_2 B_2^T = h \begin{bmatrix} h & 1 & & & \\ 1 & \ddots & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & \ddots & \ddots & 1 \\ & & & 1 & h \end{bmatrix}.$$

The eigenvalues of $h^{-1} B_2 B_2^T$ are given by

$$\lambda_j = h + 2 \cos(j\pi h^2/2), \quad j = 1, \dots, n = 2h^{-2} - 1.$$

Hence, a rough estimate of the condition number yields

$$\operatorname{cond} B_2 = \frac{\lambda_{\max}}{\lambda_{\min}} \approx \frac{h+2}{h} \approx \frac{2}{h}.$$

Note, however, that a degeneration of the mesh may lead to large deviations.

Remark 3.11 (Extension to three space dimensions). The finite element spaces V_h and Q_h of this subsection have a straightforward analogon in three space dimensions [13]. Also Algorithm 3.5 can easily be adapted by the use of tetrahedra and faces in place of triangles and edges. Hence, the given results also apply for three-dimensional simulations.

3.5.2. Continuous velocity

The second example applies a continuous approximation of the velocity but keeps the piecewise constants for the pressure as in (3.14). Since the $[\mathcal{S}_{1,0}(\mathcal{T})]^2 - \mathcal{P}_0^0(\mathcal{T})$ scheme is known to be unstable, the ansatz space V_h

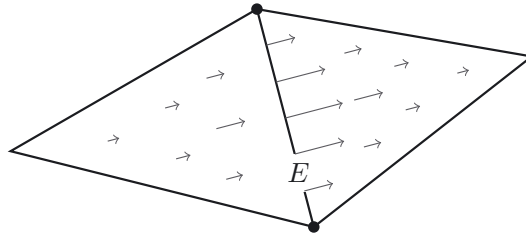


FIGURE 2. Illustration of the vector-valued function $\mathcal{Y}_E = \varphi_1 \varphi_2 \nu_E$.

is enriched by a special type of edge-bubble functions. As in [7, 15], we define for an interior edge $E \in \mathcal{E}_{\text{int}}$ the function

$$\mathcal{Y}_E := \varphi_1 \varphi_2 \nu_E \in \mathcal{V}.$$

Therein, ν_E denotes the outer normal vector and φ_1, φ_2 the hat-functions corresponding the vertices of the edge E . For an illustration of \mathcal{Y}_E see Figure 2. This yields as ansatz spaces for the velocity and the pressure,

$$V_h = [\mathcal{S}_{1,0}(\mathcal{T})]^2 \oplus \{\mathcal{Y}_E \mid E \in \mathcal{E}_{\text{int}}\}, \quad Q_h = \mathcal{P}_0^0(\mathcal{T}). \tag{3.18}$$

The proof of the corresponding inf-sup condition can be found in [7]. In order to define the subspace $V_{h,2}$, we again use the mapping $\iota : \mathcal{T} \setminus \{T_0\} \rightarrow \mathcal{E}_{\text{int}}$ given by Algorithm 3.5. Therewith, we define

$$V_{h,2} = \{\mathcal{Y}_E \mid E \in \text{Range}(\iota)\} \tag{3.19}$$

and $V_{h,1}$ as the span of all remaining basis functions of V_h .

Lemma 3.12 (Decomposition for Bernardi–Raugel). *The discretization scheme $V_h - Q_h$ from (3.18) with the given decomposition $V_h = V_{h,1} \oplus V_{h,2}$ defined in (3.19) satisfies Assumption 3.1.*

Proof. Note that the structure of B_2 is as in Lemma 3.8. Thus, it remains to show that the integral of $\text{div } \mathcal{Y}_E$ does not vanish. By definition of \mathcal{Y}_E , it holds that

$$\text{div } \mathcal{Y}_E = \nabla(\varphi_1 \varphi_2) \cdot \nu_E = \varphi_1 \nabla \varphi_2 \cdot \nu_E + \varphi_2 \nabla \varphi_1 \cdot \nu_E.$$

Hence, for a triangle T with edge E ,

$$\begin{aligned} \int_T \text{div } \mathcal{Y}_E \, dx &= \nabla \varphi_2 \cdot \nu_E \int_T \varphi_1 \, dx + \nabla \varphi_1 \cdot \nu_E \int_T \varphi_2 \, dx \\ &= \frac{|T|}{3} (\nabla \varphi_2 + \nabla \varphi_1) \cdot \nu_E. \end{aligned}$$

Let $[x_i, y_i]^T, i = 1, 2$, denote the coordinates of the nodes corresponding to φ_1 and φ_2 , respectively. Without loss of generality, we assume that the third node is located in $[0, 0]^T$. Then, the outer normal vector for E is, up to a constant, given by $\nu_E = [y_1 - y_2, x_2 - x_1]^T$. The hat-functions are defined by

$$\varphi_1(x, y) = \frac{1}{d}(y_2 x - x_2 y), \quad \varphi_2(x, y) = \frac{1}{d}(-y_1 x + x_1 y)$$

with $d = x_1 y_2 - x_2 y_1 \neq 0$ since the triangle is part of a regular triangulation. Thus, we obtain

$$(\nabla \varphi_2 + \nabla \varphi_1) \cdot \nu_E = -\frac{1}{d}((x_1 - x_2)^2 + (y_1 - y_2)^2) \neq 0$$

and therefore the claim $\int_T \text{div } \mathcal{Y}_E \, dx \neq 0$. □

Remark 3.13. As for the Crouzeix–Raviart case, the scheme (3.18) can be extended to the three-dimensional case [7]. Accordingly to Lemma 3.12, one can show that the integral of the divergence of the basis functions does not vanish on certain tetrahedra. The full-rank property of B_2 then follows as in the two-dimensional case.

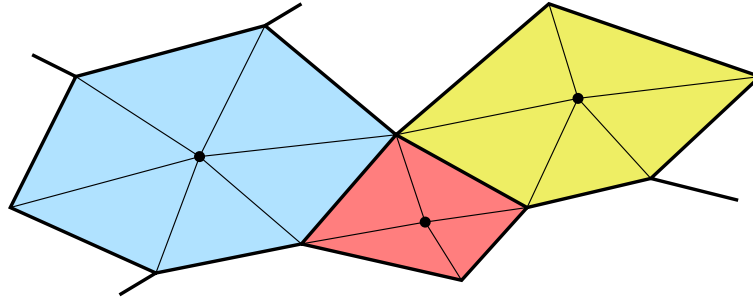


FIGURE 3. Sample of a triangulation \mathcal{T} and its decomposition into three macro elements.

3.5.3. Continuous pressure

This subsection is devoted to the decomposition of the popular Taylor–Hood element [37] in which the discretized pressure is continuous and the velocity is of higher order. The finite element spaces for this case are given in (3.15). The proof of the inf-sup condition is given in [40] or, using local arguments and macro elements, in ([15], Chap. 2.4.2). In the sequel, v_p denotes the boundary node on which the pressure has no degree of freedom.

As in ([15], Chap. II.4), we consider a triangulation \mathcal{T} which can be decomposed into macro elements in the form of node patches (of interior nodes), see Figure 3. Thus, we assume that there exist macro elements $\{\Omega_r\}_{r=1,\dots,R}$, each with exactly one interior node, which form a partition of $\bar{\Omega}$. The triangulation of Ω_r is given by the restriction of \mathcal{T} on $\bar{\Omega}_r$ and is denoted by \mathcal{T}_r . In addition, we assume that the macro elements are ordered such that v_p is a node of \mathcal{T}_1 and that \mathcal{T}_r , $2 \leq r \leq R$, has a common node with at least one \mathcal{T}_k for some $k \leq r-1$. In order to decompose the finite element space V_h such that Assumption 3.1 is fulfilled, we establish a one-to-one map $j : \mathcal{N} \setminus \{v_p\} \rightarrow \mathcal{E}_{\text{int}}$. We define j by the following algorithm, which additionally introduces sets of nodes $\mathcal{I}_r \subset \mathcal{N}(\mathcal{T}_r)$.

Algorithm 3.14 (Mapping j). Set $\mathcal{N}_R := \mathcal{N} \setminus \{v_p\}$. Iterate over macro elements, *i.e.*, over $1 \leq r \leq R$:

Step 1. Consider the nodes $\mathcal{I}_r := \mathcal{N}_R \cap \mathcal{N}(\mathcal{T}_r) = \{v_0, \dots, v_{k(r)}\}$ where v_0 denotes the middle node, as shown in Figure 4.

Step 2. Define E_j as the edge between v_0 and v_j for $j = 1, \dots, k(r)$ and E_0 as any other edge of $\mathcal{E}(\mathcal{T}_r)$ which has v_0 as an endpoint.

Step 3. Set $j(v_j) := E_j$ for $j = 0, \dots, k(r)$ and reset $\mathcal{N}_R := \mathcal{N}_R \setminus \mathcal{I}_r$. If $r \neq R$, return to Step 1 with $r := r + 1$.

Remark 3.15. The order of the macro elements and the fact that $v_p \in \mathcal{N}(\mathcal{T}_1)$ guarantees that at least one node of $\mathcal{N}(\mathcal{T}_r)$ is not included in $\mathcal{N}_R \cap \mathcal{N}(\mathcal{T}_r)$. As a consequence, the second step of Algorithm 3.14 is always realizable.

It remains to define the subspaces $V_{h,1}$ and $V_{h,2}$. Similar to the previous decompositions, let Ψ_E denote a function which vanishes in one component and equals the corresponding edge-bubble function ψ_E in the other component. The precise order of the two components depends on the geometry, see the discussion for n -gons below. We then define

$$V_{h,2} := \text{span}\{\Psi_E \mid E \in \text{Range}(j)\}$$

and $V_{h,1}$ as the span of all remaining basis functions of V_h .

In the sequel, we denote by $B_{\mathcal{I}_r}$ the submatrix of B which corresponds to the pressure nodes \mathcal{I}_r (defined in Algorithm 3.14) and the edge-bubble functions of edges in $j(\mathcal{I}_r)$.

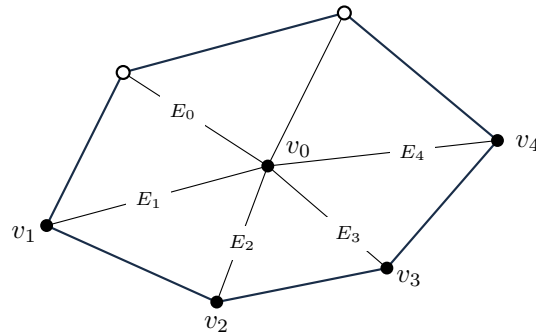


FIGURE 4. A single macro element Ω_r with an illustration of the map $j : \mathcal{N} \setminus \{v_p\} \rightarrow \mathcal{E}_{\text{int}}$, $j(v_k) = E_k$. Nodes of the type \circ are not part of \mathcal{N}_r and thus, already covered by previous macro elements.

Lemma 3.16 (Localization of Taylor–Hood). *Let \mathcal{T} be a triangulation which can be decomposed into macro elements as illustrated in Figure 3. If all submatrices $B_{\mathcal{T}_r}$, $1 \leq r \leq R$, are invertible, then the discretization scheme (3.15) with the decomposition $V_h = V_{h,1} \oplus V_{h,2}$ satisfies Assumption 3.1.*

Proof. We show that the invertibility of the matrix B_2 from Assumption 3.1 follows from the invertibility of the local matrices. For this, the essential observation is that the ordering of macro elements together with Algorithm 3.14 leads to the block structure

$$B_2 = \begin{bmatrix} B_{\mathcal{T}_1} & * & * & * \\ & B_{\mathcal{T}_2} & * & * \\ & & \ddots & * \\ 0 & & & B_{\mathcal{T}_R} \end{bmatrix}.$$

Thus, the invertibility of $B_{\mathcal{T}_r}$, $1 \leq r \leq R$, gives the assertion. □

Remark 3.17 (n -gons). On equilateral n -gons, the invertibility of the submatrix $B_{\mathcal{T}_r}$ depends on a single parameter, namely the angle enclosed by an interior edge and the x -axis. Let α_j denote the angle enclosed by the interior edge E_j and the x -axis. Then, the decision rule

$$\Psi_{E_j} = \begin{bmatrix} 0 \\ \psi_{E_j} \end{bmatrix}, \text{ if } -\frac{\pi}{6} < \alpha_j + \ell\pi < \frac{\pi}{3}, \quad \Psi_{E_j} = \begin{bmatrix} \psi_{E_j} \\ 0 \end{bmatrix} \text{ otherwise}$$

for $\ell \in \mathbb{Z}$ and $j = 0, 1, \dots, k(r)$ renders the *via* Algorithm 3.14 obtained matrix $B_{\mathcal{T}_r}$ invertible. This observation has been numerically proven correct for $3 \leq n \leq 9$ using the code available from the author’s github account [19].

Lemma 3.18 (Anisotropic scaling). *Consider a patch Ω_r with nonsingular matrix $B_{\mathcal{T}_r}$. Then, $B_{\mathcal{T}_r}$ remains nonsingular under anisotropic scalings of Ω_r , i.e., transformations of the form $S(x, y) = (ax, by)$ with $a, b > 0$.*

Proof. Let $\hat{\Omega}_r = S(\Omega_r)$ denote the transformed patch and $\hat{\varphi}_i, \hat{\Psi}_j$ the corresponding basis functions. Since $|\det DS| = ab \neq 0$, the transformation formula gives for a transformed entry of $B_{\mathcal{T}_r}$,

$$\hat{B}_{\mathcal{T}_r,ij} = \int_{\hat{\Omega}_r} \hat{\varphi}_i \operatorname{div} \hat{\Psi}_j \, dx = ab \int_{\Omega_r} \varphi_i \left(a \frac{\partial}{\partial x} \Psi_j + b \frac{\partial}{\partial y} \Psi_j \right) dx.$$

Since Ψ_j vanishes in one component, it holds $\hat{B}_{\mathcal{T}_r,ij} = c \cdot B_{\mathcal{T}_r,ij}$ with either $c = a^2b$ or $c = ab^2$. In any case, this constant is the same for the entire column of $B_{\mathcal{T}_r}$ and thus, just a nonzero factor of the determinant. □

Remark 3.19 (Extension to three space dimensions). Also the Taylor–Hood scheme (3.15) has an extension in three space dimensions. However, the algorithm to find an invertible block of the B matrix is much more involved.

3.6. Quadrilateral meshes

We close this section with a brief overview of stable finite element schemes on quadrilateral meshes and corresponding decompositions of the velocity space. Here, the triangulation \mathcal{T} is supposed to be a partition of $\bar{\Omega}$ into convex quadrilaterals. For quadrilateral elements, one considers the space of piecewise polynomials of *partial* degree k which are globally continuous.

As for the triangular case, there are finite element schemes of Taylor–Hood type ([15], Chap. II.3.2) and Bernardi–Raugel type where the velocity space is enriched by the fluxes through interior edges ([15], Chap. II.3.1). The analog of the discontinuous approach of Crouzeix–Raviart was introduced by Rannacher and Turek [34] and is given by

$$V_h = [\tilde{Q}_{1,0}(\mathcal{T})]^2, \quad Q_h = \mathcal{P}_0^0(\mathcal{T}). \quad (3.20)$$

Therein, $\tilde{Q}_{1,0}$ denotes the nonconforming space which has one degree of freedom per interior edge. In contrast to the Crouzeix–Raviart element, functions in V_h are not piecewise affine. Piecewise affine functions which are continuous in the midpoints of edges were introduced by Park and Sheen [31], see also [2]. Unfortunately, there is no known stability result for this kind of element.

In a thorough analysis by Turek, the nonconforming element (3.20) was found superior over comparable conforming elements in terms of stability, accuracy, and efficiency ([39], Chap. 3.1.1). The higher stability and accuracy of the nonconforming scheme is ascribed to the robustness of the inf-sup constant against mesh deformations.

A decomposition of V_h from (3.20) into $V_{h,1}$ and $V_{h,2}$ in the sense of Assumption 3.1 works exactly as in Section 3.5.1.

4. NUMERICAL EXAMPLES

This section illustrates the benefits of the index-1 formulation (3.4) for numerical time integration by means of two examples.

4.1. Flow in a square

As a first example we consider a flow in a square with a constructed solution. To isolate the high index effects, we consider a variant of (2.1) without the term $\frac{1}{\text{Re}}\Delta u$ which introduces stiffness to the spatially discretized system and thus, step size restrictions for explicit schemes. As exact solution for the velocity field and the pressure in time and the two spatial coordinates, we set

$$\begin{aligned} u_1(t; x_1, x_2) &= 2 \sin(8t) \cdot x_1^2(1 - x_1)^2 x_2(1 - x_2)(2x_2 - 1), \\ u_2(t; x_1, x_2) &= 2 \sin(8t) \cdot x_2^2(1 - x_2)^2 x_1(1 - x_1)(1 - 2x_1), \\ p(t; x_1, x_2) &= \sin(8t) \cdot x_1(1 - x_1)x_2(1 - x_2). \end{aligned}$$

The corresponding right-hand sides as well as the boundary and initial values are constructed accordingly. On the computational domain $(0, T) \times \Omega := (0, 1) \times [0, 1]^2$, this gives zero initial and zero Dirichlet boundary conditions.

The triangulation \mathcal{T}_N of the spatial domain is characterized by the parameter N , meaning that the unit square is uniformly divided into $(N - 1)^2$ squares that are clusters of four triangles each, see Figure 5. Besides, we choose as velocity and pressure state space the Taylor–Hood discretization (3.15).

This *criss-cross* triangulation and the Taylor–Hood elements enable the splitting $V_h = V_{h,1} \oplus V_{h,2}$ via Algorithm 3.14, cf. also Lemma 3.16 and Remark 3.17. Thus, the square matrix B_2 as defined in Assumption 3.1 is

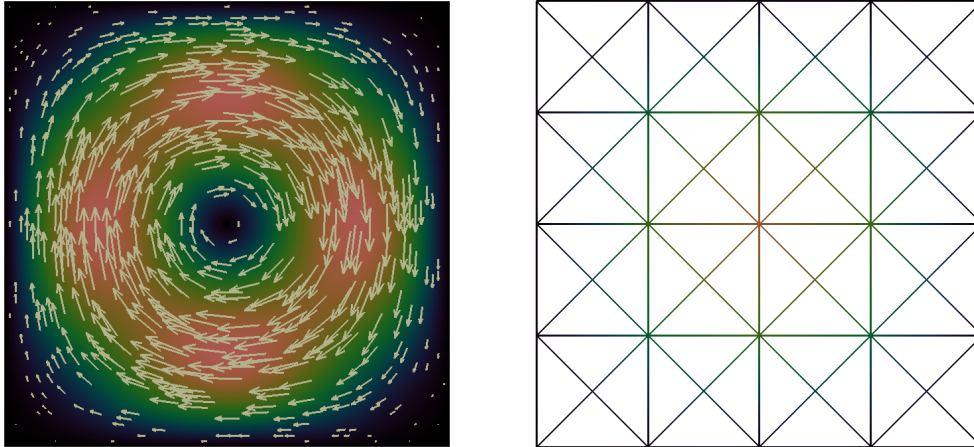


FIGURE 5. Illustration of the velocity field at $t = \frac{1}{4}$ and of the *criss-cross* triangulation for $N = 5$.

invertible and the index-1 formulation (3.4) is applicable. To investigate the time integration error, we discretize the time interval into $2^k + 1$, $k = 4, \dots, 10$, time instances and apply the semi-explicit Euler method. For the index-1 and index-2 formulations, this leads to the update formulas (3.9) and (3.8), respectively.

The resulting linear systems are solved iteratively up to an absolute residual smaller than tol . Since the solver considers relative residuals, in every iteration we corrected the tolerance by the factor $1/\|\text{rhs}^c\|$, where rhs^c is the current right-hand side. For computing $\|\text{rhs}^c\|$ and the residuals, we use the norms induced by the inner product of the discrete L^2 -spaces. For the index-2 system, this is the inner product with respect to the inverses of the mass matrices of the finite element bases. In the solution of the index-1 update (3.9), where we used the block preconditioner with $[M_D^{-1}, B_2^{-1}M_{D,2}B_2^{-T}, B_2^{-1}]$ on the diagonal with M_D denoting the diagonal of the mass matrix of the velocity space, this was approximated by the scalar product weighted by the mass matrices. Note, that the use of B_2^{-1} is cheap, because of the blockdiagonal structure, *cf.* the proof of Lemma 3.16. It has turned out that for the numerical approximation of the index-2 formulation (3.8), it is beneficial to scale the differential equation by τ .

By construction, the exact solution is known. The error of the numerical approximation is measured by taking the L^2 -norm in space and evaluating the L^2 -norm in time with the piecewise trapezoidal rule.

The numerical experiments show the improvements of the index-1 formulation for the pressure approximation, see Figure 6. As predicted by the theoretical considerations in (3.12), in the index-2 formulation, a numerical error in the algebraic constraints leads to a linear growth in the pressure error with decreasing time step sizes. A smaller residual in the continuity equation only postpones this instability. In the index-1 formulation, this systematic instability is not observed.

Remark 4.1. As it can be expected from (3.12), the pressure is better approximated for the index-1 formulation despite the fact that the residuals in the continuity constraint are larger. This difference in the residuals can be explained by two factors. First, different preconditioning leads to different residuals considered by the solver. Second, the continuity constraint residual as a part of the overall residual has a stronger weight in (3.8) than in (3.9).

The expected linear convergence in the pressure approximation, see ([18], Chap. VII.4), is not observed here. This is due to the dominance of the algebraic (for $\text{tol} = 9.8 \cdot 10^{-4}$) or the spatial discretization error

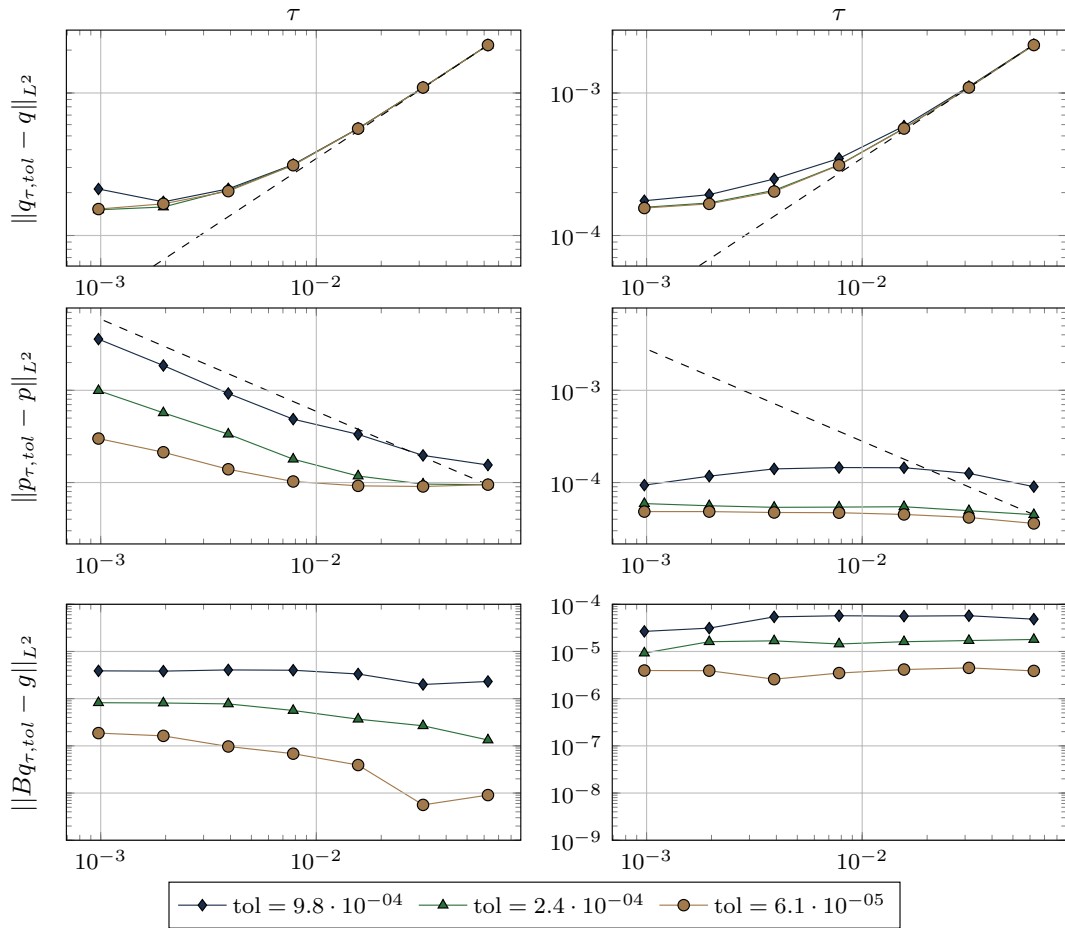


FIGURE 6. The evolution of the errors or residuals of the index-2 (*left*) and index-1 (*right*) formulation for varying time discretization parameter τ and tol for the flow on the square. The space discretization is fixed with $N = 40$. The dashed lines are the linear fit.

(for lower tolerances in the linear system solves). This guess is backed by the convergence plot for a finer spatial discretization ($N = 100$) and direct solves of the linear systems showing linear convergence for coarser time discretizations, see Figure 7.

4.2. Cylinder wake

As second example we consider the Navier-Stokes equations for the simulation of a cylinder wake as illustrated in Figure 8. As boundary conditions we set *no-slip* at the walls, a parabola as the inflow profile at the left boundary, and *do-nothing* conditions at the outflow at the right. We consider the flow at $Re = 60$, calculated with the cylinder diameter and the peak inflow velocity. We consider the time evolution of the flow in $[0, 0.2]$, starting with the steady-state Stokes solution.

For the spatial discretization, we use *Crouzeix–Raviart* elements on a nonuniform mesh with about 15 000 velocity nodes and 5000 pressure nodes. We employ Algorithm 3.5 with the modification proposed in Remark 3.9 to compute the splitting $V_h = V_{h,1} \oplus V_{h,2}$ that we need for the index-1 formulation (3.4).

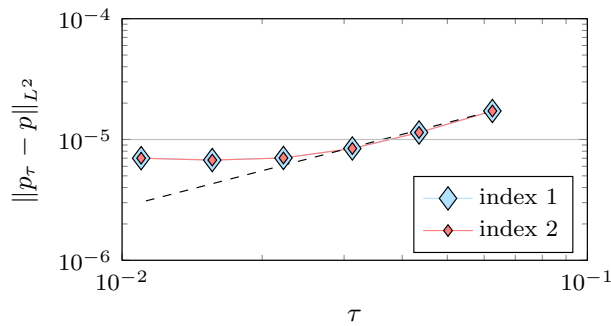


FIGURE 7. The evolution of the errors in the pressure approximation for various time step lengths τ , for the spatial discretization $N = 100$, and for direct solves of the algebraic equations.

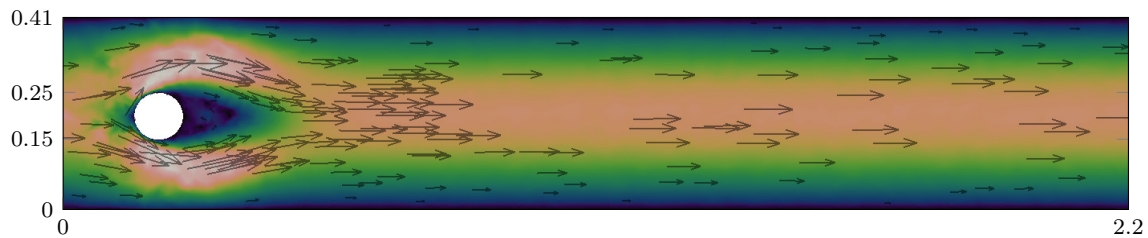


FIGURE 8. Illustration of the cylinder wake with $Re = 60$, started at the steady-state Stokes solution, at time $t = 0.2$.

To account for the stiffness, we now consider an *implicit-explicit* Euler scheme for the discretization which treats the linear diffusion implicitly and the nonlinear convection explicitly. Thus, we consider the update formulas (3.8) and (3.9) but with the discretized diffusion operator appearing in the coefficient matrix.

As in the previous example, we compute the approximation error for various time steps τ and for various accuracy levels tol for the iterative solution of the resulting linear systems. Since there is no analytical solution, we take the result of solving (1.1) with the implicit trapezoidal rule with direct solves and with $\tau = 0.2 \cdot 2^{-11} \approx 10^{-4}$ as the reference.

The results of the numerical investigation are illustrated in Figure 9. Again, the inherent instability of the index-2 formulation is obvious in the plots of the pressure error. Furthermore, since for the cylinder wake the velocity is not discretely divergence free, *i.e.*, g in (3.1b) is not zero due to the boundary conditions, the poor pressure approximation directly affects the velocity approximation. On the other hand, in the index-1 formulation, the expected linear convergence with respect to the time discretization is confirmed both for the velocity and the pressure approximation. A breakdown due to the algebraic error is only observed for a rough tolerance for the linear solver.

The difference in the residual levels in the continuity equation is due to the different preconditioning and different weighting of the overall residuals, *cf.* Remark 4.1. In the numerical tests for the index-1 case, we first observed a steady decrease with τ in the residual. This was due to fact that a factor of $1/\tau$ enters the tolerance correction $1/\|\text{rhs}^c\|$ through the third line in equation (3.9). Therefore, in the computation of the correction we have scaled this equation by $\sqrt{\tau}$ which only worsens the approximation of the linear system.

The code used for the numerical investigations is available from the author's github account [19]. The finite element implementation uses *FEniCS, Version 1.3.0*, [28], the linear systems are solved with *Krypy* [14].

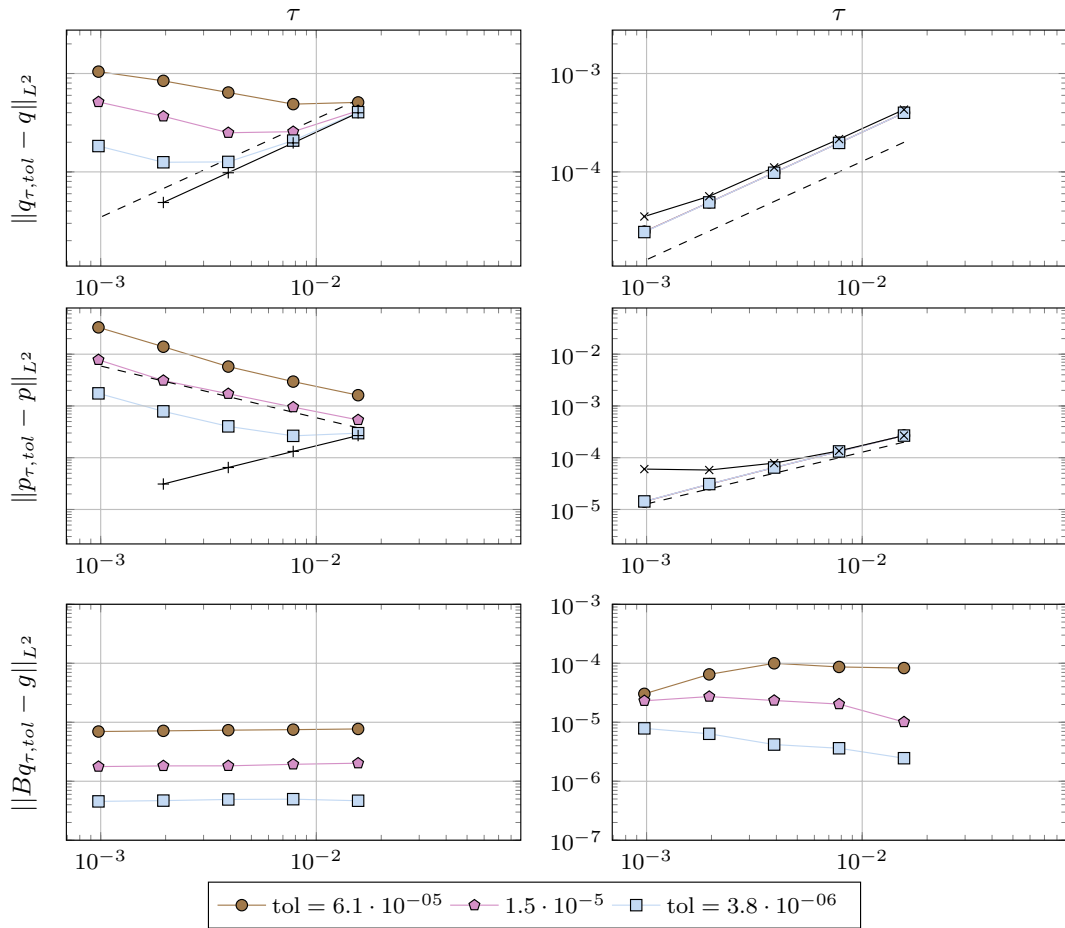


FIGURE 9. The evolution of the errors or residuals of the index-2 (*left*) and index-1 (*right*) formulation for varying time discretization parameter τ and tol for the cylinder wake. The dashed lines are the linear fit. The additional data points for the index-1 case are calculated for the much rougher tolerance $\text{tol} = 3.9 \cdot 10^{-3}$. The additional data points in the index-2 plots are the results for exact solves of the algebraic equations.

5. CONCLUSION

We have presented a new numerical approach to the unsteady NSE by a remodelling of the governing equations. Using analytical insight into the discrete spaces, we made the principles of the index reduction technique of minimal extension applicable to finite element schemes. In particular, the proposed variant preserves sparsity and maintains the physical meaning of the variables since only a permutation is applied. Unlike in penalization or in projection methods, our approach does not require time step restrictions or artificial boundary conditions for the pressure.

The necessary splitting of the finite element space is operated for commonly used Taylor–Hood and Crouzeix–Raviart finite element discretizations. We showed applicability of the proposed algorithms in two numerical examples that backed our theoretical findings. It has also turned out, that for stiff systems as the NSE for viscous fluids, one can consider a combination of our method with *IMEX* schemes that are implicit in the stiff linear part and explicit in the nonlinearity.

We remark that the presented theory ensures consistency and stability of the half-explicit method. In view of efficiency, for practical computations, suitable preconditioners will be necessary which are not yet available for the newly developed index-1 formulation (3.9). However, probably due to the gained stability, for the viscous flow around the cylinder, the *GMRes* iterations for the index-1 system converged significantly faster despite the increased system size if compared to the index-2 formulation.

REFERENCES

- [1] R. Altmann, Index reduction for operator differential-algebraic equations in elastodynamics. *Z. Angew. Math. Mech. (ZAMM)* **93** (2013) 648–664.
- [2] R. Altmann and C. Carstensen, P_1 -nonconforming finite elements on triangulations into triangles and quadrilaterals. *SIAM J. Numer. Anal.* **50** (2012) 418–438.
- [3] M. Arnold, Half-explicit Runge-Kutta methods with explicit stages for differential-algebraic systems of index 2. *BIT* **38** (1998) 415–438.
- [4] M. Arnold, K. Strehmel and R. Weiner, Half-explicit Runge-Kutta methods for semi-explicit differential-algebraic equations of index 1. *Numer. Math.* **64** (1993) 409–431.
- [5] U. M. Ascher, H. Chin, L.R. Petzold and S. Reich, Stabilization of constrained mechanical systems with DAEs and invariant manifolds. *Mech. Struct. Mach.* **23** (1995) 135–157.
- [6] R. Becker and S. Mao, Quasi-optimality of adaptive nonconforming finite element methods for the Stokes equations. *SIAM J. Numer. Anal.* **49** (2011) 970–991.
- [7] C. Bernardi and G. Raugel, Analysis of some finite elements for the Stokes problem. *Math. Comput.* **44** (1985) 71–79.
- [8] D. Braess, *Finite Elements – Theory, Fast Solvers, and Applications in Solid Mechanics*. Cambridge University Press, New York, 3rd edition (2007).
- [9] S.C. Brenner and L.R. Scott, *The Mathematical Theory of Finite Element Methods*, 3rd edition. Springer-Verlag, New York (2008).
- [10] F. Brezzi and M. Fortin, *Mixed and Hybrid Finite Element Methods*. Springer-Verlag, New York (1991).
- [11] S.L. Campbell and C.W. Gear, The index of general nonlinear DAEs. *Numer. Math.* **72** (1995) 173–196.
- [12] P.G. Ciarlet, *The Finite Element Method for Elliptic Problems*. North-Holland, Amsterdam (1978).
- [13] M. Crouzeix and P.-A. Raviart, Conforming and nonconforming finite element methods for solving the stationary Stokes equations. I. *Rev. Franc. Automat. Inform. Rech. Operat.* **7** (1973) 33–75.
- [14] A. Gaul, Krypy. Public Git Repository, Commit: 110a1fb756fb. Iterative Solvers for Linear Systems. Available at <https://github.com/andrenarchy/krypy>.
- [15] V. Girault and P.-A. Raviart, *Finite Element Methods for Navier-Stokes Equations*. Springer-Verlag, Berlin (1986).
- [16] P.M. Gresho and R.L. Sani, *Incompressible Flow and the Finite Element Method. Isothermal Laminar Flow*, vol. 2. Wiley, Chichester (2000).
- [17] E. Hairer, C. Lubich and M. Roche, The numerical solution of differential-algebraic systems by Runge-Kutta methods, vol. 1409 of *Lect. Notes Math.* Springer-Verlag, Berlin (1989).
- [18] E. Hairer and G. Wanner, *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*, 2nd edition. Springer-Verlag, Berlin (1996).
- [19] J. Heiland, TayHoodMinExtForFlowEqns. Public Git Repository, Commit: 8eb641f21d. Solution of time-dependent 2D nonviscous flow with nonconforming minimal extension. Available at <https://github.com/highlando/TayHoodMinExtForFlowEqns>.
- [20] J. Heinrich and C. Vionnet, The penalty method for the Navier-Stokes equations. *Arch. Comput. Method E.* **2** (1995) 51–65.
- [21] M. Hinze, *Optimal and instantaneous control of the instationary Navier-Stokes equations*. Habilitationsschrift, Technische Universität Berlin, Institut für Mathematik (2000).
- [22] R. Kouhia and R. Stenberg, A linear nonconforming finite element method for nearly incompressible elasticity and Stokes flow. *Comput. Methods Appl. Mech. Engrg.* **124** (1995) 195–212.
- [23] P. Kunkel and V. Mehrmann, Index reduction for differential-algebraic equations by minimal extension. *Z. Angew. Math. Mech.* **84** (2004) 579–597.
- [24] P. Kunkel and V. Mehrmann, *Differential-Algebraic Equations: Analysis and Numerical Solution*. European Mathematical Society (EMS), Zürich (2006).
- [25] S. Le Borne and D. Cook II, Construction of a discrete divergence-free basis through orthogonal factorization in \mathcal{H} -arithmetic. *Computing* **81** (2007) 215–238.
- [26] P. Lin, A sequential regularization method for time-dependent incompressible Navier-Stokes equations. *SIAM J. Numer. Anal.* **34** (1997) 1051–1071.
- [27] V.H. Linh and V. Mehrmann, Efficient integration of matrix-valued non-stiff DAEs by half-explicit methods, Technische Universität Berlin, Germany (2011) Preprint 2011–16.
- [28] A. Logg, K. Ølgaard, M. Rognes and G. Wells, Ffc: the fenics form compiler. In *Automated Solution of Differential Equations by the Finite Element Method*. Springer-Verlag, Berlin (2012) 227–238.
- [29] G. Matthies and F. Schieweck, A multigrid method for incompressible flow problems using quasi divergence free functions. *SIAM J. Sci. Comput.* **28** (2006) 141–171.

- [30] G.-P. Ostermeyer, On Baumgarte stabilization for differential-algebraic equations. In Real-Time Integration Methods for Mechanical System Simulation. In vol. 69 of *NATO ASI Series*, edited by E. Haug and R. Deyo. Springer-Verlag, Berlin (1991) 193–207.
- [31] C. Park and D. Sheen, P_1 -nonconforming quadrilateral finite element methods for second-order elliptic problems. *SIAM J. Numer. Anal.* **41** (2003) 624–640.
- [32] A. Quarteroni and A. Valli, Numerical Approximation of Partial Differential Equations. Springer-Verlag, Berlin (1994).
- [33] R. Rannacher, On the numerical solution of the incompressible Navier-Stokes equations. *Z. Angew. Math. Mech.* **73** (1993) 203–216.
- [34] R. Rannacher and S. Turek, Simple nonconforming quadrilateral Stokes element. *Numer. Methods Partial Differ. Eqs.* **8** (1992) 97–111.
- [35] J. Shen, On error estimates of the penalty method for unsteady Navier-Stokes equations. *SIAM J. Numer. Anal.* **32** (1995) 386–403.
- [36] L. Tartar, An Introduction to Navier-Stokes Equation and Oceanography. Springer-Verlag, Berlin (2006).
- [37] C. Taylor and P. Hood, A numerical solution of the Navier-Stokes equations using the finite element technique. *Int. J. Comput. Fluids* **1** (1973) 73–100.
- [38] R. Temam, Navier-Stokes Equations. Theory and Numerical Analysis. North-Holland, Amsterdam (1977).
- [39] S. Turek, Efficient Solvers for Incompressible Flow Problems. An Algorithmic and Computational Approach. Springer-Verlag, Berlin (1999).
- [40] R. Verfürth, Error estimates for a mixed finite element approximation of the Stokes equations. *RAIRO Anal. Numér.* **18** (1984) 175–182.
- [41] R. Verfürth, A Review of A Posteriori Error Estimation and Adaptive Mesh-Refinement Techniques. Wiley-Teubner, Stuttgart (1996).
- [42] J. Weickert, *Applications of the Theory of Differential-Algebraic Equations to Partial Differential Equations of Fluid Dynamics*. Ph.D. thesis, TU Chemnitz, Fakultät Mathematik, Chemnitz (1997).