

---

Konrad-Zuse-Zentrum  
für Informationstechnik Berlin

Takustraße 7  
D-14195 Berlin-Dahlem  
Germany

STEFAN VATER AND RUPERT KLEIN

# **Stability of a Cartesian Grid Projection Method for Zero Froude Number Shallow Water Flows**

Submitted to "Numerische Mathematik"

---

ZIB-Report 07-13 (June 2007)



# Stability of a Cartesian Grid Projection Method for Zero Froude Number Shallow Water Flows

Stefan Vater\*      Rupert Klein<sup>†</sup>

May 30, 2007

## Abstract

In this paper a Godunov-type projection method for computing approximate solutions of the zero Froude number (incompressible) shallow water equations is presented. It is second-order accurate and locally conserves height (mass) and momentum. To enforce the underlying divergence constraint on the velocity field, the predicted numerical fluxes, computed with a standard second order method for hyperbolic conservation laws, are corrected in two steps. First, a MAC-type projection adjusts the advective velocity divergence. In a second projection step, additional momentum flux corrections are computed to obtain new time level cell-centered velocities, which satisfy another discrete version of the divergence constraint.

The scheme features an exact and stable second projection. It is obtained by a Petrov-Galerkin finite element ansatz with piecewise bilinear trial functions for the unknown incompressible height and piecewise constant test functions. The stability of the projection is proved using the theory of generalized mixed finite elements, which goes back to NICOLAÏDES [1982]. In order to do so, the validity of three different inf-sup conditions has to be shown.

Since the zero Froude number shallow water equations have the same mathematical structure as the incompressible Euler equations of isentropic gas dynamics, the method can be easily transferred to the computation of incompressible variable density flow problems.

**Mathematics Subject Classification (2000):** 65M12, 76M12, 76M10, 35L65

**Keywords:** incompressible flows, shallow water equations, projection method, stability, mixed finite elements, inf-sup-condition

---

\*Numerische Mathematik/Scientific Computing, Freie Universität (FU) Berlin, FB Mathematik & Informatik, Institut für Mathematik, Arnimallee 6, D-14195 Berlin, E-mail: [vater@math.fu-berlin.de](mailto:vater@math.fu-berlin.de)

<sup>†</sup>Konrad-Zuse-Zentrum für Informationstechnik Berlin (ZIB), Takustraße 7, D-14195 Berlin, and Numerische Mathematik/Scientific Computing, FU Berlin, E-mail: [klein@zib.de](mailto:klein@zib.de)

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Governing Equations . . . . .	4
<b>2</b>	<b>The Numerical Method</b>	<b>5</b>
2.1	Construction of the scheme . . . . .	5
2.2	Discretization of the Projections . . . . .	8
2.3	Exact Projection Method . . . . .	12
<b>3</b>	<b>Stability of the second projection</b>	<b>13</b>
3.1	Generalized Saddle Point Problems – Theory . . . . .	14
3.2	Reformulation of the problem . . . . .	16
3.3	Stability analysis of the mixed formulation . . . . .	18
<b>4</b>	<b>Numerical Results</b>	<b>25</b>
4.1	Convergence study . . . . .	27
4.2	Advection of a vortex . . . . .	28
<b>5</b>	<b>Conclusions</b>	<b>29</b>
<b>A</b>	<b>Appendix</b>	<b>33</b>
A.1	Discretization of the new projection . . . . .	33
A.2	Properties of the Lumping-Operator . . . . .	34
	<b>References</b>	<b>38</b>

## 1 Introduction

Starting with the fundamental work of CHORIN [1968] and TEMAM [1968], the use of projection methods for the numerical solution of the incompressible flow equations has a long tradition (see e.g. VAN KAN [1986]; BELL ET AL. [1989]; BELL and MARCUS [1992]; SCHNEIDER ET AL. [1999]; ALMGREN ET AL. [2000] and references therein). In these methods, solutions are first advanced in time ignoring the solenoidal constraint of the velocity field. In a second step, the velocity field is corrected using a suitable approximation of the incompressible pressure to enforce compliance with the divergence constraint.

The stability of the projection step in exact projection methods for the incompressible Euler or shallow water equations has been an unsolved issue in the past. Difficulties arise in this context from a decoupling of the velocity and the pressure variables, which, in turn, is a consequence of using discrete gradient approximations with kernel dimension larger than one. Examples of such methods are given by BELL ET AL. [1989], BELL and MARCUS [1992] and SCHNEIDER ET AL. [1999]. To resolve this problem, approximate projection methods were introduced by ALMGREN ET AL. [1996], which use the

same discrete divergence and gradient operators as in exact projection methods, but a modified version of the discrete Laplacian. This approach results in velocity fields that satisfy the underlying divergence constraint only up to the order of accuracy of the gradient and divergence discretizations. In the present paper we propose an alternative approach that utilizes discretizations of the differential operators, which guarantee exact projections while avoiding the velocity-pressure decoupling. The discretization goes back to SÜLI [1991], and can be derived by a Petrov-Galerkin finite element ansatz with piecewise bilinear trial functions for the unknown incompressible pressure and piecewise constant test functions.

The divergence constraint on the velocity field, which arises in the zero Mach number limit of the Euler equations [KLAINERMAN and MAJDA, 1981; SCHOCHET, 1994; KLEIN, 1995] (see also the review by SCHOCHET [2005]), leads to a saddle point problem, in which the velocity is coupled with the gradient of the incompressible pressure. The fundamental theory of (discretizations of) such problems goes back to BABUŠKA [1971] and BREZZI [1974], who analyzed finite element schemes for elliptic partial differential equations with additional side constraints. This theory provides the so-called “inf-sup conditions” for existence and uniqueness of solutions and stable discretizations of such problems.

To the best of the authors knowledge, stability estimates of the Babuška-Brezzi-type have not been derived for projection methods applied to inviscid flow problems so far. This is different in the viscous case (cf. GUERMOND ET AL. [2006]). However, in contrast to the inviscid case in the incompressible Navier-Stokes equations the Laplacian of the velocity interacts with the pressure gradient, which leads to a saddle point problem of the Stokes type involving higher spatial derivatives compared to the inviscid case. Consequently, the stability proofs for methods solving the Navier-Stokes equations cannot be easily transferred.

The presented method is a non-incremental pressure-correction method for the incompressible (zero Froude number) shallow water equations. To represent advection of mass and momentum, the scheme relies on second order conservative finite volume Godunov-type methods in its predictor step. It is shown that the projection step, which corrects the cell-centered momentum to satisfy the underlying divergence constraint, is stable in the sense of mixed finite element methods, which is the main result of this paper and summarized in Theorem 3.10. The discretization features both, a compact Poisson stencil, and an exact projection. The key to achieving both of these properties at the same time lies in the fact that we let part of the in-cell slopes, which are normally determined by standard slope limiting procedures, be assigned in the projection step.

After introducing the governing equations and the consequences of the zero Froude number limit in the remainder of the introduction, we describe the construction of the numerical method in Section 2. The stability of the projection step is investigated using the theory of generalized mixed finite elements in Section 3. To demonstrate the applicability of the scheme, some basic numerical test cases are presented in Section 4. The major conclusions of this work are reported in the last Section.

## 1.1 Governing Equations

The shallow water equations are a hyperbolic system of conservation laws. In their non-dimensional form they are given by the two equations

$$\begin{aligned} \text{Sr} \frac{\partial h}{\partial t} + \nabla \cdot (h\mathbf{v}) &= 0 \\ \text{Sr} \frac{\partial(h\mathbf{v})}{\partial t} + \nabla \cdot \left( h\mathbf{v} \circ \mathbf{v} + \frac{1}{2\text{Fr}^2} h^2 \mathbf{I} \right) &= 0, \end{aligned} \quad (1)$$

which express conservation of height  $h$  and momentum  $h\mathbf{v}$ . Here, two dimensionless characteristic quantities have been introduced, namely

$$\text{Sr} := \frac{\ell'_{\text{ref}}}{t'_{\text{ref}} v'_{\text{ref}}} \quad \text{and} \quad \text{Fr} := \frac{v'_{\text{ref}}}{\sqrt{g' h'_{\text{ref}}}},$$

which are known as the *Strouhal* and the *Froude number*, respectively. The first one describes the ratio between the advection timescale  $\ell'_{\text{ref}}/v'_{\text{ref}}$  and the reference timescale  $t'_{\text{ref}}$ , whereas the latter gives the ratio between the reference velocity  $v'_{\text{ref}}$  and the gravity wave speed  $\sqrt{g' h'_{\text{ref}}}$  (celerity). In the following, we are interested in a reference time scale equal to the advection time scale of the fluid, so that  $t'_{\text{ref}} = \ell'_{\text{ref}}/v'_{\text{ref}}$  and the Strouhal number becomes one ( $\text{Sr} = 1$ ).

The zero Froude number limit of (1) can be analyzed by an asymptotic analysis with a small parameter  $\text{Fr}$  [VATER, 2005]. This is similar to the zero Mach number limit of the Euler equations (cf. KLAINERMAN and MAJDA [1981]; KLEIN [1995]), except that in the case of the Euler equations the divergence constraint arises from the energy equation, and not from the mass equation. The resulting limit equations are given by

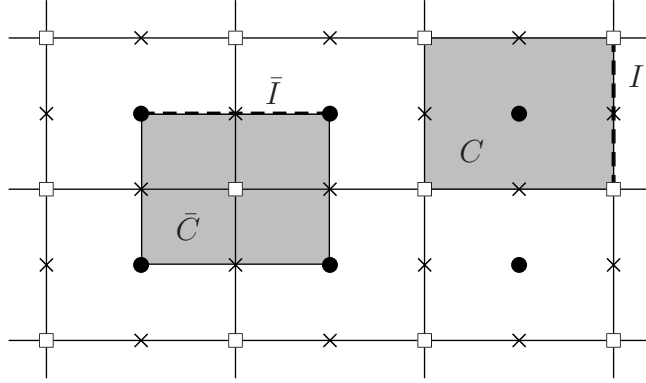
$$\begin{aligned} h_t + \nabla \cdot (h\mathbf{v}) &= 0 \\ (h\mathbf{v})_t + \nabla \cdot (h\mathbf{v} \circ \mathbf{v}) + h\nabla h^{(2)} &= \mathbf{0} \\ h &= h_0(t). \end{aligned} \quad (2)$$

This system of equations is no longer hyperbolic, but of mixed elliptic-hyperbolic type. An additional unknown  $h^{(2)}$  (the ‘‘incompressible height’’) is introduced and the height is split into a time dependent zero-gradient part  $h_0$  and a second order perturbation  $\text{Fr}^2 h^{(2)}$ . Integrating the first equation of (2) over the domain  $\Omega$  and applying the divergence theorem leads to

$$\frac{1}{h_0} \frac{dh_0}{dt} = -\frac{1}{|\Omega|} \int_{\partial\Omega} \mathbf{v} \cdot \mathbf{n} \, d\sigma. \quad (3)$$

Thus, either the change of height is given through von Neumann boundary conditions for the velocity, or the prescription of  $h_0$  implies an integral constraint on the normal velocity field on the boundary of  $\Omega$ . Furthermore, the integration over an arbitrary volume  $V \subset \Omega$  yields

$$\int_{\partial V} (h\mathbf{v}) \cdot \mathbf{n} \, d\sigma = -|V| \frac{dh_0}{dt}, \quad (4)$$



**Figure 1:** Control volume  $C$  and interface  $I$  of the primary discretization and those ( $\bar{C}$  and  $\bar{I}$ ) of the dual discretization. Cell centers are denoted by circles, nodes by squares and midpoints of the interfaces by crosses.

which implies an integral constraint on the velocity divergence in  $V$ .

In terms of optimization problems  $h^{(2)}$  can be viewed as a Lagrange multiplier, which ensures that the velocity field is in compliance with the divergence constraint (4).

## 2 The Numerical Method

The present method is a further development of the projection method presented by SCHNEIDER ET AL. [1999] for the incompressible Euler Equations, which we revisit here for the case of the zero Froude number shallow water equations. The main difference between the present scheme and that of SCHNEIDER ET AL. lies in the discretization of the projection step (see Subsection 2.2).

### 2.1 Construction of the scheme

Throughout this work we assume a Cartesian space discretization of the computational domain  $\Omega$ . In this discretization, the volume of a cell  $C$  is denoted by  $|C|$ , and two neighboring cells are separated by an interface  $I$  with area  $|I|$  (cf. Figure 1).  $\mathcal{C}$  and  $\mathcal{I}$  are defined as the collections of all cells and interfaces, respectively. We denote the set of all interfaces, which are part of the boundary of a cell  $C$ , by  $\mathcal{I}_{\partial C} \subset \mathcal{I}$ .

For the construction of the method, a finite volume scheme in conservation form is considered, i.e.

$$\mathbf{U}_C^{n+1} = \mathbf{U}_C^n - \frac{\delta t}{|C|} \sum_{I \in \mathcal{I}_{\partial C}} |I| \mathbf{F}_I \quad . \quad (5)$$

In (5)  $\mathbf{U}_C^n$  is a numerical approximation to the average of the exact solution  $\mathbf{u}(\mathbf{x}, t)$  of

the problem over cell  $C$  at time  $t^n$ :

$$\mathbf{U}_C^n \approx \frac{1}{|C|} \int_C \mathbf{u}(\mathbf{x}, t^n) d\mathbf{x} \quad , \quad \mathbf{u}(\mathbf{x}, t) := \begin{pmatrix} h \\ h\mathbf{v} \end{pmatrix} .$$

The *numerical flux*  $\mathbf{F}_I$  approximates the average of the flux function

$$\mathbf{f}(\mathbf{u}(\mathbf{x}, t), \mathbf{n}(\mathbf{x})) := \begin{pmatrix} h(\mathbf{v} \cdot \mathbf{n}) \\ h\mathbf{v}(\mathbf{v} \cdot \mathbf{n}) + h h^{(2)} \mathbf{n} \end{pmatrix}$$

of the zero Froude number shallow water equations. For these fluxes, the average is taken over one time step  $[t^n, t^{n+1}]$ , with  $t^{n+1} := t^n + \delta t$ , and over the interface  $I$  between two cells. The flux averages will be computed in three steps:

$$\mathbf{F}_I := \mathbf{F}_I^* + \mathbf{F}_I^{\text{MAC}} + \mathbf{F}_I^{\text{P2}} .$$

First, predictions of the advective fluxes  $\mathbf{F}_I^*$  are computed by the numerical solution of the hyperbolic *auxiliary system*

$$\begin{aligned} h_t^* + \nabla \cdot (h\mathbf{v})^* &= 0 \\ (h\mathbf{v})_t^* + \nabla \cdot \left( (h\mathbf{v} \circ \mathbf{v})^* + \frac{(h^*)^2}{2} \mathbf{I} \right) &= \mathbf{0} . \end{aligned} \tag{6}$$

The computation of the numerical fluxes for these equations is done using an explicit high resolution upwind method for hyperbolic conservation laws (see, e.g. [VAN LEER, 1979]). The current implementation is based on a semi-discrete method with Runge-Kutta time stepping [OSHER, 1985], which is often referred to as the *method of lines*. But the authors have been also successfully implemented a version using Lax-Wendroff-type discretizations as well as operator splitting techniques for the spatial directions, [SCHNEIDER ET AL., 1999]. The stability of the numerical solution of the auxiliary system depends on a CFL time step restriction [COURANT ET AL., 1928]. Since the eigenvalues (characteristic speeds) of this system do not depend on the Froude number, they are of order  $\mathcal{O}(1)$  as  $\text{Fr} \rightarrow 0$ , leading to  $\delta t = \mathcal{O}(\delta x)$  on a regular discretization with grid spacing  $\delta x$ .

Then, a *MAC-type projection* [HARLOW and WELCH, 1965] is applied, which corrects the advection velocity divergence by  $\mathbf{F}_I^{\text{MAC}}$  to be in compliance with the divergence constraint (4) applied to each grid cell. In a final *second projection* the non-convective components of the numerical fluxes, i.e., the pressure (or height) contributions to the momentum fluxes, are corrected by  $\mathbf{F}_I^{\text{P2}}$ , such that the new time level divergence of the cell-centered velocities satisfies (4) for another set of control volumes defined below. Furthermore, in the present new scheme this projection yields updates for the linear reconstructions of momentum in each grid cell.

To achieve second order accuracy in time for the flux components  $\mathbf{F}_I^{\text{MAC}}$  and  $\mathbf{F}_I^{\text{P2}}$ , they are evaluated at time  $t^{n+1/2} := t^n + \delta t/2$ . The construction of these quantities is motivated by a semi-discretization of the governing equations (2) in time (cf. [VATER,



2005]). Let us suppose for a moment a sufficiently smooth solution. By Taylor series expansion, height and momentum can be expressed at the new time level by

$$h(\mathbf{x}, t^{n+1}) = h(\mathbf{x}, t^n) - \delta t [\nabla \cdot (h\mathbf{v})(\mathbf{x}, t^{n+1/2})] + \mathcal{O}(\delta t^3) \quad (7)$$

and

$$\begin{aligned} (h\mathbf{v})(\mathbf{x}, t^{n+1}) &= (h\mathbf{v})(\mathbf{x}, t^n) - \delta t [\nabla \cdot (h\mathbf{v} \circ \mathbf{v})(\mathbf{x}, t^{n+1/2}) \\ &\quad + (h_0 \nabla h^{(2)})(\mathbf{x}, t^{n+1/2})] + \mathcal{O}(\delta t^3) . \end{aligned} \quad (8)$$

Assuming that appropriate approximations of the fluxes of the auxiliary system (6) have been computed, height and momentum are given at the intermediate time level by

$$\begin{aligned} (h\mathbf{v})(\mathbf{x}, t^{n+1/2}) &= (h\mathbf{v})^*(\mathbf{x}, t^{n+1/2}) - \frac{\delta t}{2} (h_0 \nabla h^{(2)})(\mathbf{x}, t^{n+1/4}) + \mathcal{O}(\delta t^3) \\ \mathbf{v}(\mathbf{x}, t^{n+1/2}) &= \mathbf{v}^*(\mathbf{x}, t^{n+1/2}) - \frac{\delta t}{2} \nabla h^{(2)}(\mathbf{x}, t^{n+1/4}) + \mathcal{O}(\delta t^3) . \end{aligned} \quad (9)$$

Here, the variables with stars denote the quantities of the auxiliary system. Note that – in order to achieve second order accuracy in time – the question at which time level the unknown  $h^{(2)}$  “lives”, can be relaxed to any point in the interval  $[t^n, t^{n+1/2}]$ . To ensure that the velocities on the left hand side of (9) satisfy the divergence constraint, we take the divergence of the first equation and obtain a first Poisson equation for  $h^{(2)}$ :

$$\frac{\delta t}{2} \nabla \cdot (h_0 \nabla h^{(2)})(\mathbf{x}, t^{n+1/4}) = \nabla \cdot (h\mathbf{v})^*(\mathbf{x}, t^{n+1/2}) + \frac{dh_0}{dt}(t^{n+1/2}) + \mathcal{O}(\delta t^3) . \quad (10)$$

With the solution of this problem the right hand side of (7) and the first term in the brackets of (8) can be calculated through (9). The second term in brackets in (8) is computed by another application of a discrete divergence constraint. Let

$$(h\mathbf{v})^{**}(\mathbf{x}) := (h\mathbf{v})(\mathbf{x}, t^n) - \delta t [\nabla \cdot (h\mathbf{v} \circ \mathbf{v})(\mathbf{x}, t^{n+1/2})] \quad (11)$$

denote a preliminary prediction of the new time level momentum that still lacks the influence of the pressure flux. Then, the momentum at the new time level is given by

$$(h\mathbf{v})(\mathbf{x}, t^{n+1}) = (h\mathbf{v})^{**}(\mathbf{x}) - \delta t (h_0 \nabla h^{(2)})(\mathbf{x}, t^{n+1/2}) + \mathcal{O}(\delta t^3) . \quad (12)$$

Imposing the divergence constraint once again at a half time step, but this time using a linear interpolation of the momentum at the full time levels, leads to

$$\frac{1}{2} [\nabla \cdot (h\mathbf{v})(\mathbf{x}, t^{n+1}) + \nabla \cdot (h\mathbf{v})(\mathbf{x}, t^n)] = -\frac{dh_0}{dt}(t^{n+1/2}) + \mathcal{O}(\delta t^2) . \quad (13)$$

Inserting (12) in (13), a second Poisson Problem for  $h^{(2)}$  is obtained:

$$\begin{aligned} \delta t \nabla \cdot (h_0 \nabla h^{(2)})(\mathbf{x}, t^{n+1/2}) &= \nabla \cdot (h\mathbf{v})^{**}(\mathbf{x}) + \nabla \cdot (h\mathbf{v})(\mathbf{x}, t^n) \\ &\quad + 2 \frac{dh_0}{dt}(t^{n+1/2}) + \mathcal{O}(\delta t^2) . \end{aligned} \quad (14)$$

Thus, by the solution of an auxiliary hyperbolic system and two Poisson problems for the incompressible height  $h^{(2)}$  numerical approximations to the fluxes of the zero Froude number shallow water equations can be computed up to second order accuracy in time.

## 2.2 Discretization of the Projections

As stated above, equations (6)–(14) are a summary of the zero-Mach-number-scheme by SCHNEIDER ET AL. [1999] applied to the shallow water case. In this section we begin to introduce deviations from earlier work. This concerns, in particular, a new discretization of the Poisson equation, which – as we will show – leads to an exact and stable projection.

The Poisson equations (10) and (14) are discretized using a method originally proposed by SÜLI [1991], who proves stability and convergence of the scheme in a mesh-dependent  $H^1$  norm. In contrast to Süli, who considers a numerical method for a scalar elliptic Dirichlet problem, we focus here on the projection step of a flow solver that results in a Poisson-type problem with von Neumann boundary conditions (cf. GRESHO [1990] and SCHNEIDER ET AL. [1999] for a discussion on that issue). The method can be either interpreted as a finite element or as a finite volume method. In the following, the scheme is introduced as a Petrov-Galerkin finite element method, which facilitates the stability proof of the projection given in the next section. Since the two Poisson equations are solved using slightly different discretizations, the method is first discussed for the second projection. Thereafter, modifications to be applied for the first Poisson problem are given.

For the derivation of the method, consider a Poisson problem with von Neumann boundary conditions:

$$\begin{cases} -\nabla \cdot \nabla p = f & \text{in } \Omega, \\ \frac{\partial p}{\partial \mathbf{n}} = 0 & \text{on } \partial\Omega. \end{cases} \quad (15)$$

Given the r.h.s.  $f \in L^2(\Omega)$  with  $\int_{\Omega} f \, d\mathbf{x} = 0$ , this problem has a unique solution  $p \in H^1(\Omega)/\mathbb{R}$ . Since the right hand side  $f$  is of the form  $-\nabla \cdot \mathbf{v}$  with a given velocity field  $\mathbf{v}$  in the equation to be solved in the projection method,  $f$  is substituted with this term in the following discussion. The weak formulation of this problem is derived by multiplication of (15) with a test function  $\psi$  and integration over the whole domain  $\Omega$ . Thus, we have to find  $p$ , such that

$$\int_{\Omega} \psi \nabla \cdot \nabla p \, d\mathbf{x} = \int_{\Omega} \psi \nabla \cdot \mathbf{v} \, d\mathbf{x} \quad \forall \psi \quad . \quad (16)$$

In (16) it is left open which trial and test spaces are considered. In contrary to the classical finite element theory, where the test function  $\psi$  is chosen to be (weakly) differentiable and Green's formula is applied to shift one derivative to the test function, here, a test space containing piecewise constant test functions is considered. In this case –

assuming for a moment that  $p$  and  $\mathbf{v}$  are sufficiently smooth – the divergence theorem can be applied.

In particular, for the construction of the test space, a *dual discretization* of the computational domain  $\Omega$  is introduced, in which  $\bar{\mathcal{C}}$  is the set of control volumes  $\bar{C}$  centered about nodes of the original grid (see Figure 1). Notice that usage of the dual cells in formulating the projection is in line with BELL and MARCUS [1992]; SCHNEIDER ET AL. [1999]. The difference will lie in how we account for piecewise linear in-cell distributions of momentum and how they are affected by the divergence correction. The interfaces between these control volumes – and the set of all such interfaces is denoted – in analogy to the primal discretization – by  $\bar{I}$  and  $\bar{T}$ , respectively. Then, the test space is given by all functions in  $L^2(\Omega)$ , which are constant on the dual control volumes. This space can be defined by

$$\mathcal{Q}^h := \{q \in L^2(\Omega) \mid \forall \bar{C} \in \bar{\mathcal{C}} : q|_{\bar{C}} \in \mathcal{P}_0(\bar{C})\} , \quad (17)$$

in which

$$\mathcal{P}_k(U) := \left\{ p \in C^\infty(U) \mid p(x, y) = \sum_{\substack{i+j \leq k \\ i, j \geq 0}} c_{ij} x^i y^j, c_{ij} \in \mathbb{R} \right\} \quad (18)$$

is the space of polynomial functions on  $U \subset \mathbb{R}^2$  of degree less than or equal to  $k$ . A basis of  $\mathcal{Q}^h$  is given by  $\bigcup_{\bar{C} \in \bar{\mathcal{C}}} \{\chi_{\bar{C}}\}$ , where  $\chi_U$  is the characteristic function on the set  $U$ . Therefore, a test function can be decomposed into  $\psi(x, y) = \sum_{\bar{C}} \psi_{\bar{C}} \chi_{\bar{C}}(x, y)$ , and equation (16) becomes

$$\sum_{\bar{C} \in \bar{\mathcal{C}}} \psi_{\bar{C}} \left( \int_{\bar{C}} \nabla \cdot \nabla p \, d\mathbf{x} - \int_{\bar{C}} \nabla \cdot \mathbf{v} \, d\mathbf{x} \right) = 0 \quad \forall \psi \in \mathcal{Q}^h .$$

Now, the divergence theorem can be applied, and we have to find  $p$ , such that

$$\sum_{\bar{C} \in \bar{\mathcal{C}}} \psi_{\bar{C}} \left( \int_{\partial \bar{C}} \nabla p \cdot \mathbf{n} \, d\sigma - \int_{\partial \bar{C}} \mathbf{v} \cdot \mathbf{n} \, d\sigma \right) = 0 \quad \forall \psi \in \mathcal{Q}^h , \quad (19)$$

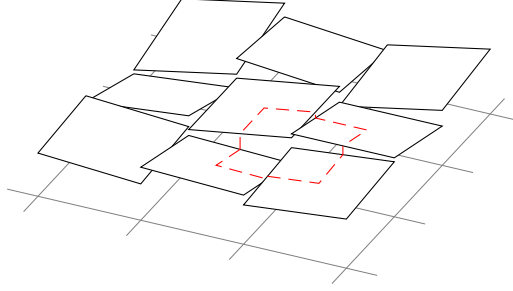
Since all of the  $\bar{C}$  are pairwise disjoint, this problem is a linear combination of the local problems to find  $p$ , such that

$$\int_{\partial \bar{C}} \nabla p \cdot \mathbf{n} \, d\sigma - \int_{\partial \bar{C}} \mathbf{v} \cdot \mathbf{n} \, d\sigma = 0 \quad \forall \bar{C} \in \bar{\mathcal{C}} , \quad (20)$$

and the solution  $p$  satisfies (19), if and only if it satisfies (20).

Using the latter formulation, the trial spaces for the unknown  $p$  and the right hand side  $\mathbf{v}$  are now defined as follows: Choosing for  $p$  a trial space of continuous functions, which are piecewise bilinear on the primal control volumes  $C \in \mathcal{C}$ , i.e.

$$\mathcal{H}^h := \{p \in H^1(\Omega) \mid \forall C \in \mathcal{C} : p|_C \in \mathcal{P}_2(C), \forall I \in \mathcal{I} : p|_I \in \mathcal{P}_1(I)\} , \quad (21)$$



**Figure 2:** Piecewise linear functions for the velocity. The red dashed line visualizes the integration path of the boundary integral, which is evaluated in the discrete divergence.

the gradient of such functions is piecewise linear in each component on a control volume of the primal discretization, but discontinuous at the interfaces. Thus, for the velocity vector  $\mathbf{v}$  a finite element space is chosen, which contains such gradients. It is defined by

$$\mathcal{U}^h := \{ \mathbf{v} = (u, v) \in [L^2(\Omega)]^2 \mid \forall C \in \mathcal{C} : \mathbf{v}|_C \in [\mathcal{P}_1(C)]^2 \} . \quad (22)$$

Note, that, although this space allows for discontinuities along cell interfaces, all the integrals in (20) are well defined. This is true, because the normal component of  $\mathbf{v}$  and  $\nabla p$  are piecewise linear along the boundaries of the dual control volumes (cf. Figure 2), and the expressions can be exactly evaluated. Notice also that piecewise linear velocity or momentum components are the natural ansatz to obtain a second order Godunov-type scheme used in the explicit predictor step.

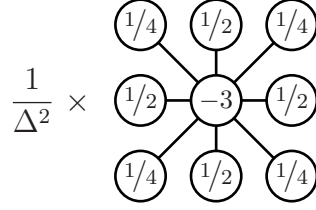
Using a suitable normalization, the integrals on the left hand side of (20) define a discrete Laplacian and divergence. Specifically, let us define the discrete Laplacian by

$$\mathbb{L} : \mathcal{H}^h \rightarrow \mathcal{Q}^h \text{ with } \mathbb{L}(p) := \sum_{\bar{C} \in \bar{\mathcal{C}}} \chi_{\bar{C}} \frac{1}{|\bar{C}|} \int_{\partial \bar{C}} \nabla p \cdot \mathbf{n} \, d\sigma \quad (23)$$

and the discrete divergence by

$$\mathbb{D} : \mathcal{U}^h \rightarrow \mathcal{Q}^h \text{ with } \mathbb{D}(\mathbf{v}) := \sum_{\bar{C} \in \bar{\mathcal{C}}} \chi_{\bar{C}} \frac{1}{|\bar{C}|} \int_{\partial \bar{C}} \mathbf{v} \cdot \mathbf{n} \, d\sigma . \quad (24)$$

Since each basis function of the test space is only nonzero on one dual control volume, the resulting stencil of the Laplacian is compact, i.e. it only uses next neighbors to the grid point for which the differential operator is discretized. As a consequence, the associated linear system can be easily computed with standard iterative methods. On a uniform Cartesian grid with the same grid spacing in both coordinate directions the stencil is given in Figure 3.



**Figure 3:** Stencil of the discrete Laplacian on a uniform Cartesian grid with the same grid spacing  $\Delta$  in both coordinate directions.

The property that the analytical gradient of  $p \in \mathcal{H}^h$  is in the space  $\mathcal{U}^h$  almost everywhere suggests that the discrete gradient operator is defined by

$$\mathbf{G} : \mathcal{H}^h \rightarrow \mathcal{U}^h \text{ with } \mathbf{G}(p) := \nabla p \quad \text{a.e.} \quad (25)$$

These discrete operators inherit from their analytic counterparts the property that they satisfy the equality  $\mathbf{L} = \mathbf{D}(\mathbf{G})$ .

The discretization of the first projection is done in a similar way. However, this time the advection velocity has to be corrected at the boundary of the primary control volumes. Thus, the test functions are chosen to be piecewise constant on each grid cell, which means that the divergence is applied to each such control volume (see Figure 4). On a Cartesian grid, the discretization is essentially shifted by half a grid cell in each coordinate direction. The resulting flux, arising from the MAC projection, is given by

$$\mathbf{F}_I^{\text{MAC}} = -\frac{\delta t}{2} \left( \begin{array}{c} h_0 \nabla h^{(2)} \cdot \mathbf{n} \\ (h\mathbf{v})^* \nabla h^{(2)} \cdot \mathbf{n} + h_0 \nabla h^{(2)} \mathbf{v}^* \cdot \mathbf{n} \end{array} \right)_I.$$

In the second projection, the local updates of the momentum are given by

$$(h\mathbf{v})^{n+1}|_C = (h\mathbf{v})^{**}(\mathbf{x})|_C - \delta t h_0 \nabla h^{(2)}(\mathbf{x})|_C \quad C \in \mathcal{C}$$

(cf. (12)). This results in the flux contribution  $\mathbf{F}_I^{\text{P2}} = (0, h_0 h^{(2)} \mathbf{n})_I^T$ , and conservation of momentum is guaranteed.

We emphasize that the update of the second projection not only involves the cell mean values, but also the gradient within a cell. This can be seen by a decomposition of the quantities into their mean value, linear and bilinear fractions, i.e.:

$$h^{(2)}(x, y)|_C = h_C^{(2)} + (x - x_C)h_{x,C}^{(2)} + (y - y_C)h_{y,C}^{(2)} + (x - x_C)(y - y_C)h_{xy,C}^{(2)},$$

where  $(x_C, y_C)$  is the center of cell  $C$ . Then, the gradient in each grid cell is given by

$$\nabla h^{(2)}(x, y)|_C = \begin{pmatrix} h_{x,C}^{(2)} \\ h_{y,C}^{(2)} \end{pmatrix} + \begin{pmatrix} y - y_C \\ x - x_C \end{pmatrix} h_{xy,C}^{(2)},$$

and the update of the mean values is given by

$$(h\mathbf{v})_C^{n+1} = (h\mathbf{v})_C^{**} - \delta t h_0 \begin{pmatrix} h_{x,C}^{(2)} \\ h_{y,C}^{(2)} \end{pmatrix},$$

whereas the correction of the gradients is computed by

$$(h\mathbf{v})_{x,C}^{n+1} = (h\mathbf{v})_{x,C}^{**} - \delta t h_0 \begin{pmatrix} 0 \\ h_{xy,C}^{(2)} \end{pmatrix}$$

and

$$(h\mathbf{v})_{y,C}^{n+1} = (h\mathbf{v})_{y,C}^{**} - \delta t h_0 \begin{pmatrix} h_{xy,C}^{(2)} \\ 0 \end{pmatrix}.$$

Additionally, a reconstruction step is introduced after the first projection, which reconstructs piecewise linear functions from cell averages of the intermediate momentum components  $(hu)_C^{**}$  and  $(hv)_C^{**}$ . The second projection is then applied to this vector field to obtain a final momentum distribution. Note that the total variation diminishing (TVD) property could be destroyed in the projection step, even if it was installed in the reconstruction step before.

### 2.3 Exact Projection Method

Using the discretization described above for the second Poisson equation, the numerical method can be formulated as an *exact projection method*. This means that the incompressibility condition on the velocity

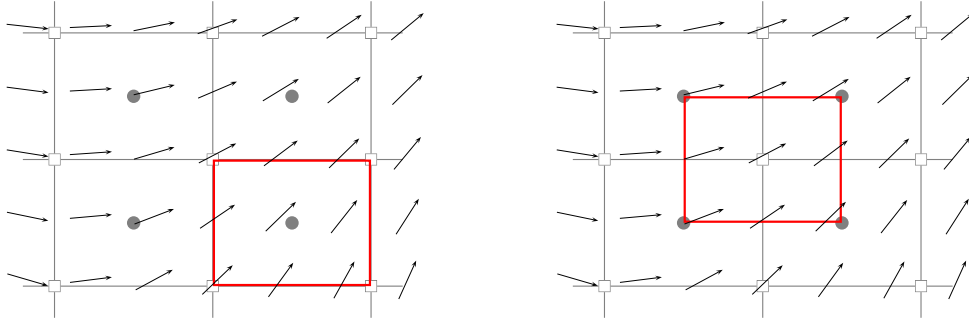
$$(\nabla \cdot \mathbf{v}^n)_{\bar{C}} := \frac{1}{\bar{C}} \int_{\partial \bar{C}} \mathbf{v} \cdot \mathbf{n} d\sigma = -\frac{1}{h_0} \frac{dh_0}{dt}$$

is theoretically satisfied to machine precision at each full time level (i.e. in practice to the precision of the iterative solver for the discrete Poisson equation). As noted above, this definition of the divergence not only incorporates the cell mean values, but also the gradients of the velocity within each cell intersecting  $\bar{C}$ .

To derive an exact projection method the piecewise linear functions for the momentum have to be used throughout the whole scheme. In the semi-discrete implementation for the solution of the auxiliary system Heun's method

$$\begin{aligned} \mathbf{U}^* &= \mathbf{U}^n + \frac{\delta t}{2} (f(\mathbf{U}^n) + f(\mathbf{U}^{*,\text{int}})) \quad \text{with} \\ \mathbf{U}^{*,\text{int}} &= \mathbf{U}^n + \delta t f(\mathbf{U}^n) \end{aligned}$$

is applied for the integration in time. This approach leads to second-order accuracy in time. To obtain second-order accuracy in space as well, the cell average values in  $\mathbf{U}^n$  and  $\mathbf{U}^{*,\text{int}}$  are reconstructed as piecewise linear functions on each grid cell. The



**Figure 4:** Application of the divergence constraint in the MAC (left) and the second projection (right).

numerical fluxes are then evaluated with the reconstructed values on the two sides of any particular interface.

Since the momentum components are already piecewise linear at time level  $t^n$ , they do not have to be reconstructed from the cell mean values and the gradients of the momentum components are used for the calculation of the numerical fluxes of the auxiliary system. These gradients are not only used for  $\mathbf{U}^n$ , but for  $\mathbf{U}^{*,\text{int}}$  as well. This does not reduce the scheme's order, because a Taylor series expansion for the the gradient of  $\mathbf{U}^{*,\text{int}}$  yields

$$\mathbf{U}_{\mathbf{x},C}^{*,\text{int}} = \mathbf{U}_{\mathbf{x},C}^n + \mathcal{O}(\delta t) .$$

In this scheme  $\mathbf{U}_{\mathbf{x},C}$  is always multiplied by  $\delta x$  to yield the numerical fluxes of the auxiliary system. Therefore, the second order accuracy in space and time is retained.

With these modifications, we have a velocity field at each time level, which satisfies the discrete divergence constraint up to the accuracy of the elliptic solver, i.e. we have constructed an exact projection method.

### 3 Stability of the second projection

In proving stability of our semi-implicit method, the stability of the second projection step is an important prerequisite. Furthermore, as stated in the introduction, the final projection often led to a velocity-pressure decoupling in former projection methods. By using the theory of mixed finite element methods, we demonstrate that such instabilities cannot occur in the presented method.

In the second projection, the height perturbation  $h^{(2)}$  is computed to correct the intermediate momentum update  $(h\mathbf{v})^{**}$  in a post-processing step (cf. (12)). Thus, we are not only interested in a stable approximation of  $h^{(2)}$ , but rather in one of the momentum at the new time step. The associated Poisson-type problem is derived by imposing the additional requirement that the momentum at the new time step shall

satisfy a discrete version of the divergence constraint

$$\int_{\partial V} (h\mathbf{v}) \cdot \mathbf{n} \, d\sigma = -|V| \frac{dh_0}{dt} . \quad (26)$$

In the context of finite element methods, this leads to the theory of *saddle point problems* (mixed finite elements), which arise from minimization problems with additional side constraints. Starting with the fundamental work of BABUŠKA [1971] and BREZZI [1974], this theory provides conditions for existence and uniqueness of solutions and for stable discretizations of such problems.

After having introduced the fundamental functional analytic framework, the discrete Poisson-type problem

$$\delta t \, \mathbf{D} \left( h_0 \mathbf{G}(h^{(2)}) \right) = \mathbf{D}((h\mathbf{v})^{**}) + \mathbf{D}((h\mathbf{v})^n) + 2 \frac{dh_0}{dt} \quad (27)$$

is reformulated for the new projection method as a generalized saddle point problem, which is the starting point for the subsequent stability analysis.

### 3.1 Generalized Saddle Point Problems – Theory

For simplicity it is always assumed that  $\Omega$  is a bounded open subset of  $\mathbb{R}^n$ , which is connected and has a Lipschitz-continuous boundary  $\partial\Omega$ . The theory of finite element methods heavily benefits from the utilization of *Sobolev spaces*. These are based on the Hilbert space  $L^2(\Omega)$ , which includes all square integrable functions on  $\Omega$ . The latter is defined by

$$L^2(\Omega) := \left\{ q \mid \int_{\Omega} |q(\mathbf{x})|^2 \, d\mathbf{x} < +\infty \right\} ,$$

and both, an inner product and a norm on this space are given by

$$(p, q)_{0,\Omega} := \int_{\Omega} p(\mathbf{x})q(\mathbf{x}) \, d\mathbf{x} , \quad \|q\|_{0,\Omega} := \sqrt{(q, q)_{0,\Omega}} .$$

Then, the first order Sobolev space is

$$H^1(\Omega) := \{ q \in L^2(\Omega) \mid \nabla q \in [L^2(\Omega)]^n \} .$$

We put  $|q|_{1,\Omega} := \|\nabla q\|_{0,\Omega}$  and  $\|q\|_{1,\Omega} := (\|q\|_{0,\Omega}^2 + |q|_{1,\Omega}^2)^{1/2}$ , which define a semi-norm and a norm on  $H^1(\Omega)$ , respectively. Note that  $|\cdot|_{1,\Omega}$  defines a norm on the quotient space  $H^1(\Omega)/\mathbb{R}$ , the space of equivalence classes of functions that differ only by a constant. We also refer to spaces of vector valued functions. For this reason, let us introduce

$$H(\text{div}; \Omega) := \{ \mathbf{v} \in [L^2(\Omega)]^n \mid \nabla \cdot \mathbf{v} \in L^2(\Omega) \} .$$



For a vector function  $\mathbf{v} \in H(\operatorname{div}; \Omega)$  it is possible to define its normal component on the boundary  $\partial\Omega$  [GIRAULT and RAVIART, 1986], and the subspace with vanishing normal component on  $\partial\Omega$  is denoted by

$$H_0(\operatorname{div}; \Omega) := \{\mathbf{v} \in H(\operatorname{div}; \Omega) \mid \mathbf{v} \cdot \mathbf{n} = 0 \text{ on } \partial\Omega\} .$$

These spaces are equipped with the Hilbertian graph norm

$$\|\mathbf{v}\|_{\operatorname{div}, \Omega} := \left( \|\mathbf{v}\|_{0, \Omega}^2 + \|\nabla \cdot \mathbf{v}\|_{0, \Omega}^2 \right)^{1/2} .$$

For the analysis of the second projection we are interested in generalized mixed formulations with three distinct bilinear forms  $a$ ,  $b_1$ ,  $b_2$ . That is, to find  $(u, p) \in \mathcal{U} \times \mathcal{H}$ , such that

$$\begin{cases} a(u, v) + b_1(p, v) = \langle v', v \rangle & \forall v \in \mathcal{V} \\ b_2(u, q) = \langle q', q \rangle & \forall q \in \mathcal{Q} . \end{cases} \quad (28)$$

In this formulation,  $\mathcal{H}$ ,  $\mathcal{Q}$ ,  $\mathcal{U}$  and  $\mathcal{V}$  are four Hilbert spaces (or, more generally, reflexive Banach spaces) with norms  $\|\cdot\|_{\mathcal{H}}$ ,  $\|\cdot\|_{\mathcal{Q}}$ ,  $\|\cdot\|_{\mathcal{U}}$  and  $\|\cdot\|_{\mathcal{V}}$ . The bilinear form  $a$  is defined on  $\mathcal{U} \times \mathcal{V}$ ,  $b_1$  on  $\mathcal{H} \times \mathcal{V}$  and  $b_2$  on  $\mathcal{U} \times \mathcal{Q}$ . Furthermore,  $v'$  and  $q'$  are elements of  $\mathcal{V}'$  and  $\mathcal{Q}'$ , the dual spaces of  $\mathcal{V}$  and  $\mathcal{Q}$ . The abstract theory of such problems is given in NICOLAÏDES [1982] and developed further in BERNARDI ET AL. [1988].

To obtain conditions for existence, uniqueness and stability of problem (28), let us introduce for any  $r' \in \mathcal{H}'$  and  $q' \in \mathcal{Q}'$  the closed affine spaces

$$\mathcal{K}_1(r') := \{v \in \mathcal{V} \mid \forall r \in \mathcal{H} : b_1(r, v) = \langle r', r \rangle\}$$

and

$$\mathcal{K}_2(q') := \{w \in \mathcal{U} \mid \forall q \in \mathcal{Q} : b_2(w, q) = \langle q', q \rangle\} .$$

We denote by  $\mathcal{K}_i := \mathcal{K}_i(0)$  ( $i = 1, 2$ ) the kernel of the operator induced by  $b_i$ . With these definitions the following Theorem can be stated:

**Theorem 3.1** (NICOLAÏDES [1982]) *Let  $a(\cdot, \cdot)$  and  $b_i(\cdot, \cdot)$  ( $i = 1, 2$ ) be bounded. Assume that there exists a constant  $\alpha > 0$ , such that*

$$\inf_{w \in \mathcal{K}_2} \sup_{v \in \mathcal{K}_1} \frac{a(w, v)}{\|w\|_{\mathcal{U}} \|v\|_{\mathcal{V}}} \geq \alpha \quad (29)$$

and

$$\sup_{w \in \mathcal{K}_2} a(w, v) > 0 \quad \forall v \in \mathcal{K}_1 \setminus \{0\} \quad . \quad (30)$$

Furthermore, assume that the  $b_i(\cdot, \cdot)$  ( $i = 1, 2$ ) satisfy the inf-sup conditions

$$\inf_{r \in \mathcal{H}} \sup_{v \in \mathcal{V}} \frac{b_1(r, v)}{\|r\|_{\mathcal{H}} \|v\|_{\mathcal{V}}} \geq \beta_1 > 0 \quad (31)$$

and

$$\inf_{q \in \mathcal{Q}} \sup_{w \in \mathcal{U}} \frac{b_2(w, q)}{\|w\|_{\mathcal{U}} \|q\|_{\mathcal{Q}}} \geq \beta_2 > 0 \quad . \quad (32)$$

Then, problem (28) has a unique solution  $(u, p)$  for all  $v' \in \mathcal{V}'$  and  $q' \in \mathcal{Q}'$  and the following estimate holds:

$$\|u\|_{\mathcal{U}} + \|p\|_{\mathcal{H}} \leq c \left( \|v'\|_{\mathcal{V}'} + \|q'\|_{\mathcal{Q}'} \right) \quad . \quad (33)$$

For the discretization of problem (28), it is assumed that there are finite-dimensional subspaces  $\mathcal{H}^h \subset \mathcal{H}$ ,  $\mathcal{Q}^h \subset \mathcal{Q}$ ,  $\mathcal{U}^h \subset \mathcal{U}$  and  $\mathcal{V}^h \subset \mathcal{V}$  and bilinear forms  $a_h : \mathcal{U}^h \times \mathcal{V}^h \rightarrow \mathbb{R}$ ,  $b_{1h} : \mathcal{H}^h \times \mathcal{V}^h \rightarrow \mathbb{R}$  and  $b_{2h} : \mathcal{U}^h \times \mathcal{Q}^h \rightarrow \mathbb{R}$ . Given the linear functionals  $v'_h \in (\mathcal{V}^h)'$  and  $q'_h \in (\mathcal{Q}^h)'$ , we are looking for the solution  $(u_h, p_h) \in \mathcal{U}^h \times \mathcal{H}^h$  of the discrete problem

$$\begin{cases} a_h(u_h, v_h) + b_{1h}(p_h, v_h) = \langle v'_h, v_h \rangle & \forall v_h \in \mathcal{V}^h \\ b_{2h}(u_h, q_h) = \langle q'_h, q_h \rangle & \forall q_h \in \mathcal{Q}^h \end{cases} \quad (34)$$

approximating the solution of the continuous problem. With the definition of the discrete affine spaces  $\mathcal{K}_1^h$  and  $\mathcal{K}_2^h$ , in analogy to the continuous case, Theorem 3.1 can be applied to problem (34), and existence, uniqueness and stability are obtained given the constants  $\alpha$ ,  $\beta_1$  and  $\beta_2$  in (29), (31) and (32) are independent of the grid parameter  $h$ . Examples of mixed finite element discretizations of such type are given in NICOLAÏDES [1982] and BERNARDI ET AL. [1988]. A nonconforming discretization, where  $\mathcal{U}^h \not\subset \mathcal{U}$ , is constructed in ANGERMANN [2003]. Moreover, error estimates are provided in these references for both, the conforming and the nonconforming situation.

In the following, such a formulation is derived for the new projection in order to analyze its stability concerning the corrected momentum field.

### 3.2 Reformulation of the problem

The derivation of a mixed formulation equivalent to the Poisson-type problem (27) is easily established. The continuous counterpart of this equation is obtained by a combination of the momentum update and the divergence constraint, i.e.,

$$\begin{aligned} (h\mathbf{v})^{n+1} &= (h\mathbf{v})^{**} - \delta t (h_0 \nabla h^{(2)}) \\ \frac{1}{2} [\nabla \cdot (h\mathbf{v})^{n+1} + \nabla \cdot (h\mathbf{v})^n] &= -\frac{dh_0}{dt} \quad . \end{aligned} \quad (35)$$

A variational formulation of these two equations is derived by the usual procedure: (35)<sub>1</sub> and (35)<sub>2</sub> are multiplied with test functions  $\varphi$  and  $\psi$ , and the resulting equations are integrated over the whole domain  $\Omega$ . This leads to

$$\begin{aligned} ((h\mathbf{v})^{n+1}, \varphi)_{0,\Omega} + \left( \delta t h_0 \nabla h^{(2)}, \varphi \right)_{0,\Omega} &= ((h\mathbf{v})^{**}, \varphi)_{0,\Omega} \\ (\nabla \cdot (h\mathbf{v})^{n+1}, \psi)_{0,\Omega} &= - \left( \nabla \cdot (h\mathbf{v})^n + 2 \frac{dh_0}{dt}, \psi \right)_{0,\Omega} \quad . \end{aligned} \quad (36)$$

Note that this formulation can be already interpreted as a generalized problem as formulated in (28). The discrete method, equivalent to the Poisson-type problem (27), is derived by introducing appropriate finite dimensional trial and test spaces. For the choice of the trial spaces, we are confined to our selection for the momentum ( $h\mathbf{v}$ ) and the height  $h^{(2)}$ . In the projection method, the momentum distribution is approximated by discontinuous piecewise linear functions belonging to the space  $\mathcal{U}^h$  defined in (22). The height perturbation  $h^{(2)} \in \mathcal{H}^h$  is given by continuous piecewise bilinear functions (cf. (21)).

To obtain the same divergence as in (27), also the test functions  $\psi$  for the second equation of (36) are fixed to be piecewise constant on dual control volumes, forming the space  $\mathcal{Q}^h$  defined in (17). The selection of the test space  $\mathcal{V}^h$  for the first equation is yet undetermined. Let us choose  $\mathcal{V}^h = \mathcal{U}^h$ , the space which is also used for the momentum variable. A basis of  $\mathcal{V}^h$  is given by

$$\bigcup_{C \in \mathcal{C}} \left\{ \begin{pmatrix} \chi_C \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ \chi_C \end{pmatrix}, \begin{pmatrix} (x - x_C)\chi_C \\ 0 \end{pmatrix}, \begin{pmatrix} (y - y_C)\chi_C \\ 0 \end{pmatrix}, \right. \\ \left. \begin{pmatrix} 0 \\ (x - x_C)\chi_C \end{pmatrix}, \begin{pmatrix} 0 \\ (y - y_C)\chi_C \end{pmatrix} \right\}, \quad (37)$$

where  $(x_C, y_C)$  is the center of the cell  $C$ .

The following discussion is focused on Cartesian grids with grid cells  $C_{i,j}$ ,  $i = 1, \dots, m$ ,  $j = 1, \dots, n$ , and cell centers  $(x_i, y_j)$ . Because of the linearity of the equations (36) in  $\varphi$  and  $\psi$ , it is sufficient to “test” them with only a basis of  $\mathcal{U}^h$  and  $\mathcal{Q}^h$ , respectively. Let us consider the first equation in conjunction with the test function  $\varphi = (\chi_{C_{i,j}}, 0)^T$ . Because the second component of  $\varphi$  is zero and its support is  $C_{i,j}$ , this yields

$$\int_{C_{i,j}} (hu)^{n+1} d\mathbf{x} + \delta t h_0 \int_{C_{i,j}} \frac{\partial h^{(2)}}{\partial x} d\mathbf{x} = \int_{C_{i,j}} (hu)^{**} d\mathbf{x} \quad . \quad (38)$$

Furthermore, by expanding the height  $h^{(2)}$  in a volumewise representation, i.e.

$$h^{(2)}(x, y)|_{C_{i,j}} = h_{i,j}^{(2)} + (x - x_i)h_{x,i,j}^{(2)} + (y - y_j)h_{y,i,j}^{(2)} + (x - x_i)(y - y_j)h_{xy,i,j}^{(2)} \quad , \quad (39)$$

the calculation of the second integral in (38) leads to

$$\int_{C_{i,j}} \frac{\partial h^{(2)}}{\partial x} d\mathbf{x} = \int_{C_{i,j}} \left( h_{x,i,j}^{(2)} + (y - y_j)h_{xy,i,j}^{(2)} \right) d\mathbf{x} = \delta x \delta y h_{x,i,j}^{(2)} \quad .$$

The integral of the second term vanishes, because it is an odd function in  $y$  with respect to  $y_j$ . With similar results for the other terms in (38), we finally obtain

$$(hu)_{i,j}^{n+1} + \delta t h_0 h_{x,i,j}^{(2)} = (hu)_{i,j}^{**} \quad . \quad (40)$$

By using the other five test functions in (37), this procedure yields the equations

$$\begin{aligned}
(hv)_{i,j}^{n+1} + \delta t h_0 h_{y,i,j}^{(2)} &= (hv)_{i,j}^{**} \\
(hu)_{x,i,j}^{n+1} &= (hu)_{x,i,j}^{**} \\
(hu)_{y,i,j}^{n+1} + \delta t h_0 h_{xy,i,j}^{(2)} &= (hu)_{y,i,j}^{**} \\
(hv)_{x,i,j}^{n+1} + \delta t h_0 h_{xy,i,j}^{(2)} &= (hv)_{x,i,j}^{**} \\
(hv)_{y,i,j}^{n+1} &= (hv)_{y,i,j}^{**} .
\end{aligned} \tag{41}$$

Therefore, six equations are obtained for each cell  $C_{i,j}$ . They represent the discretization of (36)<sub>1</sub>.

The discretization of the second equation in (36) is done as follows. The application of the test function  $\psi = \chi_{\bar{C}}$  and the divergence theorem yields for the terms involving the momentum the key ingredient of the discrete divergence  $D(\cdot)$ . Thus, multiplying this equation by  $\chi_{\bar{C}}/|\bar{C}|$  and summation over  $\bar{C} \in \bar{\mathcal{C}}$  leads to

$$D((h\mathbf{v})^{n+1}) = -D((h\mathbf{v})^n) - 2 \frac{dh_0}{dt} . \tag{42}$$

Let us recall that  $h^{(2)}$  is uniquely defined by its nodal values and that each velocity component has three degrees of freedom per grid cell. Then there are  $7 \cdot m \cdot n$  unknowns in case of periodic boundary conditions, where  $m$  and  $n$  are the number of cells in  $x$  and  $y$  direction, respectively. The analysis above yielded the same number of linear equations. Finally, by inserting the equations from (40) and (41) into (42), the second discrete Poisson-type problem from our new projection method is obtained. We have derived a *Petrov-Galerkin* mixed formulation, which utilizes different trial and test spaces for the scalar variables.

### 3.3 Stability analysis of the mixed formulation

In order to apply the theory from Section 3.1 to the mixed formulation (36), the corresponding continuous problem is defined which has been shown to have a unique solution in VATER [2005]. Here, the main investigation will be on the stability of the discrete mixed formulation.

For the analysis of the continuous problem appropriate function spaces for the trial and test functions have to be chosen. In the Poisson-type problem

$$\delta t \nabla \cdot (h_0 \nabla h^{(2)}) = \nabla \cdot (h\mathbf{v})^{**} + \nabla \cdot (h\mathbf{v})^n + 2 \frac{dh_0}{dt}$$

– the continuous counterpart of (27) – the height perturbation  $h^{(2)}$  is only determined up to an additive constant. This constant can be fixed by the additional condition of

a zero mean value, i.e.,  $\int_{\Omega} h^{(2)} d\mathbf{x} = 0$ . Thus, a suitable trial space for  $h^{(2)}$  is given by  $\mathcal{H} := H^1(\Omega)/\mathbb{R}$ . An appropriate space for the momentum should also bound the divergence of the unknown variable. Furthermore, the boundary conditions are given by the integral constraint (26). For simplicity, let us assume that there is no flux across the boundary, i.e., impermeable rigid walls and  $dh_0/dt \equiv 0$ . Then, the momentum is sought in the space  $\mathcal{U} = H_0(\text{div}; \Omega)$ . The test functions of the discrete problem are discontinuous at the interfaces either of the primal or of the dual discretization. Therefore, no particular regularity is assumed for the test spaces in the continuous problem as well, and they are defined by  $\mathcal{V} = [L^2(\Omega)]^2$  and  $\mathcal{Q} = L^2(\Omega)$ , respectively.

With the definition of the bilinear forms

$$\begin{aligned} a : \mathcal{U} \times \mathcal{V} &\rightarrow \mathbb{R} & \text{with} & & a(\mathbf{w}, \mathbf{v}) &:= (\mathbf{w}, \mathbf{v})_{0,\Omega} \\ b_1 : \mathcal{H} \times \mathcal{V} &\rightarrow \mathbb{R} & \text{with} & & b_1(r, \mathbf{v}) &:= (\nabla r, \mathbf{v})_{0,\Omega} \\ b_2 : \mathcal{U} \times \mathcal{Q} &\rightarrow \mathbb{R} & \text{with} & & b_2(\mathbf{w}, q) &:= (\nabla \cdot \mathbf{w}, q)_{0,\Omega} \end{aligned} \quad (43)$$

problem (36) can be reformulated to obtain the following continuous saddle point problem. Find  $((h\mathbf{v})^{n+1}, \delta t h_0 h^{(2)}) \in \mathcal{U} \times \mathcal{H}$ , such that

$$\begin{aligned} a((h\mathbf{v})^{n+1}, \boldsymbol{\varphi}) + b_1(\delta t h_0 h^{(2)}, \boldsymbol{\varphi}) &= ((h\mathbf{v})^{**}, \boldsymbol{\varphi})_{0,\Omega} \quad \forall \boldsymbol{\varphi} \in \mathcal{V} \\ b_2((h\mathbf{v})^{n+1}, \psi) &= -b_2((h\mathbf{v})^n, \psi) \quad \forall \psi \in \mathcal{Q} \end{aligned} \quad (44)$$

This obviously defines a problem of the form (28). The formulation is also referred to as a *primal-dual* formulation [THOMAS and TRUJILLO, 1999; ANGERMANN, 2003]. In VATER [2005] it is shown that the given bilinear forms are bounded and that the inf-sup conditions (29)–(32) are satisfied. Thus, the following theorem can be stated:

**Theorem 3.2 (VATER [2005])** *The generalized saddle point problem defined by (44) has a unique solution  $((h\mathbf{v})^{n+1}, \delta t h_0 h^{(2)})$  in  $\mathcal{U} \times \mathcal{H}$ .*

Since  $\mathcal{H}^h \subset \mathcal{H}$ ,  $\mathcal{U}^h \subset [L^2(\Omega)]^2$  and  $\mathcal{V}^h \subset \mathcal{V}$ , and the discrete gradient  $\mathbf{G}$  is equal to its continuous counterpart on each grid cell, the bilinear forms  $a$  and  $b_1$  are well defined on  $\mathcal{U}^h \times \mathcal{V}^h$  and  $\mathcal{H}^h \times \mathcal{V}^h$ , respectively. This is different for  $b_{2h}$ , since  $\mathcal{U}^h \not\subset \mathcal{U}$ . The bilinear form represents the discrete divergence from (24), which motivates the definition

$$b_{2h} : \mathcal{U}^h \times \mathcal{Q}^h \rightarrow \mathbb{R} \quad \text{with} \quad b_{2h}(\mathbf{v}_h, q_h) := \sum_{\bar{C} \in \bar{\mathcal{C}}} q_{h,\bar{C}} \int_{\partial \bar{C}} \mathbf{v}_h \cdot \mathbf{n} d\sigma, \quad (45)$$

where  $q_{h,\bar{C}}$  is the (constant) value of  $q_h$  on  $\bar{C}$ . This definition is consistent with the definition of its continuous counterpart  $b_2$ , since for functions  $\mathbf{v} \in H(\text{div}; \Omega)$  they both give the same result. Furthermore, the  $H(\text{div}; \Omega)$  norm is no longer appropriate for the space  $\mathcal{U}^h$ , and a suitable mesh dependent norm  $\|\cdot\|_{\mathcal{U}^h}$  has to be introduced (cf. [BRAESS, 2003]).

**Proposition 3.3** *A norm on the finite element space  $\mathcal{U}^h$  is given by*

$$\|\mathbf{w}_h\|_{\mathcal{U}^h} := \|\mathbf{w}_h\|_{0,\Omega} + \sup_{z_h \in \mathcal{Q}^h} \frac{b_{2h}(\mathbf{w}_h, z_h)}{\|z_h\|_{\mathcal{Q}}} \quad \text{for } \mathbf{w}_h \in \mathcal{U}^h .$$

**Proof.** *We have to show definiteness, homogeneity, and the triangle inequality of  $\|\cdot\|_{\mathcal{U}^h}$ :*

- *First, it follows by the definition of the norm that for  $\mathbf{w}_h \in \mathcal{U}^h$  with  $\|\mathbf{w}_h\|_{\mathcal{U}^h} = 0$  one obtains  $\|\mathbf{w}_h\|_{0,\Omega} = 0$ . Since  $\mathbf{w}_h$  is piecewise linear, i.e., piecewise continuous,  $\mathbf{w}_h$  has to be zero almost everywhere.*
- *For  $\lambda \in \mathbb{R}$  and  $\mathbf{w}_h \in \mathcal{U}^h$  we have*

$$\|\lambda \mathbf{w}_h\|_{\mathcal{U}^h} = \|\lambda \mathbf{w}_h\|_{0,\Omega} + \sup_{z_h \in \mathcal{Q}^h} \frac{b_{2h}(\lambda \mathbf{w}_h, z_h)}{\|z_h\|_{\mathcal{Q}}} = |\lambda| \|\mathbf{w}_h\|_{\mathcal{U}^h} .$$

- *The triangle inequality holds for  $\mathbf{w}_h, \tilde{\mathbf{w}}_h \in \mathcal{U}^h$ , since*

$$\begin{aligned} \|\mathbf{w}_h + \tilde{\mathbf{w}}_h\|_{\mathcal{U}^h} &= \|\mathbf{w}_h + \tilde{\mathbf{w}}_h\|_{0,\Omega} + \sup_{z_h \in \mathcal{Q}^h} \frac{b_{2h}(\mathbf{w}_h + \tilde{\mathbf{w}}_h, z_h)}{\|z_h\|_{\mathcal{Q}}} \\ &\leq \|\mathbf{w}_h\|_{0,\Omega} + \|\tilde{\mathbf{w}}_h\|_{0,\Omega} + \sup_{z_h \in \mathcal{Q}^h} \frac{b_{2h}(\mathbf{w}_h, z_h) + b_{2h}(\tilde{\mathbf{w}}_h, z_h)}{\|z_h\|_{\mathcal{Q}}} \\ &\leq \|\mathbf{w}_h\|_{\mathcal{U}^h} + \|\tilde{\mathbf{w}}_h\|_{\mathcal{U}^h} \quad \square \end{aligned}$$

In this norm, the bilinear form  $b_{2h}$  is continuous, since for arbitrary  $q_h \in \mathcal{Q}^h$  and  $\mathbf{w}_h \in \mathcal{U}^h$  it follows that

$$\begin{aligned} b_{2h}(\mathbf{w}_h, q_h) &= \frac{\|q_h\|_{\mathcal{Q}} b_{2h}(\mathbf{w}_h, q_h)}{\|q_h\|_{\mathcal{Q}}} \\ &\leq \|q_h\|_{\mathcal{Q}} \sup_{z_h \in \mathcal{Q}^h} \frac{b_{2h}(\mathbf{w}_h, z_h)}{\|z_h\|_{\mathcal{Q}}} \\ &\leq \|q_h\|_{\mathcal{Q}} \|\mathbf{w}_h\|_{\mathcal{U}^h} \end{aligned}$$

**Proposition 3.4** *For  $\mathbf{w}_h \in \mathcal{U}^h$  one has*

$$\sup_{z_h \in \mathcal{Q}^h} \frac{b_{2h}(\mathbf{w}_h, z_h)}{\|z_h\|_{\mathcal{Q}}} = \left( \sum_{\bar{C} \in \bar{\mathcal{C}}} \frac{1}{|\bar{C}|} \left( \int_{\partial \bar{C}} \mathbf{w}_h \cdot \mathbf{n} \, d\sigma \right)^2 \right)^{1/2} .$$

**Proof.** Taking  $\mathbf{w}_h \in \mathcal{U}^h$  and  $z_h \in \mathcal{Q}^h$  it follows from the Cauchy-Schwarz inequality that

$$\begin{aligned} \frac{b_{2h}(\mathbf{w}_h, z_h)}{\|z_h\|_{\mathcal{Q}}} &= \frac{\sum_{\bar{C}} z_h|_{\bar{C}} \int_{\partial\bar{C}} \mathbf{w}_h \cdot \mathbf{n} \, d\sigma}{\left(\sum_{\bar{C}} |\bar{C}| (z_h|_{\bar{C}})^2\right)^{1/2}} \\ &= \frac{\sum_{\bar{C}} (|\bar{C}|^{1/2} z_h|_{\bar{C}}) \left(|\bar{C}|^{-1/2} \int_{\partial\bar{C}} \mathbf{w}_h \cdot \mathbf{n} \, d\sigma\right)}{\left(\sum_{\bar{C}} |\bar{C}| (z_h|_{\bar{C}})^2\right)^{1/2}} \\ &\leq \frac{\left(\sum_{\bar{C}} |\bar{C}| (z_h|_{\bar{C}})^2\right)^{1/2} \left(\sum_{\bar{C}} |\bar{C}|^{-1} \left(\int_{\partial\bar{C}} \mathbf{w}_h \cdot \mathbf{n} \, d\sigma\right)^2\right)^{1/2}}{\left(\sum_{\bar{C}} |\bar{C}| (z_h|_{\bar{C}})^2\right)^{1/2}} \\ &= \left(\sum_{\bar{C}} \frac{1}{|\bar{C}|} \left(\int_{\partial\bar{C}} \mathbf{w}_h \cdot \mathbf{n} \, d\sigma\right)^2\right)^{1/2} \end{aligned}$$

Since  $z_h$  is arbitrary, this gives the proof in one direction. On the other hand, setting  $z_h|_{\bar{C}} := |\bar{C}|^{-1} \int_{\partial\bar{C}} \mathbf{w}_h \cdot \mathbf{n} \, d\sigma$  gives

$$\frac{b_{2h}(\mathbf{w}_h, z_h)}{\|z_h\|_{\mathcal{Q}}} = \left(\sum_{\bar{C}} \frac{1}{|\bar{C}|} \left(\int_{\partial\bar{C}} \mathbf{w}_h \cdot \mathbf{n} \, d\sigma\right)^2\right)^{1/2}$$

Taking the supremum over all  $z_h \in \mathcal{Q}^h$  leads to the desired result.  $\square$

With the definition of the bilinear form in (45), the discrete mixed formulation derived in Section 3.2 is to find  $((h\mathbf{v})^{n+1}, \delta t h_0 h^{(2)}) \in \mathcal{U}^h \times \mathcal{H}^h$ , such that

$$\begin{aligned} a((h\mathbf{v})^{n+1}, \varphi_h) + b_1(\delta t h_0 h^{(2)}, \varphi_h) &= ((h\mathbf{v})^{**}, \varphi_h)_{0,\Omega} \quad \forall \varphi_h \in \mathcal{V}^h \\ b_{2h}((h\mathbf{v})^{n+1}, \psi_h) &= -b_{2h}((h\mathbf{v})^n, \psi_h) \quad \forall \psi_h \in \mathcal{Q}^h. \end{aligned} \quad (46)$$

Note that the trial space  $\mathcal{U}^h$  is not contained in its continuous counterpart  $\mathcal{U}$ . Therefore, the discrete problem (46) is an approximation using *nonconforming finite elements*.

Now, the verification of the inf-sup-conditions can be carried out. The proof for the  $b_1$  form is nearly identical to the continuous case (cf. [VATER, 2005]).

**Proposition 3.5** *There exists a constant  $\beta_1^* > 0$  independent of the mesh size,  $h$ , such that*

$$\inf_{r_h \in \mathcal{H}^h} \sup_{\mathbf{v}_h \in \mathcal{V}^h} \frac{b_1(r_h, \mathbf{v}_h)}{\|r_h\|_{\mathcal{H}} \|\mathbf{v}_h\|_{\mathcal{V}}} \geq \beta_1^*$$

**Proof.** It has been already pointed out that  $r_h \in \mathcal{H}^h$  implies  $\nabla r_h \in \mathcal{V}^h$ . Thus, we have for arbitrary  $r_h \in \mathcal{H}^h$

$$\sup_{\mathbf{v}_h \in \mathcal{V}^h} \frac{b_1(r_h, \mathbf{v}_h)}{\|\mathbf{v}_h\|_{\mathcal{V}}} \geq \frac{b_1(r_h, \nabla r_h)}{\|\nabla r_h\|_{0,\Omega}} = \frac{|r_h|_{1,\Omega}^2}{|r_h|_{1,\Omega}} = \|r_h\|_{\mathcal{H}} \quad . \quad \square$$

Next, it is proved what is normally known as coercivity for the bilinear form  $a$ . Since we deal with a Petrov-Galerkin method, the characterization has to be generalized to the two conditions (29) and (30). Let us define the subspaces

$$\begin{aligned}\mathcal{K}_1^h &:= \{\mathbf{v}_h \in \mathcal{V}^h \mid \forall r_h \in \mathcal{H}^h : b_1(r_h, \mathbf{v}_h) = 0\} \\ \mathcal{K}_2^h &:= \{\mathbf{w}_h \in \mathcal{U}^h \mid \forall q_h \in \mathcal{Q}^h : b_{2h}(\mathbf{w}_h, q_h) = 0\}.\end{aligned}$$

In the following, it is shown that there is a one-to-one mapping between these spaces, and an estimate can be given between corresponding elements. To characterize the spaces, it suffices to test the bilinear forms that are used in defining them against a complete set of basis functions of the test spaces. Thus, let  $r_h \in \mathcal{H}^h$  with  $r_h(x_{k+1/2}, y_{l+1/2}) = \delta_{ik}\delta_{jl}$  for a given node  $(x_{i+1/2}, y_{j+1/2})$ . Assuming a cell wise representation of  $r_h$  (cf. (39)), a careful investigation of such a basis function reveals that  $r_{x,l,k} = \pm \frac{1}{2\delta x}$ ,  $r_{y,l,k} = \pm \frac{1}{2\delta y}$  and  $r_{xy,l,k} = \pm \frac{1}{\delta x \delta y}$  for  $l \in \{i, i+1\}$ ,  $k \in \{j, j+1\}$ . Thus,  $\mathbf{v}_h = (u, v) \in \mathcal{V}^h$  is in  $\mathcal{K}_1^h$ , if and only if for all possible  $(i, j)$

$$\begin{aligned}0 &= b_1(r_h, \mathbf{v}_h) = \sum_{l,k} \int_{C_{lk}} \nabla r_h \cdot \mathbf{v}_h \, d\mathbf{x} \\ &= \sum_{l=i}^{i+1} \sum_{k=j}^{j+1} \delta x \delta y \left( u_{l,k} r_{x,l,k} + v_{l,k} r_{y,l,k} + \frac{1}{12} (\delta y^2 u_{y,l,k} + \delta x^2 v_{x,l,k}) r_{xy,l,k} \right) \\ &= -\frac{\delta y}{2} u_{i+1,j+1} - \frac{\delta x}{2} v_{i+1,j+1} + \frac{\delta y^2}{12} u_{y,i+1,j+1} + \frac{\delta x^2}{12} v_{x,i+1,j+1} \\ &\quad + \frac{\delta y}{2} u_{i,j+1} - \frac{\delta x}{2} v_{i,j+1} - \frac{\delta y^2}{12} u_{y,i,j+1} - \frac{\delta x^2}{12} v_{x,i,j+1} \\ &\quad + \frac{\delta y}{2} u_{i,j} + \frac{\delta x}{2} v_{i,j} + \frac{\delta y^2}{12} u_{y,i,j} + \frac{\delta x^2}{12} v_{x,i,j} \\ &\quad - \frac{\delta y}{2} u_{i+1,j} + \frac{\delta x}{2} v_{i+1,j} - \frac{\delta y^2}{12} u_{y,i+1,j} - \frac{\delta x^2}{12} v_{x,i+1,j}\end{aligned}\tag{47}$$

Similarly, let  $q_h \in \mathcal{Q}^h$  with  $q_h = \chi_{\bar{C}_{i+1/2,j+1/2}}$  be arbitrary. Then,  $\mathbf{w}_h = (u, v) \in \mathcal{U}^h$  is in  $\mathcal{K}_2^h$ , if and only if for all possible  $(i, j)$

$$\begin{aligned}0 &= -b_{2h}(\mathbf{w}_h, q_h) = - \int_{\partial \bar{C}_{i+1/2,j+1/2}} \mathbf{w}_h \cdot \mathbf{n} \, d\sigma \\ &= -\frac{\delta y}{2} u_{i+1,j+1} - \frac{\delta x}{2} v_{i+1,j+1} + \frac{\delta y^2}{8} u_{y,i+1,j+1} + \frac{\delta x^2}{8} v_{x,i+1,j+1} \\ &\quad + \frac{\delta y}{2} u_{i,j+1} - \frac{\delta x}{2} v_{i,j+1} - \frac{\delta y^2}{8} u_{y,i,j+1} - \frac{\delta x^2}{8} v_{x,i,j+1} \\ &\quad + \frac{\delta y}{2} u_{i,j} + \frac{\delta x}{2} v_{i,j} + \frac{\delta y^2}{8} u_{y,i,j} + \frac{\delta x^2}{8} v_{x,i,j} \\ &\quad - \frac{\delta y}{2} u_{i+1,j} + \frac{\delta x}{2} v_{i+1,j} - \frac{\delta y^2}{8} u_{y,i+1,j} - \frac{\delta x^2}{8} v_{x,i+1,j}\end{aligned}\tag{48}$$



Comparing (47) and (48), we observe that these conditions only differ by a constant factor in the terms, which include partial derivatives of the velocity components. This means that a one-to-one mapping between  $\mathcal{K}_1^h$  and  $\mathcal{K}_2^h$  can be defined by multiplying the partial derivatives of an element with  $8/12 = 2/3$ , and the spaces have the same dimension. Furthermore, the following estimates can be given for corresponding elements  $\mathbf{v}_h \in \mathcal{K}_1^h$  and  $\mathbf{w}_h \in \mathcal{K}_2^h$  (i.e. with the same mean values  $\bar{\mathbf{w}}_h = \bar{\mathbf{v}}_h$ , and linear variations  $\nabla \tilde{\mathbf{w}}_h = 2/3 \nabla \tilde{\mathbf{v}}_h$ ):

$$\begin{aligned} a(\mathbf{w}_h, \mathbf{v}_h) &= a(\bar{\mathbf{w}}_h, \bar{\mathbf{v}}_h) + a(\tilde{\mathbf{w}}_h, \tilde{\mathbf{v}}_h) \\ &= a(\bar{\mathbf{v}}_h, \bar{\mathbf{v}}_h) + \frac{2}{3} a(\tilde{\mathbf{v}}_h, \tilde{\mathbf{v}}_h) \geq \frac{2}{3} a(\mathbf{v}_h, \mathbf{v}_h) \end{aligned}$$

and

$$a(\mathbf{w}_h, \mathbf{v}_h) = a(\bar{\mathbf{w}}_h, \bar{\mathbf{w}}_h) + \frac{3}{2} a(\tilde{\mathbf{w}}_h, \tilde{\mathbf{w}}_h) \geq a(\mathbf{w}_h, \mathbf{w}_h)$$

and

$$\begin{aligned} a(\mathbf{v}_h, \mathbf{v}_h) &\leq \frac{3}{2} a(\mathbf{w}_h, \mathbf{v}_h) = \frac{3}{2} \left( a(\bar{\mathbf{w}}_h, \bar{\mathbf{w}}_h) + \frac{3}{2} a(\tilde{\mathbf{w}}_h, \tilde{\mathbf{w}}_h) \right) \\ &\leq \frac{9}{4} a(\mathbf{w}_h, \mathbf{w}_h) \end{aligned}$$

With these estimates, we can also prove the desired properties for the  $a$  form in the discrete case:

**Proposition 3.6** *There exists a constant  $\alpha^* > 0$  independent of the mesh size,  $h$ , such that*

$$\inf_{\mathbf{w}_h \in \mathcal{K}_2^h} \sup_{\mathbf{v}_h \in \mathcal{K}_1^h} \frac{a(\mathbf{w}_h, \mathbf{v}_h)}{\|\mathbf{w}_h\|_{\mathcal{U}^h} \|\mathbf{v}_h\|_{\mathcal{V}}} \geq \alpha^* \quad . \quad (49)$$

Furthermore,

$$\sup_{\mathbf{w}_h \in \mathcal{K}_2^h} a(\mathbf{w}_h, \mathbf{v}_h) > 0 \quad \forall \mathbf{v}_h \in \mathcal{K}_1^h \setminus \{0\} \quad . \quad (50)$$

**Proof.** For  $\mathbf{w}_h \in \mathcal{K}_2^h$ ,  $\|\mathbf{w}_h\|_{\mathcal{U}^h} = \|\mathbf{w}_h\|_{0,\Omega} \neq 0$ . Thus, using the estimates derived from the one-to-one mapping above, for each such  $\mathbf{w}_h$  we have

$$\sup_{\mathbf{v}_h \in \mathcal{K}_1^h} \frac{a(\mathbf{w}_h, \mathbf{v}_h)}{\|\mathbf{v}_h\|_{\mathcal{V}^h}} \geq \frac{a(\mathbf{w}_h, \mathbf{w}_h)}{\frac{3}{2} \|\mathbf{w}_h\|_{0,\Omega}} = \frac{2}{3} \frac{\|\mathbf{w}_h\|_{0,\Omega}^2}{\|\mathbf{w}_h\|_{0,\Omega}} = \frac{2}{3} \|\mathbf{w}_h\|_{\mathcal{U}^h} \quad ,$$

and for  $\mathbf{v}_h \in \mathcal{K}_1^h \setminus \{0\}$

$$\sup_{\mathbf{w}_h \in \mathcal{K}_2^h} a(\mathbf{w}_h, \mathbf{v}_h) \geq \frac{2}{3} a(\mathbf{v}_h, \mathbf{v}_h) > 0 \quad .$$

Therefore, the conditions (49) and (50) are satisfied.  $\square$

Before the inf-sup condition for the bilinear form  $b_{2h}$  is also proved, a *lumping operator*  $\mathbf{L} : \mathcal{H}^h \rightarrow \mathcal{Q}^h$  is introduced, which is given by

$$\mathbf{L}r_h := \sum_{\bar{C} \in \bar{\mathcal{C}}} \chi_{\bar{C}} r_h(x_{\bar{C}}, y_{\bar{C}}) \quad \forall r_h \in \mathcal{H}^h,$$

where  $(x_{\bar{C}}, y_{\bar{C}})$  again is the midpoint of  $\bar{C}$ , i.e. the coordinate of the grid node around which  $\bar{C}$  is centered. Thus, in each dual control volume, the value of  $\mathbf{L}r_h$  is the value of  $r_h$  at the corresponding node in the middle of the control volume. This operator has the following properties:

**Proposition 3.7** *For  $r_h \in \mathcal{H}^h$  with  $\nabla r_h \cdot \mathbf{n} \equiv 0$  on  $\partial\Omega$  we have*

$$\|\nabla r_h\|_{0,\Omega}^2 \leq -b_{2h}(\nabla r_h, \mathbf{L}r_h)$$

**Proof.** See Appendix A.2. □

**Proposition 3.8** *For  $r_h \in \mathcal{H}^h$  the estimate*

$$\|\mathbf{L}r_h\|_{0,\Omega} \leq C \|r_h\|_{0,\Omega}$$

where  $C$  is a constant, holds.

**Proof.** See Appendix A.2. □

Now, we are in the position to prove the inf-sup condition for  $b_{2h}$ . The general idea is adapted from a proof of a similar problem in ANGERMANN [2003].

**Proposition 3.9** *There exists a constant  $\beta_2^* > 0$  independent of the mesh size,  $h$ , such that*

$$\inf_{q_h \in \mathcal{Q}^h} \sup_{\mathbf{w}_h \in \mathcal{U}^h} \frac{b_{2h}(\mathbf{w}_h, q_h)}{\|\mathbf{w}_h\|_{\mathcal{U}^h} \|q_h\|_{\mathcal{Q}}} \geq \beta_2^*$$

**Proof.** To show the inf-sup condition for the  $b_{2h}(\cdot, \cdot)$  form an auxiliary mapping  $G_h : \mathcal{Q}^h \rightarrow \mathcal{U}^h$  is introduced. It is defined by the solution of the Poisson problem

$$r_h \in \mathcal{H}^h : \quad -\mathbf{L}(r_h) = q_h$$

for  $q_h \in \mathcal{Q}^h$ , where  $G_h q_h := \nabla r_h \in \mathcal{U}^h$ . This is to find  $r_h \in \mathcal{H}^h$ , such that

$$b_{2h}(\nabla r_h, z_h) = (q_h, z_h)_{0,\Omega} \quad \forall z_h \in \mathcal{Q}^h. \quad (51)$$

Using the properties of the lumping operator  $\mathbf{L}$  and the Poincaré inequality, the following estimate can be given for the solution  $r_h$  of the Poisson problem (51) (which can

be shown to have a unique solution for fixed mesh size  $h$ ):

$$\begin{aligned}
\|\nabla r_h\|_{0,\Omega}^2 &\leq -b_{2h}(\nabla r_h, \mathbf{L}r_h) && \text{(Proposition 3.7)} \\
&= (q_h, \mathbf{L}r_h)_{0,\Omega} && \text{(Poisson problem (51))} \\
&\leq \|q_h\|_{0,\Omega} \|\mathbf{L}r_h\|_{0,\Omega} && \text{(Cauchy-Schwarz inequality)} \\
&\leq C_1 \|q_h\|_{0,\Omega} \|r_h\|_{0,\Omega} && \text{(Proposition 3.8)} \\
&\leq C_2 \|q_h\|_{0,\Omega} \|\nabla r_h\|_{0,\Omega} . && \text{(Poincaré inequality)}
\end{aligned}$$

Thus, we have

$$\|\nabla r_h\|_{0,\Omega} \leq C_2 \|q_h\|_{0,\Omega} .$$

Furthermore, this solution satisfies

$$\sup_{z_h \in \mathcal{Q}^h} \frac{b_{2h}(\nabla r_h, z_h)}{\|z_h\|_{\mathcal{Q}}} = \sup_{z_h \in \mathcal{Q}^h} \frac{(q_h, z_h)_{0,\Omega}}{\|z_h\|_{\mathcal{Q}}} = \|q_h\|_{\mathcal{Q}} .$$

By the definition of the norm on  $\mathcal{U}^h$ , it then follows that

$$\|G_h q_h\|_{\mathcal{U}^h} = \|\nabla r_h\|_{0,\Omega} + \sup_{z_h \in \mathcal{Q}^h} \frac{b_{2h}(\nabla r_h, z_h)}{\|z_h\|_{\mathcal{Q}}} \leq C \|q_h\|_{\mathcal{Q}}$$

where  $C = 1 + C_2$ , and

$$\|G_h q_h\|_{\mathcal{U}^h} \|q_h\|_{\mathcal{Q}} \leq C \|q_h\|_{\mathcal{Q}}^2 = C b_{2h}(\nabla r_h, q_h) = C b_{2h}(G_h q_h, q_h)$$

which leads to

$$\frac{1}{C} \leq \frac{b_{2h}(G_h q_h, q_h)}{\|G_h q_h\|_{\mathcal{U}^h} \|q_h\|_{\mathcal{Q}}} \leq \sup_{\mathbf{w}_h \in \mathcal{U}^h} \frac{b_{2h}(\mathbf{w}_h, q_h)}{\|\mathbf{w}_h\|_{\mathcal{U}^h} \|q_h\|_{\mathcal{Q}}}$$

Since  $q_h$  was chosen arbitrarily, this proves the inf-sup condition for  $b_{2h}(\cdot, \cdot)$ .  $\square$

As a summary of this section, we can conclude with the following:

**Theorem 3.10** *The generalized mixed formulation (46) has a unique and stable solution  $((h\mathbf{v})^{n+1}, \delta t h_0 h^{(2)})$  in  $\mathcal{U}^h \times \mathcal{H}^h$ .*

We have successfully established a mixed formulation equivalent to the second projection of the new scheme. Using this formulation for the stability analysis of the projection step, stability has been shown for the discrete problem. This gives approximations, in which the solution of the Poisson problem  $h^{(2)}$  and the momentum update  $(h\mathbf{v})^{n+1}$  cannot decouple.

## 4 Numerical Results

To illustrate the performance of the described projection method, the results of two test cases are presented. The main goal is to assess its accuracy and to compare it

with a previous version of the method introduced by [SCHNEIDER ET AL., 1999] which rests on standard discretizations for the differential operators used in the projection step. Furthermore, the differences between an exact and an approximate projection formulation are assessed. In the first test case, the second-order convergence of the method is demonstrated for smooth solutions. The second test deals with the translation of a vortex.

For both test cases the exact solution of the particular problem is known, and the error of the numerical approximation can be computed. The computations are performed on a uniform Cartesian grid with equal grid spacing  $\delta x = \delta y$ . The boundary conditions are those discussed in VATER [2005]. So far, we have only investigated the case of constant background height  $h_0 \equiv 1$ . Thus, in all calculations, the term  $dh_0/dt$  is set to zero. To start with initial data, which have zero divergence, i.e.,

$$\mathbf{D}(\mathbf{v}^0) = -\frac{1}{h_0} \frac{dh_0}{dt} \Big|_{t=0} = 0 ,$$

the given values for the momentum are corrected by the solution of the Poisson problem

$$\mathbf{L}(\varphi) = \mathbf{D}((h\mathbf{v})^{0,r})$$

for  $\varphi \in \mathcal{H}^h$ . Here,  $(h\mathbf{v})^{0,r}$  is a linear reconstruction of the exact solution  $(h\mathbf{v})$  at time  $t = 0$ . The initial momentum distribution is then given by

$$(h\mathbf{v})^0 = (h\mathbf{v})^{0,r} - \mathbf{G}(\varphi) \quad .$$

As mentioned earlier, the auxiliary system is solved using an explicit standard second-order Godunov-type method for hyperbolic conservation laws. Since the stability of this method strongly relies on a CFL time step restriction, in all the computations presented in this chapter a time step has been chosen, which is at least  $C = 0.8$  times smaller than the maximum allowed by the CFL condition.

The discrete divergence and gradient operators, which are used in the two elliptic correction steps, are those given in Appendix A.1. The linear systems for computing the height  $h^{(2)}$  on the primary and on the dual discretizations are solved using the Bi-CGSTAB algorithm [VAN DER VORST, 1992]. In each iteration, the Euclidean norm

$$\|r_{\mathcal{C}}\|_2 := \sqrt{\sum_{C \in \mathcal{C}} r_C^2}$$

(similarly for the second Poisson problem with  $\|r_{\bar{\mathcal{C}}}\|_2$ ) of the residual vector

$$\begin{aligned} r_{\mathcal{P}_1}(h^{(2)}) &:= \mathbf{D}((h\mathbf{v})^*) - \frac{\delta t}{2} \mathbf{D}(h_0^{n+1/4} \mathbf{G}(h^{(2)})) \\ r_{\mathcal{P}_2}(h^{(2)}) &:= \mathbf{D}((h\mathbf{v})^{**}) + \mathbf{D}((h\mathbf{v})^n) - \delta t \mathbf{D}(h_0^{n+1/2} \mathbf{G}(h^{(2)})) \end{aligned}$$

is calculated. The algorithm is terminated when either this absolute value or the ratio between the norm of the current residual and that of the initial residual is less than  $10^{-11}$ .

#### 4.1 Convergence study

The first test case demonstrates the second-order convergence of numerical solutions to the exact solution for smooth data. This test, which involves a Taylor vortex being translated at a constant speed, was originally proposed in MINION [1996] and ALMGREN ET AL. [1998] for the incompressible flow equations. Here it has been adapted for the zero Froude number shallow water equations.

For constant height  $h_0$  and an initial velocity distribution

$$\begin{aligned} u_0(x, y) &= 1 - 2 \cos(2\pi x) \sin(2\pi y) \\ v_0(x, y) &= 1 + 2 \sin(2\pi x) \cos(2\pi y) \quad , \end{aligned}$$

the exact solution of the zero Froude number shallow water equations is given by

$$\begin{aligned} u(x, y, t) &= 1 - 2 \cos(2\pi(x - t)) \sin(2\pi(y - t)) \\ v(x, y, t) &= 1 + 2 \sin(2\pi(x - t)) \cos(2\pi(y - t)) \\ h^{(2)}(x, y, t) &= -\cos(4\pi(x - t)) - \cos(4\pi(y - t)) \quad . \end{aligned}$$

The problem is solved on the unit square with  $(x, y) \in [0, 1]^2$  and periodic boundary conditions. It describes the advection of four vortices in the  $(1, 1)$  direction. The piecewise linear reconstruction of the momentum field components is done using central differences with no slope limiter.

The numerical solution is computed on three different grids with  $32 \times 32$ ,  $64 \times 64$  and  $128 \times 128$  cells. We start the calculation at  $t = 0$ , and the error vector in the velocity  $\mathbf{e}^N$  with elements

$$e_{i,j}^N := \left| \overline{u(x, y, t^N)^{C_{i,j}}} - u_{i,j}^N \right| + \left| \overline{v(x, y, t^N)^{C_{i,j}}} - v_{i,j}^N \right|$$

is evaluated at time  $t^N = 3$ . This corresponds to 750, 1500 and 3000 time steps, respectively. Note that we could have also incorporated the linear variation of the velocity on each grid cell in the error analysis of the new projection. We do not choose this alternative in favor of a better comparison with the original method. The global error is measured using a discrete  $L^2$  norm and the  $L^\infty$  norm. These are defined by

$$\|\mathbf{e}^N\|_0 := \left( \sum_{i,j} (|C_{i,j}| e_{i,j}^N)^2 \right)^{1/2} \quad \text{and} \quad \|\mathbf{e}^N\|_\infty := \max_{i,j} \{e_{i,j}^N\} \quad .$$

We have summarized these error measures for the original projection method by SCHNEIDER ET AL. [1999] as well as for the approximate and the exact projection methods in Table 1. Here, the ‘‘approximate projection method’’ utilizes the same stencil as the exact projection method, but it leaves slope computations for the in-cell distributions of momentum entirely to classical slope limiting procedures instead of letting several components of these derivatives be determined by the projection step.

Method	Norm	32x32	Rate $\gamma$	64x64	Rate $\gamma$	128x128
projection by SCHNEIDER ET AL. [1999]	$L^2$	0.292096	2.16	0.065415	2.17	0.014566
	$L^\infty$	0.419370	2.16	0.094106	2.18	0.020747
approximate projection	$L^2$	0.291967	2.16	0.065412	2.17	0.014566
	$L^\infty$	0.419130	2.16	0.094098	2.18	0.020747
exact projection	$L^2$	0.082379	2.65	0.013129	2.23	0.002796
	$L^\infty$	0.126207	2.46	0.022999	2.33	0.004573

**Table 1:** Errors and convergence rates for the different projection methods.

Additionally, the corresponding convergence rate  $\gamma$  is given, which is calculated by

$$\gamma := \frac{\log(\|e_c^N\| / \|e_f^N\|)}{\log(\delta x_c / \delta x_f)} . \quad (52)$$

In this definition,  $e_c^N$  and  $e_f^N$  are the computed error vectors of the solution on the coarse and the fine grid and  $\delta x_c$  and  $\delta x_f$  are the corresponding grid spacings. Clearly, second order accuracy is obtained in the  $L^2$  as well as in the  $L^\infty$  norm. Also note that the absolute error obtained with the exact projection is about four times smaller than the one obtained with the approximate projection method and with the scheme by SCHNEIDER ET AL. [1999].

## 4.2 Advection of a vortex

Let us consider the advection of a vortex by a constant background flow. For the implementation of this test case, originally proposed by GRESHO and CHAN [1990], a rectangular domain with size  $[0, 4] \times [0, 1]$  is examined. The domain has periodic boundary conditions at the short sides and walls at the long sides. The initial conditions are defined to be

$$u_0(x, y) = 1 - v_\theta(r) \sin \theta \quad \text{and} \quad v_0(x, y) = v_\theta(r) \cos \theta \quad ,$$

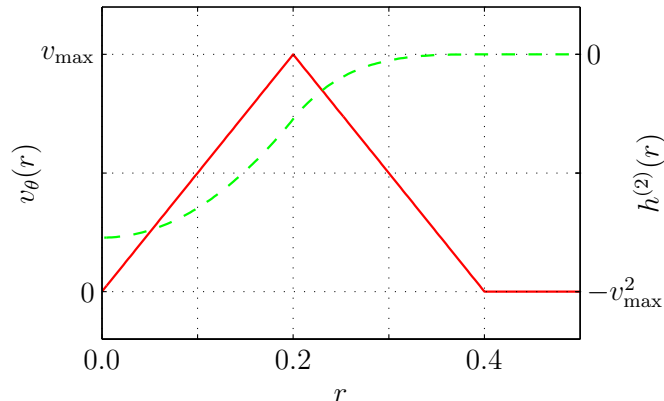
in which

$$v_\theta(r) = \begin{cases} 5r v_{\max} & \text{for } 0 \leq r < \frac{1}{5} \\ (2 - 5r) v_{\max} & \text{for } \frac{1}{5} \leq r < \frac{2}{5} \\ 0 & \text{for } \frac{2}{5} \leq r \end{cases} \quad (53)$$

and

$$r = \sqrt{\left(x - \frac{1}{2}\right)^2 + \left(y - \frac{1}{2}\right)^2} .$$

In equation (53)  $v_{\max}$  is the maximum tangential velocity of the vortex. The height  $h^{(2)}$  must then satisfy the constraint  $\partial_r h^{(2)} = v_\theta^2 / r$ . This relationship is visualized in Figure 5.



**Figure 5:** Advection of a vortex: tangential velocity (solid red) and height profile (dashed green) with respect to the distance  $r$  from the center of the vortex.

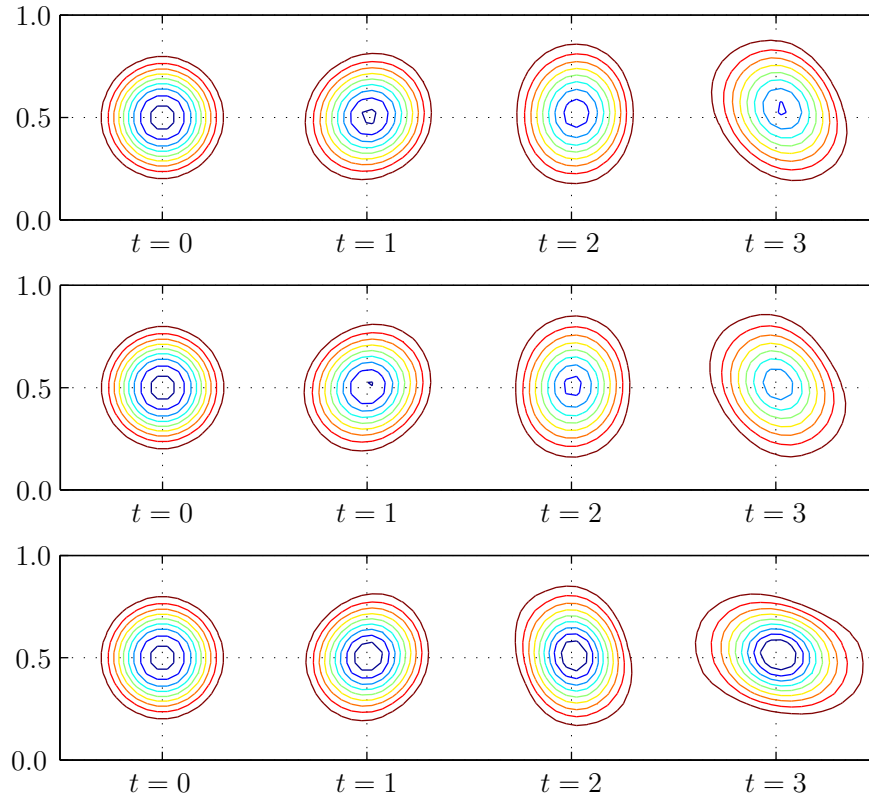
The test is set up with  $v_{\max} = 1$  and background height  $h_0 \equiv 1$ . The computational domain consists of  $80 \times 20$  grid cells. Three different strategies for the linear reconstruction of the components in the momentum variable are investigated. In particular, we consider central differences (no limiter), the *monotonized central difference (MC)* limiter and *Sweby's* limiter [SCHULZ-RINNE, 1993] with  $k = 1.8$ , the latter being a convex combination of the *minmod* ( $k = 1$ ) and the *superbee* limiter ( $k = 2$ ).

For comparison, the results for the scheme by SCHNEIDER ET AL. [1999] are given in Figure 6, in which the stream function of the velocity distribution is displayed at four different times of the simulation. Similar to the results in SCHNEIDER ET AL. [1999] for the incompressible Euler equations, the core is advected almost along the center line of the channel. Also, the vortex experiences a considerable deformation due to the coarse discretization we have chosen for this test.

As in the convergence studies, the new exact projection method shows a significant improvement in the numerical results for this test (cf. Figure 7). All reconstruction strategies show less deviation from the center line of the channel than in the original method. Furthermore, the loss in vorticity is slightly reduced. Again, the results of the approximate projection method (not shown) are comparable to the ones obtained by the method of SCHNEIDER ET AL. [1999].

## 5 Conclusions

In this paper, we demonstrate that it is possible to formulate a finite volume projection method for incompressible flows with an exact *and* stable projection step. No further stabilization techniques are required to prevent a velocity-pressure decoupling, which is often observed in former exact projection methods. This is achieved by using a Petrov-Galerkin finite element discretization of associated Poisson problem, originally proposed in SÜLI [1991]. Furthermore, the method locally conserves mass and momentum.



**Figure 6:** Advection of a vortex at times  $t = 0, 1, 2$  and  $3$  for the method by SCHNEIDER ET AL. [1999]. Contour lines of the stream function are shown at  $[-0.02, -0.04, \dots, -0.18]$  starting from outside of the vortex. Top: unlimited slopes, middle: monotonized central difference (MC) limiter, bottom: Sweby's limiter ( $k = 1.8$ ).



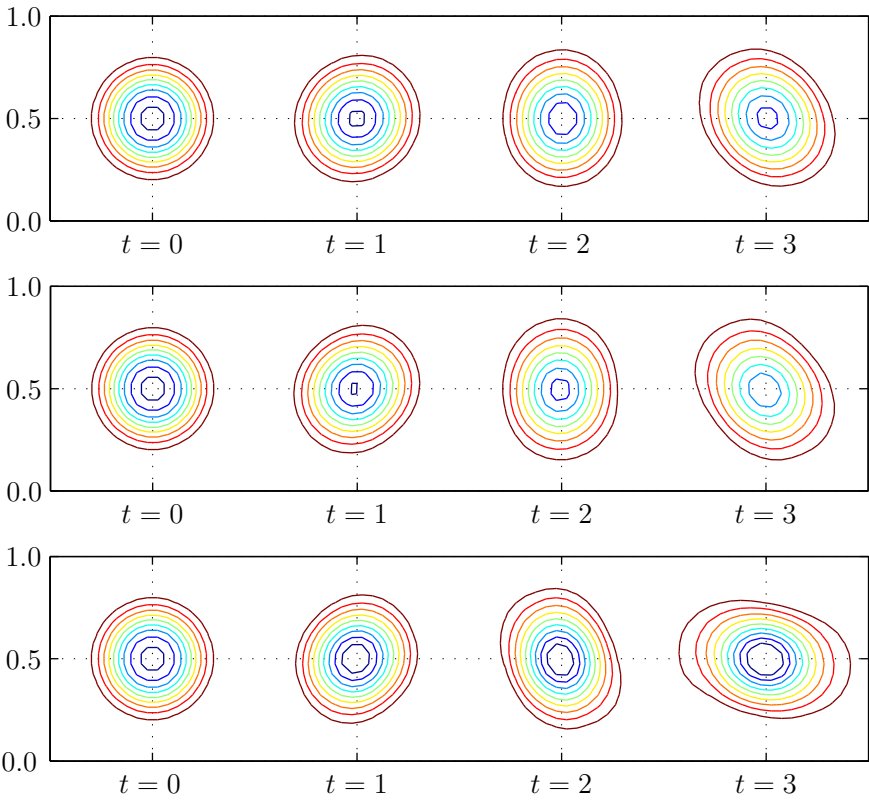


Figure 7: Same as Figure 6 for the new exact projection method.

In order to prove stability of the second projection step, which corrects the cell-centered momentum to be in compliance with the divergence constraint, we have used the theory of mixed finite element methods, the latter providing strong results about the stability of discretizations. This technique is well known from finite element methods for *viscous* incompressible flows, where the Laplacian of the velocity field interacts with the pressure gradient. Here, the theory is applied in the case of a finite volume method for *inviscid* incompressible flows, which means that the velocity directly interacts with the pressure gradient.

The numerical results, obtained from the application of the new method, show practical accuracy improvements on fixed grids compared to the method presented in SCHNEIDER ET AL. [1999], with both methods being second order accurate. The discretization for the new projection can be also used for the first projection of the method, yielding a unified discretization for both Poisson-type problems. Furthermore, the linear systems associated with the Poisson equations can be solved with the same algorithms that are used for standard second order discretizations of the differential operators.

However, there are still some open questions, and the analysis of them is ongoing research. It was mentioned that the second projection adjusts the piecewise linear portions of the momentum field, which, in turn, results in a possible loss of the TVD property of the whole method. This is a delicate issue, because it concerns the stability of the predictor step. So far, we have only numerical evidence that this still results in stable approximations.

In the present paper, it was only proved that the projection step yields a stable approximation and does not admit any pressure-velocity decoupling. One of the next steps is to investigate the convergence of this mixed formulation to the continuous solution (mentioned in Theorem 3.2). Furthermore, it would be desirable to extend our results to zero Mach number variable density flows, where the Projection results in a Poisson-type problem with a weighted Laplace operator.

The overall motivation for this work stems from meteorological and combustion applications. In such problems, we have small, but non-zero Mach numbers (resp. Froude numbers). The extension of the current method to allow for smooth transitions from fully compressible to zero Froude number flows would hopefully yield favorable results for these application areas. Such attempts were already reported in KLEIN [1995] for one space dimension and in GERATZ [1997] for higher dimensional problems. We are planning to advance the ideas outlined in these references.

## Acknowledgements

The authors are grateful to Eberhard Bänsch and Nicola Botta for helpful discussions and remarks that improved the content of the paper. Furthermore, Nicola Botta provided the software environment (the compressible flow solver), in which the algorithm was implemented and the presented calculations were accomplished. This work benefited greatly from free software products. Without these tools – such as L<sup>A</sup>T<sub>E</sub>X, the GNU C/C++ compiler and the Linux operating system – a lot of tasks would not have been

so easy to realize. It is our pleasure to thank all developers for their excellent products.

The authors thank Deutsche Forschungsgemeinschaft for their partial support of this work through grants KL 611/6 and KL 611/14.

## A Appendix

### A.1 Discretization of the new projection

Here, the discrete gradient, divergence and Laplacian of the second projection are given for a two-dimensional Cartesian grid with constant grid spacings  $\delta x$  and  $\delta y$ . The operators for the first projection are derived by shifting the indices by one half. The double index  $(i, j)$  is used to refer to a cell value, while the index  $(i + 1/2, j + 1/2)$  is used for node values.

Let us define

$$\begin{aligned} p_{x,i,j} &:= \frac{1}{2\delta x} (p_{i+1/2,j+1/2} - p_{i-1/2,j+1/2} + p_{i+1/2,j-1/2} - p_{i-1/2,j-1/2}) \\ p_{y,i,j} &:= \frac{1}{2\delta y} (p_{i+1/2,j+1/2} - p_{i+1/2,j-1/2} + p_{i-1/2,j+1/2} - p_{i-1/2,j-1/2}) \\ p_{xy,i,j} &:= \frac{1}{\delta x \delta y} (p_{i+1/2,j+1/2} - p_{i-1/2,j+1/2} - p_{i+1/2,j-1/2} + p_{i-1/2,j-1/2}) . \end{aligned}$$

The discrete gradient  $\mathbf{G}$  is then given by

$$\mathbf{G}(p) = \begin{pmatrix} p_{x,i,j} \\ p_{y,i,j} \end{pmatrix} + \begin{pmatrix} y - y_j \\ x - x_i \end{pmatrix} p_{xy,i,j} .$$

The divergence  $\mathbf{D}$  is defined by

$$\begin{aligned} \mathbf{D}(\mathbf{v}) &= \frac{1}{2\delta x} (u_{i+1,j+1} - u_{i,j+1} + u_{i+1,j} - u_{i,j}) \\ &\quad + \frac{\delta y}{8\delta x} (-u_{y,i+1,j+1} + u_{y,i,j+1} + u_{y,i+1,j} - u_{y,i,j}) \\ &\quad + \frac{1}{2\delta y} (v_{i+1,j+1} - v_{i+1,j} + v_{i,j+1} - v_{i,j}) \\ &\quad + \frac{\delta x}{8\delta y} (-v_{x,i+1,j+1} + v_{x,i+1,j} + v_{x,i,j+1} - v_{x,i,j}) . \end{aligned}$$

With the above definitions  $\mathbf{D}(\mathbf{G}(\cdot))$  is the 9-points Laplacian proposed by SÜLI [1991] (cf. Figure 3):

$$\begin{aligned} \mathbf{L}(p) &= \mathbf{D}(\mathbf{G}(p)) \\ &= \frac{1}{8} (\Delta_{xx,i+1/2,j+3/2}(p) + 6\Delta_{xx,i+1/2,j+1/2}(p\bar{c}) + \Delta_{xx,i+1/2,j-1/2}(p)) \\ &\quad + \frac{1}{8} (\Delta_{yy,i+3/2,j+1/2}(p) + 6\Delta_{yy,i+1/2,j+1/2}(p\bar{c}) + \Delta_{yy,i-1/2,j+1/2}(p)) \end{aligned}$$

with

$$\begin{aligned}\Delta_{xx,i+1/2,j+1/2}(p) &:= \frac{1}{\delta x^2} (p_{i+3/2,j+1/2} - 2p_{i+1/2,j+1/2} + p_{i-1/2,j+1/2}) \\ \Delta_{yy,i+1/2,j+1/2}(p) &:= \frac{1}{\delta y^2} (p_{i+1/2,j+3/2} - 2p_{i+1/2,j+1/2} + p_{i+1/2,j-1/2}) .\end{aligned}$$

## A.2 Properties of the Lumping-Operator

**Proposition A.1** For  $p_h \in \mathcal{H}^h$  with  $\nabla p_h \cdot \mathbf{n} \equiv 0$  on  $\partial\Omega$  we have

$$\|\nabla p_h\|_{0,\Omega}^2 \leq -b_{2h}(\nabla p_h, \mathbb{L}p_h)$$

**Proof.** Let us consider a cell-wise representation of  $p_h$ , i.e. on a control volume  $C_{i,j}$  of the primary discretization  $p_h$  can be also represented by

$$p_h(x, y)|_{C_{i,j}} = p_{i,j} + (x - x_i)p_{x,i,j} + (y - y_j)p_{y,i,j} + (x - x_i)(y - y_j)p_{xy,i,j} ,$$

in which  $p_{i,j}$  is the mean value of  $p_h$  on  $C_{i,j}$ , and  $p_{x,i,j}$ ,  $p_{y,i,j}$  and  $p_{xy,i,j}$  are the partial and mixed derivatives of  $p_h$  in  $(x_i, y_j)$ , respectively. With this definition, we have

$$\begin{aligned}[\nabla p_h(x, y)]^2|_{C_{i,j}} &= p_{x,i,j}^2 + 2(y - y_j)p_{x,i,j}p_{xy,i,j} + (y - y_j)^2p_{xy,i,j}^2 \\ &\quad + p_{y,i,j}^2 + 2(x - x_i)p_{y,i,j}p_{xy,i,j} + (x - x_i)^2p_{xy,i,j}^2\end{aligned}$$

Furthermore, we obtain

$$\begin{aligned}\|\nabla p_h\|_{0,\Omega}^2 &= \sum_{i,j} \int_{C_{i,j}} [\nabla p_h]^2 dx \\ &= \delta x \delta y \sum_{i,j} \left( p_{x,i,j}^2 + p_{y,i,j}^2 + \frac{\delta x^2 \delta y^2}{12} p_{xy,i,j}^2 \right)\end{aligned}$$

To compare this result with the expression in the  $b_{2h}(\cdot, \cdot)$  form, the bilinear form has to be written as sum over the primary cells. Using partial summation, this leads to

$$\begin{aligned}b_{2h}(\nabla p_h, \mathbb{L}p_h) &= p_{1/2,1/2} \left[ \frac{\delta y}{2} p_{x,1,1} + \frac{\delta x}{2} p_{y,1,1} - \frac{\delta x^2 + \delta y^2}{8} p_{xy,1,1} \right] \\ &\quad + \sum_{j=1}^{n-1} p_{1/2,j+1/2} \left[ \frac{\delta y}{2} (p_{x,1,j+1} + p_{x,1,j}) + \frac{\delta x}{2} (p_{y,1,j+1} - p_{y,1,j}) + \frac{\delta x^2 + \delta y^2}{8} (-p_{xy,1,j+1} + p_{xy,1,j}) \right] \\ &\quad + p_{1/2,n+1/2} \left[ \frac{\delta y}{2} p_{x,1,n} - \frac{\delta x}{2} p_{y,1,n} + \frac{\delta x^2 + \delta y^2}{8} p_{xy,1,n} \right]\end{aligned}$$

$$\begin{aligned}
& + \sum_{i=1}^{m-1} \left( p_{i+1/2,1/2} \left[ \frac{\delta y}{2} (p_{x,i+1,1} - p_{x,i,1}) + \frac{\delta x}{2} (p_{y,i+1,1} + p_{y,i,1}) + \frac{\delta x^2 + \delta y^2}{8} (-p_{xy,i+1,1} + p_{xy,i,1}) \right] \right. \\
& \quad + \sum_{j=1}^{n-1} p_{i+1/2,j+1/2} \left[ \frac{\delta y}{2} (p_{x,i+1,j+1} - p_{x,i,j+1} + p_{x,i+1,j} - p_{x,i,j}) \right. \\
& \quad \quad + \frac{\delta x}{2} (p_{y,i+1,j+1} + p_{y,i,j+1} - p_{y,i+1,j} - p_{y,i,j}) \\
& \quad \quad \left. \left. + \frac{\delta x^2 + \delta y^2}{8} (-p_{xy,i+1,j+1} + p_{xy,i,j+1} + p_{xy,i+1,j} - p_{xy,i,j}) \right] \right) \\
& \quad + p_{i+1/2,n+1/2} \left[ \frac{\delta y}{2} (p_{x,i+1,n} - p_{x,i,n}) + \frac{\delta x}{2} (-p_{y,i+1,n} - p_{y,i,n}) + \frac{\delta x^2 + \delta y^2}{8} (p_{xy,i+1,n} - p_{xy,i,n}) \right] \\
& + p_{m+1/2,1/2} \left[ -\frac{\delta y}{2} p_{x,m,1} + \frac{\delta x}{2} p_{y,m,1} + \frac{\delta x^2 + \delta y^2}{8} p_{xy,m,1} \right] \\
& + \sum_{j=1}^{n-1} p_{m+1/2,j+1/2} \left[ \frac{\delta y}{2} (-p_{x,m,j+1} - p_{x,m,j}) + \frac{\delta x}{2} (p_{y,m,j+1} - p_{y,m,j}) + \frac{\delta x^2 + \delta y^2}{8} (p_{xy,m,j+1} - p_{xy,m,j}) \right] \\
& + p_{m+1/2,n+1/2} \left[ -\frac{\delta y}{2} p_{x,m,n} - \frac{\delta x}{2} p_{y,m,n} - \frac{\delta x^2 + \delta y^2}{8} p_{xy,m,n} \right] \\
= & \sum_{j=1}^n \left[ \frac{\delta y}{2} p_{x,1,j} (p_{1/2,j-1/2} + p_{1/2,j+1/2}) + \frac{\delta x}{2} p_{y,1,j} (p_{1/2,j-1/2} - p_{1/2,j+1/2}) \right. \\
& \quad \left. + \frac{\delta x^2 + \delta y^2}{8} p_{xy,1,j} (-p_{1/2,j-1/2} + p_{1/2,j+1/2}) \right] \\
& + \sum_{i=1}^{m-1} \left( \sum_{j=1}^n \left[ \frac{\delta y}{2} p_{x,i,j} (-p_{i+1/2,j-1/2} - p_{i+1/2,j+1/2}) + \frac{\delta y}{2} p_{x,i+1,j} (p_{i+1/2,j-1/2} + p_{i+1/2,j+1/2}) \right. \right. \\
& \quad + \frac{\delta x}{2} p_{y,i,j} (p_{i+1/2,j-1/2} - p_{i+1/2,j+1/2}) + \frac{\delta x}{2} p_{y,i+1,j} (p_{i+1/2,j-1/2} - p_{i+1/2,j+1/2}) \\
& \quad + \frac{\delta x^2 + \delta y^2}{8} p_{xy,i,j} (p_{i+1/2,j-1/2} - p_{i+1/2,j+1/2}) \\
& \quad \left. \left. + \frac{\delta x^2 + \delta y^2}{8} p_{xy,i+1,j} (-p_{i+1/2,j-1/2} + p_{i+1/2,j+1/2}) \right] \right) \\
& + \sum_{j=1}^n \left[ \frac{\delta y}{2} p_{x,m,j} (-p_{m+1/2,j-1/2} - p_{m+1/2,j+1/2}) + \frac{\delta x}{2} p_{y,m,j} (p_{m+1/2,j-1/2} - p_{m+1/2,j+1/2}) \right. \\
& \quad \left. + \frac{\delta x^2 + \delta y^2}{8} p_{xy,m,j} (p_{m+1/2,j-1/2} - p_{m+1/2,j+1/2}) \right] \\
= & \sum_{i=1}^m \sum_{j=1}^n \left[ \frac{\delta y}{2} p_{x,i,j} (-p_{i+1/2,j-1/2} - p_{i+1/2,j+1/2} + p_{i-1/2,j-1/2} + p_{i-1/2,j+1/2}) \right. \\
& \quad + \frac{\delta x}{2} p_{y,i,j} (p_{i+1/2,j-1/2} - p_{i+1/2,j+1/2} + p_{i-1/2,j-1/2} - p_{i-1/2,j+1/2}) \\
& \quad \left. + \frac{\delta x^2 + \delta y^2}{8} p_{xy,i,j} (p_{i+1/2,j-1/2} - p_{i+1/2,j+1/2} - p_{i-1/2,j-1/2} + p_{i-1/2,j+1/2}) \right] \\
= & \sum_{i=1}^m \sum_{j=1}^n \left[ \frac{\delta y}{2} p_{x,i,j} (-2\delta x p_{x,i,j}) + \frac{\delta x}{2} p_{y,i,j} (-2\delta y p_{y,i,j}) + \frac{\delta x^2 + \delta y^2}{8} p_{xy,i,j} (-\delta x \delta y p_{xy,i,j}) \right]
\end{aligned}$$

$$= -\delta x \delta y \sum_{i,j} \left( p_{x,i,j}^2 + p_{y,i,j}^2 + \frac{\delta x^2 \delta y^2}{8} p_{xy,i,j}^2 \right)$$

These results lead to the desired estimate:

$$\begin{aligned} \|\nabla p_h\|_{0,\Omega} &= \delta x \delta y \sum_{i,j} \left( p_{x,i,j}^2 + p_{y,i,j}^2 + \frac{\delta x^2 \delta y^2}{12} p_{xy,i,j}^2 \right) \\ &\leq \delta x \delta y \sum_{i,j} \left( p_{x,i,j}^2 + p_{y,i,j}^2 + \frac{\delta x^2 \delta y^2}{8} p_{xy,i,j}^2 \right) \\ &= -b_{2h}(\nabla p_h, \mathbb{L} p_h) \end{aligned} \quad \square$$

**Proposition A.2** For  $p_h \in \mathcal{H}^h$  the estimate

$$\|\mathbb{L} p_h\|_{0,\Omega} \leq C \|p_h\|_{0,\Omega}$$

where  $C$  is a constant, is true.

**Proof.** Since  $p_h$  is piecewise bilinear, its  $L^2$ -norm can be rewritten as

$$\begin{aligned} \|p_h\|_{0,\Omega}^2 &= \int_{\Omega} p_h^2 dx \\ &= \sum_{i,j} \int_{C_{i,j}} [p_{i,j} + (x - x_i)p_{x,i,j} + (y - y_j)p_{y,i,j} + (x - x_i)(y - y_j)p_{xy,i,j}]^2 dx \\ &= \sum_{i,j} \int_{C_{i,j}} [p_{i,j}^2 + (x - x_i)^2 p_{x,i,j}^2 + (y - y_j)^2 p_{y,i,j}^2 + (x - x_i)^2 (y - y_j)^2 p_{xy,i,j}^2] dx \\ &= \delta x \delta y \sum_{i,j} \left[ p_{i,j}^2 + \frac{2 \delta x^2}{3 \cdot 8} p_{x,i,j}^2 + \frac{2 \delta y^2}{3 \cdot 8} p_{y,i,j}^2 + \frac{4 \delta x^2 \delta y^2}{9 \cdot 64} p_{xy,i,j}^2 \right] \\ &= \delta x \delta y \sum_{i,j} \left[ \frac{1}{16} (p_{i+1/2,j+1/2} + p_{i+1/2,j-1/2} + p_{i-1/2,j+1/2} + p_{i-1/2,j-1/2})^2 \right. \\ &\quad + \frac{1}{48} (p_{i+1/2,j+1/2} + p_{i+1/2,j-1/2} - p_{i-1/2,j+1/2} - p_{i-1/2,j-1/2})^2 \\ &\quad + \frac{1}{48} (p_{i+1/2,j+1/2} - p_{i+1/2,j-1/2} + p_{i-1/2,j+1/2} - p_{i-1/2,j-1/2})^2 \\ &\quad \left. + \frac{1}{144} (p_{i+1/2,j+1/2} - p_{i+1/2,j-1/2} - p_{i-1/2,j+1/2} + p_{i-1/2,j-1/2})^2 \right] \end{aligned}$$

$$\begin{aligned}
&= \delta x \delta y \sum_{i,j} \left[ \frac{1}{9} (p_{i+1/2,j+1/2}^2 + p_{i+1/2,j-1/2}^2 + p_{i-1/2,j+1/2}^2 + p_{i-1/2,j-1/2}^2) \right. \\
&\quad + \frac{1}{9} (p_{i+1/2,j+1/2} p_{i+1/2,j-1/2} + p_{i+1/2,j+1/2} p_{i-1/2,j+1/2} \\
&\quad\quad + p_{i+1/2,j-1/2} p_{i-1/2,j-1/2} + p_{i-1/2,j+1/2} p_{i-1/2,j-1/2}) \\
&\quad\quad \left. + \frac{1}{18} (p_{i+1/2,j+1/2} p_{i-1/2,j-1/2} + p_{i+1/2,j-1/2} p_{i-1/2,j+1/2}) \right] \\
&= \frac{\delta x \delta y}{18} \sum_{i,j} \left[ (p_{i+1/2,j+1/2}^2 + p_{i+1/2,j-1/2}^2 + p_{i-1/2,j+1/2}^2 + p_{i-1/2,j-1/2}^2) \right. \\
&\quad + (p_{i+1/2,j+1/2} + p_{i+1/2,j-1/2} + p_{i-1/2,j+1/2} + p_{i-1/2,j-1/2})^2 \\
&\quad \left. - (p_{i+1/2,j+1/2} p_{i-1/2,j-1/2} + p_{i+1/2,j-1/2} p_{i-1/2,j+1/2}) \right]
\end{aligned}$$

Since

$$\begin{aligned}
&p_{i+1/2,j+1/2} p_{i-1/2,j-1/2} + p_{i+1/2,j-1/2} p_{i-1/2,j+1/2} \\
&\leq \frac{1}{2} \left( p_{i+1/2,j+1/2}^2 + p_{i+1/2,j-1/2}^2 + p_{i-1/2,j+1/2}^2 + p_{i-1/2,j-1/2}^2 \right)
\end{aligned}$$

it follows that

$$\begin{aligned}
\|p_h\|_{0,\Omega}^2 &\geq \frac{\delta x \delta y}{18} \sum_{i=1}^m \sum_{j=1}^n \frac{1}{2} (p_{i+1/2,j+1/2}^2 + p_{i+1/2,j-1/2}^2 + p_{i-1/2,j+1/2}^2 + p_{i-1/2,j-1/2}^2) \\
&= \frac{\delta x \delta y}{36} \left[ p_{1/2,1/2}^2 + \sum_{j=1}^{n-1} 2p_{1/2,j+1/2}^2 + p_{1/2,n+1/2}^2 \right. \\
&\quad + 2 \sum_{i=1}^{m-1} \left( p_{i+1/2,1/2}^2 + \sum_{j=1}^{n-1} 2p_{i+1/2,j+1/2}^2 + p_{i+1/2,n+1/2}^2 \right) \\
&\quad \left. + p_{1/2,1/2}^2 + \sum_{j=1}^{n-1} 2p_{1/2,j+1/2}^2 + p_{1/2,n+1/2}^2 \right] \\
&= \frac{1}{9} \|\mathbb{L}p_h\|_{0,\Omega}^2
\end{aligned}$$

□

## References

- ALMGREN, A. S., BELL, J. B., COLELLA, P., HOWELL, L. H., and WELCOME, M. L. [1998]. A Conservative Adaptive Projection Method for the Variable Density Incompressible Navier-Stokes Equations. *Journal of Computational Physics*, 142 (1): pp. 1–46.
- ALMGREN, A. S., BELL, J. B., and CRUTCHFIELD, W. Y. [2000]. Approximate Projection Methods: Part I. Inviscid Analysis. *SIAM Journal on Scientific Computing*, 22 (4): pp. 1139–1159.
- ALMGREN, A. S., BELL, J. B., and SZYMCZAK, W. G. [1996]. A Numerical Method for the Incompressible Navier-Stokes Equations Based on an Approximate Projection. *SIAM Journal on Scientific Computing*, 17 (2): pp. 358–369.
- ANGERMANN, L. [2003]. Node-Centered Finite Volume Schemes and Primal-Dual Mixed Formulations. *Communications in Applied Analysis*, 7 (4): pp. 529–566.
- BABUŠKA, I. [1971]. Error-Bounds for Finite Element Method. *Numerische Mathematik*, 16: pp. 322–333.
- BELL, J. B., COLELLA, P., and GLAZ, H. M. [1989]. A second-order projection method for the incompressible navier-stokes equations. *Journal of Computational Physics*, 85 (2): pp. 257–283.
- BELL, J. B. and MARCUS, D. L. [1992]. A second-order projection method for variable-density flows. *Journal of Computational Physics*, 101: pp. 334–348.
- BERNARDI, C., CANUTO, C., and MADAY, Y. [1988]. Generalized Inf-Sup conditions for the Chebyshev spectral approximation of the Stokes problem. *SIAM Journal on Numerical Analysis*, 25 (6): pp. 1237–1271.
- BRAESS, D. [2003]. *Finite Elemente: Theorie, schnelle Löser und Anwendungen in der Elastizitätstheorie*. Springer, Berlin, 3rd edition.
- BREZZI, F. [1974]. On the existence, uniqueness and approximation of saddle point problems arising from Lagrangian multipliers. *RAIRO Analyse numérique*, 8: pp. 129–151.
- CHORIN, A. J. [1968]. Numerical Solution of the Navier-Stokes Equations. *Mathematics of Computation*, 22 (104): pp. 745–762.
- COURANT, R., FRIEDRICHS, K. O., and LEWY, H. [1928]. Über die partiellen Differenzgleichungen der mathematischen Physik. *Mathematische Annalen*, 100: pp. 32–74.
- GERATZ, K. J. [1997]. *Erweiterung eines Godunov-Typ-Verfahrens für zwei-dimensionale kompressible Strömungen auf die Fälle kleiner und verschwindender Machzahl*. PhD dissertation, Rheinisch-Westfälische Technische Hochschule Aachen.
- GIRAULT, V. and RAVIART, P.-A. [1986]. *Finite Element Methods for Navier-Stokes Equations*, volume 5 of *Springer Series in Computational Mathematics*. Springer, Berlin.
- GRESHO, P. M. [1990]. On the theory of semi-implicit projection methods for viscous incompressible flow and its implementation via a finite element method that also introduces a nearly consistent mass matrix. Part 1: Theory. *International Journal for Numerical Methods in Fluids*, 11 (5): pp. 587–620.
- GRESHO, P. M. and CHAN, S. T. [1990]. On the theory of semi-implicit projection methods for viscous incompressible flow and its implementation via a finite element method that also



- introduces a nearly consistent mass matrix. Part 2: Implementation. *International Journal for Numerical Methods in Fluids*, 11 (5): pp. 621–659.
- GUERMOND, J.-L., MINEV, P., and SHEN, J. [2006]. An overview of projection methods for incompressible flows. *Computer Methods in Applied Mechanics and Engineering*, 195 (44–47): pp. 6011–6045.
- HARLOW, F. H. and WELCH, J. E. [1965]. Numerical Calculation of Time-Dependent Viscous Incompressible Flow of Fluid with Free Surface. *The Physics of Fluids*, 8 (12): pp. 2182–2189.
- VAN KAN, J. [1986]. A Second-Order Accurate Pressure-Correction Scheme for Viscous Incompressible Flow. *SIAM Journal on Scientific and Statistical Computing*, 7 (3): pp. 870–891.
- KLAINERMAN, S. and MAJDA, A. [1981]. Singular limits of quasilinear hyperbolic systems with large parameters and the incompressible limit of compressible fluids. *Communications in Pure Applied Mathematics*, 34: pp. 481–524.
- KLEIN, R. [1995]. Semi-Implicit Extension of a Godunov-Type Scheme Based on Low Mach Number Asymptotics I: One-Dimensional Flow. *Journal of Computational Physics*, 121: pp. 213–237.
- VAN LEER, B. [1979]. Towards the Ultimate Conservative Difference Scheme. V. A Second-Order Sequel to Godunov’s Method. *Journal of Computational Physics*, 32 (1): pp. 101–136.
- MINION, M. L. [1996]. A projection method for locally refined grids. *Journal of Computational Physics*, 127 (1): pp. 158–178.
- NICOLAÏDES, R. A. [1982]. Existence, uniqueness and approximation for generalized saddle point problems. *SIAM Journal on Numerical Analysis*, 19 (2): pp. 349–357.
- OSHER, S. [1985]. Convergence of generalized MUSCL schemes. *SIAM Journal on Numerical Analysis*, 22 (5): pp. 947–961.
- SCHNEIDER, T., BOTTA, N., GERATZ, K. J., and KLEIN, R. [1999]. Extension of Finite Volume Compressible Flow Solvers to Multi-dimensional, Variable Density Zero Mach Number Flows. *Journal of Computational Physics*, 155: pp. 248–286.
- SCHOCHET, S. [1994]. Fast Singular Limits of Hyperbolic PDEs. *Journal of Differential Equations*, 114 (2): pp. 476–512.
- SCHOCHET, S. [2005]. The mathematical theory of low Mach number flows. *RAIRO Modélisation Mathématique et Analyse Numérique*, 39 (3): pp. 441–458.
- SCHULZ-RINNE, C. W. [1993]. *The Riemann problem for two-dimensional gas dynamics and new limiters for high-order schemes*. PhD dissertation, Eidgenössische Technische Fachhochschule (ETH) Zürich. Diss. ETH No. 10297.
- SÜLI, E. [1991]. Convergence of finite volume schemes for Poisson’s equation on nonuniform meshes. *SIAM Journal on Numerical Analysis*, 28 (5): pp. 1419–1430.
- TEMAM, R. [1968]. Une méthode d’approximation de la solution des équations de Navier-Stokes. *Bulletin de la Société Mathématique de France*, 96: pp. 115–152.
- THOMAS, J.-M. and TRUJILLO, D. [1999]. Mixed finite volume methods. *International Journal for Numerical Methods in Engineering*, 46 (9): pp. 1351–1366.
- VATER, S. [2005]. A New Projection Method for the Zero Froude Number Shallow Water Equations. PIK Report 97, Potsdam Institute for Climate Impact Research.

VAN DER VORST, H. A. [1992]. Bi-CGSTAB: A fast and smoothly converging variant of Bi-CG for the solution of nonsymmetric linear systems. *SIAM Journal on Scientific and Statistical Computing*, 13 (2): pp. 631–644.