

## STUDY OF PHYSICS-BASED PRECONDITIONING WITH HIGH-ORDER GALERKIN DISCRETIZATION FOR HYPERBOLIC WAVE PROBLEMS

CLÉMENTINE COURTÈS<sup>1</sup>, EMMANUEL FRANCK<sup>2</sup>, PHILIPPE HELLUY<sup>2</sup> AND HERBERT OBERLIN<sup>3</sup>

**Abstract.** In this article, we detail the construction of a physics-based preconditioner. The Schur decomposition is the key point of the method which is tested on two hyperbolic systems : acoustic wave equations and shallow water equations without source term. Some conserved properties between preconditioner and initial operator are discussed, especially the propagation speeds of a plane wave.

**Résumé.** Dans cet article, nous détaillons la construction d'un préconditionneur basé sur la physique sous-jacente des équations considérées. La décomposition de Schur est le point clé de la méthode, qui sera testée sur deux systèmes hyperboliques : les équations des ondes acoustiques et les équations de Saint Venant sans terme source. Nous étudions certaines propriétés conservées entre le préconditionneur et l'opérateur initial, notamment les vitesses de propagation d'une onde plane.

### INTRODUCTION

Hyperbolic systems are able to model complex physics through nonlinear conservation laws. However, describing complex physical phenomena can be difficult, since those problems prove to be strongly multi-scaled. A good example is the passive advection of pollutant in a river or coastal areas. The flow is modeled by the shallow water equations :

$$\begin{cases} \partial_t h + \nabla \cdot (h \mathbf{u}) = 0, & (1a) \\ \partial_t (h \mathbf{u}) + \nabla \cdot (h \mathbf{u} \otimes \mathbf{u}) + \nabla p = -gh \nabla b, & (1b) \end{cases}$$

with  $g$  the gravitational constant,  $h$  the height of the fluid,  $p$  its pressure (defined by  $p = \frac{gh^2}{2}$ ) and  $\mathbf{u}$  its velocity. This hyperbolic system, first obtained by [11], governs morphodynamics flows caused by the movement of a fluid in contact with the bottom topography  $b$ . It can be derived from the Navier-Stokes equation as proved by Gerbeau and Perthame in [15].

A third equation is considered to take into account the passive transport of the pollutant :

$$\partial_t c + \mathbf{u} \cdot \nabla c = s, \quad (2)$$

with  $c$  the concentration of the pollutant and  $s$  the source term. Without source term in (1b), we recognize Euler isentropic equations coupled to transport of suspended matter such as volcanic dust during an eruption

<sup>1</sup> LMO, Université Paris Sud, Orsay, France

<sup>2</sup> INRIA Nancy Grand-Est and IRMA Strasbourg, TONUS Team, France

<sup>3</sup> Max-Planck-Institut für Plasmaphysik, Boltzmannstraße 2D-85748 Garching, Germany

or radioactive particles during a nuclear accident.

In one dimension, the system (1)-(2) is hyperbolic with the three eigenvalues  $u + \sqrt{gh}$ ,  $u - \sqrt{gh}$  (for Saint-Venant system) and  $u$  (for the transport equation), see [1]. If  $u \ll \sqrt{gh}$  (small Froude number), the characteristic time for flow and pollutant transfert are very different and two time scales have to be considered. The time step for pollutant is determined thanks to the Courant-Friedrichs-Lewy condition  $(|u| + \sqrt{gh})\Delta t \leq \Delta x$  whereas gravity waves' speeds are  $\sqrt{gh}$ . Arises then new layers of numerical complications, which are usually dealt with using implicit schemes : for instance, solving implicitly the Saint-Venant part may be a good option.

Sediment transport problem is a second example of a multi-scale problem, and may be described by the shallow water system coupled to Exner equation, as explained in [16] or [18]. In this system, the topology  $b$  of shallow water equations (1) is described by Exner equation

$$\partial_t b + \zeta \nabla \cdot \mathbf{q} = 0, \quad (3)$$

with  $\mathbf{q} = \mathbf{q}(\mathbf{u})$  the sediment flux and  $\zeta$  a constant which depends on the sediment coefficient porosity. Again, gravity waves' whose speeds are  $\sqrt{gh}$  and sedimentation behaviors range on different time-scales, which leads Bilanceri *et al* in [3] to design an implicit time-advancing scheme.

**Main purpose** To take into account the ill-conditioning of the system, we aim to study the efficiency of a preconditioned implicit algorithm for hyperbolic systems with Continuous Galerkin high-order method proposed in [7], [9], [8] and to understand the difficulties associated. In this paper, we focus the study of this method on two simpler models than the complete morphodynamics system (1). The first one corresponds to the so-called acoustic wave equations (4) and models the propagation through a material medium of acoustic waves

$$\begin{cases} \partial_t p + c \nabla \cdot \mathbf{u} = 0, \\ \partial_t \mathbf{u} + c \nabla p = 0, \end{cases} \quad (4)$$

where  $p$  stands for the pressure,  $\mathbf{u}$  the velocity and  $c$  the sound speed. The second one is composed of shallow water equations without source term.

$$\begin{cases} \partial_t h + \nabla \cdot (h \mathbf{u}) = 0, \\ \partial_t (h \mathbf{u}) + \nabla \cdot (h \mathbf{u} \otimes \mathbf{u}) + \nabla p = 0. \end{cases} \quad (5)$$

**Numerical problem** Implicit schemes make use of the inversion of a linear system through iterative solvers (exact solvers being too greedy for 2D or 3D problems). However, waves' speeds may correspond to different scales (for example slow and fast MHD waves) which make the ratio between the smallest eigenvalue of the implicit matrix and the biggest one blow up. Moreover, one of the waves' speeds may converge to zero and so, for large time step, the state may be close to the stationary one : discrete model has thus eigenvalues near zero (*e.g.* in low Mach or low Froude number regime). To explain the second case, we propose to use the wave acoustic equation (closed to the shallow water or Euler equations in the low Mach or low Froude regime) for a large time step. The implicit system obtained can be written

$$\begin{cases} p + \frac{1}{\epsilon} \nabla \cdot \mathbf{u} = p_0, \\ \mathbf{u} + \frac{1}{\epsilon} \nabla p = \mathbf{u}_0, \end{cases} \quad (6)$$

with a small parameter  $\epsilon$  which corresponds to large time step or large sound velocity. When  $\epsilon$  is small, the limit model is given by

$$\begin{cases} \nabla \cdot \mathbf{u} = 0, \\ \nabla p = 0. \end{cases} \quad (7)$$

For this limit system, with standard boundary conditions, there is generally no unique solution. One of the problems is that, when  $\epsilon$  is small, the full system tends to a system which doesn't have uniqueness of the solution. Consequently, the condition number increases when  $\epsilon$  tends to zero. This non-uniqueness problem can be highly amplified by a non adapted space discretization. To illustrate, we consider the hydrostatic mode for the wave problem which corresponds to a stationary solution having a constant pressure and a zero velocity. Its existence is due to the fact that the pressure is determined within an arbitrary constant of integration associated with the solution of the continuous system. For the continuous mathematical problem, the hydrostatic mode is the only one present. However, this may not be the case in the discretized problem. For example, additional numerical modes may occur when the velocity field and the pressure are not discretized in the good function spaces. These spurious pressure modes correspond to stationary but non constant (in space) pressures associated with zero velocities. For this family of spurious modes, the discrete pressure belongs to the null space of the discrete gradient operator and solution uniqueness is lost since any multiple of a spurious mode can be added to any solution of the discrete equations and still satisfy them (this problem is explained in [6]). A way of addressing this problem is to use the compatible finite element spaces given by the Inf-Sup estimation. In this paper we apply our preconditioner to a standard continuous Galerkin discretization, which does not satisfy the inf-sup condition. Of course it could also be applied to a smarter space discretization and then would give even better results. To address the problem (6), we propose the same guidelines as those used in compressible resistive magnetohydrodynamics model (MHD) in order to derive a physics-based preconditioner (outlined first in [9]). The idea is to find a problem which approximates the solution of (6) for a given  $\epsilon$  but for which the implicit operator has a condition number independent of  $\epsilon$  and is easy to solve with classical solvers as multi-grids. Using this approximate problem as a preconditioning, we normally obtain at the end an efficient solver. Moreover, for lot of applications, it is interesting to use high-order spatial method which captures some fine scales. It is known that high-order methods generate ill-conditioned matrices (see, for example, the Table 5.2 in [2]). This is why, implicit scheme often works hand-in-hand with efficient preconditioners because of the ill-conditioning of the multi-scaled (or stiff-) hyperbolic systems with high-order representation.

**Outline** At the beginning, we gather results on acoustic wave equations (4) in Section 1. Describing the continuous study model and its propagation wave velocities should be the goal of Subsection 1.1, whereas Subsection 1.2 condenses time and spatial discretization (*i.e.* the  $\theta$ -scheme and continuous Galerkin method). A significant proportion of the first section is devoted to the study of the physics-based preconditioner : Subsections 1.3 and 1.4. Eventually, numerical results are summarized in Subsection 1.5. Shallow water equations (5) is the topic of Section 2. Subsection 2.1 consists in studying standard equations before a linearization and its discretization (Subsection 2.3). Last but not least, Subsections 2.4, 2.5 and 2.6 collect all results on wave propagation for the Schur complement. Subsection 2.7 gathers numerical results for shallow water equations.

## 1. ACOUSTIC WAVE EQUATIONS

In this section, acoustic wave equations (4) are considered in 2D : the velocity  $\mathbf{u}$  is defined in both directions (*i.e.*  $\mathbf{u} = (u_1, u_2)$ ). A feature of those equations is that the Schur complement may be computed either on the velocity or on the pressure, which leads to exhibit two possible preconditioners. Both will be studied in this section.

### 1.1. Study model for wave equations

**Boundary conditions** We propose to add some admissible boundary conditions to acoustic wave equations (4). Three admissible boundary conditions are available :

- case (a) :  $\mathbf{u} \cdot \mathbf{n} \equiv \mathbf{0}$  on  $\partial\Omega$  with  $\mathbf{n}$  the unit normal vector at the boundary,
- case (b) :  $p \equiv 0$  on  $\partial\Omega$ ,
- case (c) :  $p - \mathbf{u} \cdot \mathbf{n} = p^0 - \mathbf{u}^0 \cdot \mathbf{n}$  on  $\partial\Omega$ .

The three cases are studied in this paragraph.

**Proposition 1.1.** *For all the boundary conditions proposed, the wave model admits a unique solution.*

*Proof.* To prove the uniqueness, we compute the energy estimate associated with the model. We multiply the first equation of system (4) by  $p$ , the second one by  $\mathbf{u}$  and add both equations to obtain

$$\frac{1}{2} \partial_t \int_{\Omega} (p^2 + \mathbf{u} \cdot \mathbf{u}) \, d\mathbf{x} = - \int_{\Omega} c \nabla \cdot (p\mathbf{u}) \, d\mathbf{x} = - \int_{\partial\Omega} cp\mathbf{u} \cdot \mathbf{n} \, d\sigma,$$

and so is the energy conserved in the cases (a) or (b). For the third case, the previous equation may be rewritten in the following form

$$\frac{1}{2} \partial_t \int_{\Omega} (p^2 + \mathbf{u} \cdot \mathbf{u}) \, d\mathbf{x} = -\frac{c}{4} \int_{\partial\Omega} (p + \mathbf{u} \cdot \mathbf{n})^2 \, d\sigma + \frac{c}{4} \int_{\partial\Omega} (p - \mathbf{u} \cdot \mathbf{n})^2 \, d\sigma \leq \frac{c}{4} \int_{\partial\Omega} (p_0 - \mathbf{u}_0 \cdot \mathbf{n})^2 \, d\sigma. \quad (8)$$

We consider two solutions  $(p_1, \mathbf{u}_1)$  and  $(p_2, \mathbf{u}_2)$  with the same initial data and the same boundary conditions and denote  $\bar{p} = p_1 - p_2$  and  $\bar{\mathbf{u}} = \mathbf{u}_1 - \mathbf{u}_2$ . Those variables satisfy the same system (4) with boundary conditions equal to zero. Using the previous energy estimate, we obtain that

$$\frac{1}{2} \partial_t \int_{\Omega} (\bar{p}^2 + \bar{\mathbf{u}} \cdot \bar{\mathbf{u}}) \, d\mathbf{x} \leq 0.$$

Since initial energy vanishes, the energy vanishes for all time. Consequently the solution is unique.  $\square$

**Remark 1.2.** The existence of a solution is based on the Hille-Yosida theorem [5].

**Remark 1.3.** The classical boundary conditions consist in imposing a matrix system on  $\partial\Omega$  and the aim of this remark is to rewrite boundary conditions in this standard form. Let us define  $\mathbf{V} = (p, u_1, u_2)^t$  to rewrite acoustic wave equations (4) on a classical hyperbolic form

$$\partial_t \mathbf{V} + \underbrace{\begin{pmatrix} 0 & c & 0 \\ c & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}}_{=:A_1} \partial_x \mathbf{V} + \underbrace{\begin{pmatrix} 0 & 0 & c \\ 0 & 0 & 0 \\ c & 0 & 0 \end{pmatrix}}_{=:A_2} \partial_y \mathbf{V} = 0. \quad (9)$$

According to [4,10], the suitable boundary conditions are equivalent to the following relation

$$\mathbb{B} \begin{pmatrix} p \\ \mathbf{u} \end{pmatrix} = \mathbb{B} \begin{pmatrix} p^0 \\ \mathbf{u}^0 \end{pmatrix}, \quad \text{on } \partial\Omega, \quad (10)$$

where the matrix  $\mathbb{B}$  is chosen such that, for any exterior unit normal vector  $\mathbf{n} = (n_1, n_2)$ , the matrix

$$\frac{1}{2} (n_1 A_1 + n_2 A_2) + \mathbb{B}$$

must be non negative, with  $A_1$  and  $A_2$  defined in the generic hyperbolic form (9). By adapting computations of [4,10] to our current problem, the corresponding boundary conditions matrix is defined such as  $\mathbb{B} = -PD^-P^{-1}$ , where  $D^-$  is the non positive part of the diagonal matrix which appears in the diagonalization of  $n_1 A_1 + n_2 A_2$  and  $P$  the basis change matrix (*i.e.*  $n_1 A_1 + n_2 A_2 = PD^-P^{-1} + PD^+P^{-1}$ ). A straightforward computation gives

$$\mathbb{B} = \frac{1}{2} \begin{pmatrix} c & -cn_1 & -cn_2 \\ -cn_1 & cn_1^2 & cn_1 n_2 \\ -cn_2 & cn_1 n_2 & cn_2^2 \end{pmatrix}$$

and thus, combining definition of  $\mathbb{B}$  and Condition (10) leads to

$$\mathbb{B} \begin{pmatrix} p \\ \mathbf{u} \end{pmatrix} = \frac{1}{2} \begin{pmatrix} c(p - \mathbf{u} \cdot \mathbf{n}) \\ -n_1 c(p - \mathbf{u} \cdot \mathbf{n}) \\ -n_2 c(p - \mathbf{u} \cdot \mathbf{n}) \end{pmatrix} = \frac{1}{2} c(p^0 - \mathbf{u}^0 \cdot \mathbf{n}) \begin{pmatrix} 1 \\ -n_1 \\ -n_2 \end{pmatrix}.$$

We recognize the third condition (case (c)) :  $p - \mathbf{u} \cdot \mathbf{n} = p^0 - \mathbf{u}^0 \cdot \mathbf{n}$  on  $\partial\Omega$ . For case (a) and (b) the boundary conditions matrix  $\mathbb{B}$  is given by

$$\mathbb{B}_{\mathbf{u}} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & cn_1^2 & cn_1 n_2 \\ 0 & cn_1 n_2 & cn_2^2 \end{pmatrix} \quad \text{and} \quad \mathbb{B}_p = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

**Propagation velocities : eigenvalues, Riemann invariants and advection equations** The propagation velocities correspond to advection velocities for Riemann invariants.

**Definition 1.4** (Plane wave). A generic plane wave is defined by

$$\mathbf{V}(t, \mathbf{x}) = \mathbf{V}^0 f(\mathbf{k} \cdot \mathbf{x} - \omega t), \quad (11)$$

with  $\mathbf{k} = (k_1, k_2)$ ,  $\mathbf{x} = (x, y)$  and  $\mathbf{V}^0$  a vector independent of  $t$  and  $\mathbf{x}$ .

The idea is to consider the solution of (9) as a plane wave, which leads to the generic system

$$A(\mathbf{k}, \omega) \mathbf{V}^0 = 0, \quad (12)$$

where  $\mathbf{V}^0 = (p^0, u_1^0, u_2^0)^t$ . A condition to have a non-trivial solution is for the kernel of  $A(\mathbf{k}, \omega)$  to be non-trivial. Having both an eigenvalue equals to zero (which leads to the dispersion relation) and the associated eigenvectors (which leads to the Riemann invariants) ensures a non-trivial kernel.

**Proposition 1.5.** *The Riemann invariants are  $(0, -k_2, k_1)^t$ ,  $(1, \frac{k_1}{\|\mathbf{k}\|}, \frac{k_2}{\|\mathbf{k}\|})^t$  and  $(1, -\frac{k_1}{\|\mathbf{k}\|}, -\frac{k_2}{\|\mathbf{k}\|})^t$ , traveling respectively with a null velocity, a velocity equals to  $c\|\mathbf{k}\|$  and a velocity equals to  $-c\|\mathbf{k}\|$ .*

Before proving Proposition 1.5, we fix the convention for the space-time Fourier transform of a function  $\mathbf{v}$

$$\mathcal{F}(\mathbf{v})(\omega, \mathbf{k}) = \int_{\mathbb{R}^3} \mathbf{v}(t, \mathbf{x}) e^{i(\mathbf{k} \cdot \mathbf{x} - \omega t)} d\mathbf{x} dt, \quad \forall (\mathbf{k}, \omega) \in \mathbb{R}^2 \times \mathbb{R}.$$

*Proof.* From system (4), we apply the Fourier transform in space and time to obtain  $A_{\text{wave}}(\mathbf{k}, \omega) \mathcal{F}(\mathbf{V}) = 0$  (equation equivalent to the generic equation (12)) with

$$A_{\text{wave}}(\mathbf{k}, \omega) = \begin{pmatrix} i\omega & -ic k_1 & -ic k_2 \\ -ic k_1 & i\omega & 0 \\ -ic k_2 & 0 & i\omega \end{pmatrix}.$$

and the corresponding spectrum

$$\sigma(A_{\text{wave}}) = \{i\omega, i(\omega - c\|\mathbf{k}\|), i(\omega + c\|\mathbf{k}\|)\}.$$

The plane wave is not equal to zero if  $\text{Ker}(A_{\text{wave}}) \neq \{\mathbf{0}\}$ . Canceling one of the previous eigenvalues leads to the following dispersion relations

$$\omega = 0, \quad \omega = c\|\mathbf{k}\|, \quad \omega = -c\|\mathbf{k}\|$$

and the associated kernel space of  $A_{\text{wave}}$  (cf Table 1).  $\square$

Dispersion relations (canceling the eigenvalue)	Riemann invariants (basis vector of the eigenspace)
$\omega = 0$	$\text{Ker}(A_{\text{wave}}) = \text{Span} \left\{ \begin{pmatrix} 0 \\ -k_2 \\ k_1 \end{pmatrix} \right\}$
$\omega = c\ \mathbf{k}\ $	$\text{Ker}(A_{\text{wave}}) = \text{Span} \left\{ \begin{pmatrix} 1 \\ \frac{k_1}{\ \mathbf{k}\ } \\ \frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$
$\omega = -c\ \mathbf{k}\ $	$\text{Ker}(A_{\text{wave}}) = \text{Span} \left\{ \begin{pmatrix} 1 \\ -\frac{k_1}{\ \mathbf{k}\ } \\ -\frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$

TABLE 1. [Acoustic wave equations] Dispersion relations and Riemann invariants for the acoustic wave model

**Remark 1.6.** Instead of using the Fourier transform, we could have studied directly the jacobian matrix. However, this method will not work for a second order operator (Subsection 1.4). In order to simplify the comparison, we have chosen to use the same proof for any operator.

## 1.2. Time and spatial discretization

**Time discretization** In low Mach number regime for instance, we would like to define time step  $\Delta t$  such as  $|\mathbf{u}|\Delta t \sim \Delta x$  which implies  $c\Delta t \gg \Delta x$ . This is why, the classical explicit scheme for this problem is not adapted for large sound speed because of the restriction on the time step  $\Delta t$  for this scheme. Indeed, the time step is constrained by this velocity. Therefore, using an implicit scheme (and more generally a  $\theta$ -scheme) is an interesting way to circumvent the problem.

Let us define  $\Delta t$  (resp.  $\Delta x$ ) the constant time step (resp. the uniform size mesh) and  $p^n$  (resp.  $\mathbf{u}^n$ ) the pressure (resp. the velocity) at time  $t^n = n\Delta t$ . The selected  $\theta$ -scheme together with the discretized boundary conditions write as follows, with  $\theta \in [0, 1]$ ,

$$\begin{cases} p^{n+1} + \theta c \Delta t \nabla \cdot \mathbf{u}^{n+1} = p^n - (1 - \theta)c \Delta t \nabla \cdot \mathbf{u}^n, & \text{in } \Omega, \\ \mathbf{u}^{n+1} + \theta c \Delta t \nabla p^{n+1} = \mathbf{u}^n - (1 - \theta)c \Delta t \nabla p^n, & \text{in } \Omega, \\ \mathbf{u}^{n+1} \cdot \mathbf{n} \equiv 0 \text{ or } p^{n+1} \equiv 0 \text{ or } p^{n+1} - \mathbf{u}^{n+1} \cdot \mathbf{n} = p^0 - \mathbf{u}^0 \cdot \mathbf{n}, & \text{on } \partial\Omega. \end{cases} \quad (13)$$

**Remark 1.7.** A good choice of  $\theta$  enables us to bypass the problem of a too restrictive CFL condition. Namely, the Crank-Nicholson scheme (associated to  $\theta = \frac{1}{2}$ ) is unconditionally-stable and of second-order.

**Continuous Galerkin scheme and spatial discretization** This paragraph is devoted to the spatial discretization used : the Continuous Galerkin scheme, which is based on a polynomial approximation of the weak formulation of the equations.

Let us take two test functions  $\nu_p, \boldsymbol{\nu}_u = (\nu_{u_1}, \nu_{u_2})$ , smooth enough to obtain weak form of the problem by multiplying each equation of (13) by the corresponding test function and performing integrations by part. Eventually, the weak form yields

$$\begin{cases} \int_{\Omega} p^{n+1} \nu_p - c \theta \Delta t \left( \int_{\Omega} \mathbf{u}^{n+1} \cdot \nabla \nu_p - \int_{\partial\Omega} (\mathbf{u}^{n+1} \cdot \mathbf{n}) \nu_p \right) = \int_{\Omega} p^n \nu_p + (1 - \theta)c \Delta t \left( \int_{\Omega} \mathbf{u}^n \cdot \nabla \nu_p - \int_{\partial\Omega} (\mathbf{u}^n \cdot \mathbf{n}) \nu_p \right), \\ \int_{\Omega} \mathbf{u}^{n+1} \cdot \boldsymbol{\nu}_u - c \theta \Delta t \left( \int_{\Omega} p^{n+1} \nabla \cdot \boldsymbol{\nu}_u - \int_{\partial\Omega} (\boldsymbol{\nu}_u \cdot \mathbf{n}) p^{n+1} \right) = \int_{\Omega} \mathbf{u}^n \cdot \boldsymbol{\nu}_u + (1 - \theta)c \Delta t \left( \int_{\Omega} p^n \nabla \cdot \boldsymbol{\nu}_u - \int_{\partial\Omega} (\boldsymbol{\nu}_u \cdot \mathbf{n}) p^n \right). \end{cases}$$

To this weak formulation, we join the boundary term associated with the matrix  $\mathbb{B}$ . Since we are enforcing these boundaries through penalization, we multiply this term by a large coefficient  $\frac{1}{\epsilon}$ .

Let us define the spatial degree of freedom  $j$  as Gauss-Lobatto grid point  $\mathbf{x}_j$  (*i.e.* the endpoints of the spatial interval are included in the set of  $\{\mathbf{x}_j\}_{j \in \llbracket 1, J \rrbracket}$ ). We denote  $\Phi_j$  the associated basis function. The approximate pressure decomposes on this basis thanks to  $p^n = \sum_{j=1}^J p_j^n \Phi_j(\mathbf{x})$ , and the same holds true for the velocity. To simplify the notations, we denote  $p_\Delta^n$  the vector of size  $J$  defined by  $p_\Delta^n = (p_1^n, \dots, p_J^n)$  and so does the same for  $\mathbf{u}_\Delta^n = (u_{1,\Delta}^n, u_{2,\Delta}^n)$ . Lastly, test functions are defined such as  $\nu_p = \nu_{u_1} = \nu_{u_2} = \Phi_i$ , for  $i \in \llbracket 1, J \rrbracket$ , to obtain the following matrix system,  $J_{ac} V_\Delta^{n+1} = \overline{B} V_\Delta^n + \text{Boundary Conditions}$ , given by

$$\begin{aligned} & \left( \begin{pmatrix} M & 0 & \theta(\mathcal{D}_1 + B_p(n_1)) \\ 0 & M & \theta(\mathcal{D}_2 + B_p(n_2)) \\ \theta(\mathcal{D}_1 + B_{\mathbf{u}}(n_1)) & \theta(\mathcal{D}_2 + B_{\mathbf{u}}(n_2)) & M \end{pmatrix} + \frac{1}{\epsilon} \mathbb{B}_\Delta \right) \begin{pmatrix} u_{1,\Delta}^{n+1} \\ u_{2,\Delta}^{n+1} \\ p_\Delta^{n+1} \end{pmatrix} \\ &= \begin{pmatrix} M & 0 & -(1-\theta)(\mathcal{D}_1 + B_p(n_1)) \\ 0 & M & -(1-\theta)(\mathcal{D}_2 + B_p(n_2)) \\ -(1-\theta)(\mathcal{D}_1 + B_{\mathbf{u}}(n_1)) & -(1-\theta)(\mathcal{D}_2 + B_{\mathbf{u}}(n_2)) & M \end{pmatrix} \begin{pmatrix} u_{1,\Delta}^n \\ u_{2,\Delta}^n \\ p_\Delta^n \end{pmatrix} + \frac{1}{\epsilon} \mathbb{B}_\Delta \begin{pmatrix} u_{1,\Delta}^0 \\ u_{2,\Delta}^0 \\ p_\Delta^0 \end{pmatrix} \end{aligned} \quad (14)$$

with the matrices' definitions :  $M$  the mass matrix  $M_{i,j} = \int_\Omega \Phi_i \Phi_j$ ,  $(\mathcal{D}_1)_{i,j} = -c \Delta t (\int_\Omega \Phi_j \partial_x \Phi_i)$ ,  $(\mathcal{D}_2)_{i,j} = -c \Delta t (\int_\Omega \Phi_j \partial_y \Phi_i)$ ,  $(B_p(n_k))_{i,j} = c \Delta t (\int_{\partial\Omega} \Phi_i \Phi_j n_k)$  and  $(B_{\mathbf{u}}(n_k))_{i,j} = c \Delta t (\int_{\partial\Omega} n_k \Phi_i \Phi_j)$  and  $\mathbb{B}_h = (\mathbb{B}_h^{k_1 k_2})_{k_1, k_2}$  with  $(\mathbb{B}_h^{k_1 k_2})_{i,j} = \int_{\partial\Omega} \mathbb{B}_{k_1, k_2} \Phi_i \Phi_j$ .

**Remark 1.8.** Using Gauss-Lobatto points is an appropriate choice to simplify computations, because it converts the mass matrix to a diagonal matrix, easy to invert. Moreover, Gauss-Lobatto quadrature rule is accurate for polynomials up to degree  $2J - 1$  ( however the integration is not exact for these polynomial), where  $J$  is the number of grid points [12].

### 1.3. Preconditioning for the wave equation

According to the inf-sup condition for Friedrichs' systems [13, 14], the time-discrete variational formulation of the wave equation is well-posed for  $\mathbf{u}$  in  $\mathbb{H}(\text{div}, \Omega)$  and  $p$  in  $\mathbb{H}^1(\Omega)$ . The space  $\mathbb{H}(\text{div}, \Omega)$  stands for  $\{\mathbf{v} \in \mathbb{L}^2(\Omega), \text{ such that } \nabla \cdot \mathbf{v} \in \mathbb{L}^2(\Omega)\}$ . The well-posedness ensures the existence and the uniqueness of the solution and thus the jacobian matrix inversion. However, in practice, this uniqueness may be lost for  $\Delta t$  and  $\Delta x$  asymptotically large, but we will not go into details.

In this subsection, instead of inverting directly the jacobian matrix  $J_{ac}$  in equation (14), we plan on constructing a preconditioning for the acoustic wave equations (4). The main idea here is to treat the stiffness of our problem by designing a preconditioner using a simplified form of the equations. We adjust the guidelines used by Chacon for the resistive MHD model [7] : we design a diffusive operator which is time-consistent with a second order hyperbolic operator. According to Chacon's terminology, those operator will be called "parabolic".

This new parabolic equation provides us with the preconditioner and can be solved easily with multi-grid methods or preconditioning conjugate gradient [7]. As explained earlier, this preconditioner may be determined for the pressure or for the velocity.

**Schur complement on the pressure** To find relative parabolic form of the equations, we start off from a semi-discrete scheme (Equation (13)) summarized in the compact form

$$\begin{cases} \begin{pmatrix} \mathbb{D}_2 & L \\ U & D_1 \end{pmatrix} \begin{pmatrix} \mathbf{u}^{n+1} \\ p^{n+1} \end{pmatrix} = \begin{pmatrix} \mathbf{R}_u \\ \mathbf{R}_p \end{pmatrix}, & \text{in } \Omega, \\ \mathbf{u}^{n+1} \cdot \mathbf{n} \equiv 0 \text{ or } p^{n+1} \equiv 0 \text{ or } p^{n+1} - \mathbf{u}^{n+1} \cdot \mathbf{n} = p^0 - \mathbf{u}^0 \cdot \mathbf{n}, & \text{on } \partial\Omega, \end{cases} \quad (15)$$

$$(16)$$

with  $D_1 = I_1$ ,  $\mathbb{D}_2 = \begin{pmatrix} I_1 & 0 \\ 0 & I_1 \end{pmatrix}$ ,  $U = \theta c \Delta t \nabla \cdot I_2$ ,  $L = \theta c \Delta t \nabla I_1$  and  $\begin{pmatrix} \mathbf{R}_u \\ R_p \end{pmatrix} = \begin{pmatrix} \mathbf{u}^n - (1-\theta)c \Delta t \nabla p^n \\ p^n - (1-\theta)c \Delta t \nabla \cdot \mathbf{u}^n \end{pmatrix}$ .

In the previous equation,  $I_d$  stands for the  $d \times d$ -identity matrix. We have denoted  $\nabla I_1 = \begin{pmatrix} \partial_x \\ \partial_y \end{pmatrix}$  the gradient operator,  $\Delta I_1 = \partial_x^2 + \partial_y^2$  the Laplacien operator and  $\nabla \cdot I_2 = \text{tr} \begin{pmatrix} \partial_x & 0 \\ 0 & \partial_y \end{pmatrix}$  the divergence operator (with  $\text{tr}$  the trace of the matrix).

Let us first focus on the first system (15). The so-called Schur block decomposition applies and the previous matrix system turns into

$$\begin{pmatrix} I_2 & 0 \\ U \mathbb{D}_2^{-1} & I_1 \end{pmatrix} \begin{pmatrix} \mathbb{D}_2 & 0 \\ 0 & P_{p,\text{schur}} \end{pmatrix} \begin{pmatrix} I_2 & \mathbb{D}_2^{-1} L \\ 0 & I_1 \end{pmatrix} \begin{pmatrix} \mathbf{u}^{n+1} \\ p^{n+1} \end{pmatrix} = \begin{pmatrix} \mathbf{R}_u \\ R_p \end{pmatrix},$$

where

$$P_{p,\text{schur}} = D_1 - U \mathbb{D}_2^{-1} L = I_1 - c^2 \theta^2 \Delta t^2 \Delta I_1 \quad (17)$$

stands for the so-called Schur complement on the pressure. Solving the system (15) with this decomposition is equivalent to the following algorithm

$$\text{(Pressure algorithm)} \begin{cases} \text{Predictor : } \mathbb{D}_2 \mathbf{u}^* = \mathbf{R}_u, \\ \text{Potential evolution : } P_{p,\text{schur}} p^{n+1} = -U \mathbf{u}^* + R_p, \\ \text{Corrector : } \mathbb{D}_2 \mathbf{u}^{n+1} = \mathbb{D}_2 \mathbf{u}^* - L p^{n+1}. \end{cases}$$

Now it is necessary to add boundary conditions for each operator. Using the second equation of (4) we obtain that  $\partial_t(\mathbf{u} \cdot \mathbf{n}) = -c(\nabla p \cdot \mathbf{n})$

- case (a) :  $\mathbf{u} \cdot \mathbf{n} \equiv \mathbf{0}$  on  $\partial\Omega$  : Dirichlet boundary condition  $(\mathbf{u} \cdot \mathbf{n}) = 0$  on  $\mathbb{D}_2$  and homogeneous Neumann  $(\nabla p^{n+1} \cdot \mathbf{n}) = 0$  for  $P_{p,\text{schur}}$ ,
- case (b) :  $p \equiv 0$  on  $\partial\Omega$  : no boundary condition on  $\mathbb{D}_2$  and homogeneous Dirichlet for  $P_{p,\text{schur}}$ ,
- case (c) :  $p - \mathbf{u} \cdot \mathbf{n} = g$  on  $\partial\Omega$  with  $g$  constant in time : Dirichlet boundary condition  $(\mathbf{u}^* \cdot \mathbf{n}) = p^n - g$  (prediction) and  $(\mathbf{u}^{n+1} \cdot \mathbf{n}) = p^{n+1} - g$  (correction) and Robin boundary condition  $c \Delta t (\nabla p^{n+1} \cdot \mathbf{n}) + p^{n+1} = p^n$  for  $P_{p,\text{schur}}$ .

**Remark 1.9.** By splitting the second equation of the semi-discrete scheme (13), we recover the same result as Schur's theory

$$\begin{cases} \frac{\mathbf{u}^* - \mathbf{u}^n}{\Delta t} = -(1-\theta)c \nabla p^n, \\ \frac{p^{n+1} - p^n}{\Delta t} + c \theta \nabla \cdot \mathbf{u}^{n+1} = -(1-\theta)c \nabla \cdot \mathbf{u}^n, \\ \frac{\mathbf{u}^{n+1} - \mathbf{u}^*}{\Delta t} + c \theta \nabla p^{n+1} = 0. \end{cases} \quad (18)$$

Plugging the third equation into the second one provides

$$\begin{cases} \mathbf{u}^* = \mathbf{u}^n - \Delta t (1-\theta)c \nabla p^n, \\ p^{n+1} - \Delta t^2 c^2 \theta^2 \Delta p^{n+1} = -c \Delta t \theta \nabla \cdot \mathbf{u}^* + p^n - \Delta t (1-\theta)c \nabla \cdot \mathbf{u}^n, \\ \mathbf{u}^{n+1} = \mathbf{u}^* - c \Delta t \theta \nabla p^{n+1}, \end{cases}$$

which is exactly the (Pressure algorithm).

**Schur complement on the velocity** The same guidelines can be followed to obtain a Schur complement on the velocity  $\mathbf{u}$  instead of the pressure in order to anticipate the Schur complement for shallow water equations.



Equation (13) turns into

$$\begin{cases} \begin{pmatrix} D_1 & U \\ L & \mathbb{D}_2 \end{pmatrix} \begin{pmatrix} p^{n+1} \\ \mathbf{u}^{n+1} \end{pmatrix} = \begin{pmatrix} R_p \\ \mathbf{R}_u \end{pmatrix}, & \text{in } \Omega, \\ \mathbf{u}^{n+1} \cdot \mathbf{n} \equiv 0 \text{ or } p^{n+1} \equiv 0 \text{ or } p^{n+1} - \mathbf{u}^{n+1} \cdot \mathbf{n} = p^0 - \mathbf{u}^0 \cdot \mathbf{n}, & \text{on } \partial\Omega, \end{cases}$$

where  $D_1$ ,  $\mathbb{D}_2$ ,  $L$ ,  $U$ ,  $\mathbf{R}_u$  and  $R_p$  are introduced in Equation (15). Applying Schur theory enables the following equation to be derived

$$\begin{pmatrix} I_1 & 0 \\ LD_1^{-1} & I_2 \end{pmatrix} \begin{pmatrix} D_1 & 0 \\ 0 & P_{u,\text{schur}} \end{pmatrix} \begin{pmatrix} I_1 & D_1^{-1}U \\ 0 & I_2 \end{pmatrix} \begin{pmatrix} p^{n+1} \\ \mathbf{u}^{n+1} \end{pmatrix} = \begin{pmatrix} \mathbf{R}_u \\ R_p \end{pmatrix}.$$

The operator

$$P_{u,\text{schur}} = \mathbb{D}_2 - LD_1^{-1}U = I_2 - c^2\theta^2 \Delta t^2 \nabla (\nabla \cdot I_2) \quad (19)$$

corresponds to the Schur complement on the velocity and leads to the following algorithm

$$\text{(Velocity algorithm)} \begin{cases} \text{Predictor : } D_1 p^* = R_p, & \text{in } \Omega, \\ \text{Propagation step : } P_{u,\text{schur}} \mathbf{u}^{n+1} = -Lp^* + \mathbf{R}_u, & \text{in } \Omega, \\ \text{Corrector : } D_1 p^{n+1} = D_1 p^* - U \mathbf{u}^{n+1}, & \text{in } \Omega. \end{cases}$$

As previously, for the Schur complement on the pressure, we need to write the boundary conditions for each operator. Using the first equation of (4) we obtain that  $\partial_t p = -c \nabla \cdot \mathbf{u}$ ,

- case (a) :  $\mathbf{u} \cdot \mathbf{n} \equiv 0$  on  $\partial\Omega$  : no boundary condition for  $D_1$ , homogeneous Dirichlet boundary condition  $(\mathbf{u} \cdot \mathbf{n}) = 0$  for  $P_{u,\text{schur}}$ ,
- case (b) :  $p \equiv 0$  on  $\partial\Omega$  : homogeneous Dirichlet for  $D_1$  and homogenous Neumann boundary condition  $\nabla \cdot \mathbf{u} = 0$  for  $P_{u,\text{schur}}$ ,
- case (c) :  $p - \mathbf{u} \cdot \mathbf{n} = g$  on  $\partial\Omega$  with  $g$  constant in time : Dirichlet boundary condition  $p^* = \mathbf{u}^n \cdot \mathbf{n} + g$  (prediction) and  $p^{n+1} = (\mathbf{u}^{n+1} \cdot \mathbf{n}) + g$  (correction) and Robin boundary condition  $c \Delta t (\nabla \cdot \mathbf{u}^{n+1}) + (\mathbf{u}^{n+1} \cdot \mathbf{n}) = (\mathbf{u}^n \cdot \mathbf{n})$  for  $P_{u,\text{schur}}$ .

**Remark 1.10.** Without any Schur complement, splitting the time discretization (13) with respect to

$$\begin{cases} p^* = p^n - (1 - \theta)c \Delta t \nabla \cdot \mathbf{u}^n, & \text{in } \Omega, \\ \mathbf{u}^{n+1} + \theta c \Delta t \nabla p^{n+1} = \mathbf{u}^n - (1 - \theta)c \Delta t \nabla p^n, & \text{in } \Omega, \\ p^{n+1} + \theta c \Delta t \nabla \cdot \mathbf{u}^{n+1} = p^*, & \text{in } \Omega, \end{cases}$$

and plugging the third equation in the second one brings the same velocity algorithm

$$\begin{cases} p^* = p^n - (1 - \theta)c \Delta t \nabla \cdot \mathbf{u}^n, & \text{in } \Omega, \\ \mathbf{u}^{n+1} - \theta^2 c^2 \Delta t^2 \nabla (\nabla \cdot \mathbf{u}^{n+1}) = -\theta c \Delta t \nabla p^* + \mathbf{u}^n - (1 - \theta)c \Delta t \nabla p^n, & \text{in } \Omega, \\ p^{n+1} = p^* - \theta c \Delta t \nabla \cdot \mathbf{u}^{n+1}, & \text{in } \Omega. \end{cases}$$

#### 1.4. Plane wave study for the different preconditioners

To retrieve from those systems of equations the underlying physics, preconditioning has to follow some properties : the preconditioned system should keep as most physical properties from the original problem as possible. For instance, for both systems to be equivalent at the spectral level, preconditioning should have the same propagation speeds as the full model. More precisely, the same method used for proposition 1.5 is performed here, in addition, the wave model and the preconditioners have the same propagation properties

Dispersion relations	Kernel for velocity-preconditioner	Kernel for pressure-preconditioner
$\omega = 0$	$\text{Ker}(A_{\mathbf{u}}) = \text{Span} \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ -k_2 \\ k_1 \end{pmatrix} \right\}$	$\text{Ker}(A_p) = \text{Span} \left\{ \begin{pmatrix} 0 \\ k_1 \\ k_2 \end{pmatrix}, \begin{pmatrix} 0 \\ -k_2 \\ k_1 \end{pmatrix} \right\}$
$\omega = c\ \mathbf{k}\ $	$\text{Ker}(A_{\mathbf{u}}) = \text{Span} \left\{ \begin{pmatrix} 1 \\ \frac{k_1}{\ \mathbf{k}\ } \\ \frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$	$\text{Ker}(A_p) = \text{Span} \left\{ \begin{pmatrix} 1 \\ \frac{k_1}{\ \mathbf{k}\ } \\ \frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$
$\omega = -c\ \mathbf{k}\ $	$\text{Ker}(A_{\mathbf{u}}) = \text{Span} \left\{ \begin{pmatrix} 1 \\ -\frac{k_1}{\ \mathbf{k}\ } \\ -\frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$	$\text{Ker}(A_p) = \text{Span} \left\{ \begin{pmatrix} 1 \\ -\frac{k_1}{\ \mathbf{k}\ } \\ -\frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$

TABLE 2. [Acoustic wave equations] Dispersion relations and Riemann invariants for the velocity- and the pressure-preconditioner

provided they have the same dispersion relations and the same associated kernel. We propose to compare the wave model to the two preconditioners in the homogeneous case. The preconditioners based on Schur complement on the pressure and Schur complement on the velocity are consistent with the two following models when the time step  $\Delta t$  is small

$$\begin{cases} \partial_t p + c \nabla \cdot \mathbf{u} = 0, \\ \partial_{tt} \mathbf{u} - c^2 \nabla (\nabla \cdot \mathbf{u}) = 0, \end{cases} \quad \text{and} \quad \begin{cases} \partial_{tt} p - c^2 \Delta p = 0, \\ \partial_t \mathbf{u} + c \nabla p = 0. \end{cases} \quad (20)$$

**Proposition 1.11.** *The models (20) associated with the preconditioners have exactly the same propagation properties as the wave operator. They both are spectrally equivalent.*

*Proof.* Once again, we apply the space-and time-Fourier transform to the two preconditioners. We obtain Equation (12) with

$$A_{\mathbf{u}} = \begin{pmatrix} i\omega & -ik_1 c & -ik_2 c \\ 0 & -\omega^2 + c^2 k_1^2 & c^2 k_1 k_2 \\ 0 & c^2 k_1 k_2 & -\omega^2 + c^2 k_2^2 \end{pmatrix}, \quad A_p = \begin{pmatrix} -\omega^2 + c^2(k_1^2 + k_2^2) & 0 & 0 \\ -ik_1 c & i\omega & 0 \\ -ik_2 c & 0 & i\omega \end{pmatrix}.$$

The spectrum of the matrices are given by

$$\sigma(A_{\mathbf{u}}) = \{i\omega, -\omega^2, -\omega^2 + c^2\|\mathbf{k}\|^2\}, \quad \sigma(A_p) = \{i\omega, i\omega, -\omega^2 + c^2\|\mathbf{k}\|^2\}$$

and so are the dispersion relations (same for the two models) and the corresponding kernels summarized in Table 2. The dispersion relations are exactly the same as those of the wave operator. The kernel associated with the velocity preconditioner has an additional constant pressure solution, which corresponds to the basis vector  $(1, 0, 0)^t$  in the eigenspace associated to  $\omega = 0$ . However, this additional solution is killed by the boundary conditions when we impose the pressure at the boundary for the equation  $\partial_t p + c \nabla \cdot \mathbf{u} = 0$ . Following the same line of thought, the kernel associated with the pressure preconditioner has an additional solution too which corresponds to the basis vector  $(0, k_1, k_2)^t$  in the eigenspace associated to  $\omega = 0$ . However, this additional solution is killed by the boundary conditions when we impose the normal velocity at the boundary in the equation  $\partial_t \mathbf{u} + c \nabla p = 0$ .

□

### 1.5. Numerical results

First of all, we propose to compare the classical GMRES method with and without preconditioning for the two physics-based preconditioners. Our test case consists in studying the following solution of the wave equations on  $\Omega = [0, 1]^2$

$$\begin{cases} p(t, x, y) = -2\sqrt{2}\pi \sin(2\sqrt{2}\pi c t) \cos(2\pi x) \cos(2\pi y), \\ u_1(t, x, y) = 2\pi \cos(2\sqrt{2}\pi c t) \sin(2\pi x) \cos(2\pi y), \\ u_2(t, x, y) = 2\pi \cos(2\sqrt{2}\pi c t) \cos(2\pi x) \sin(2\pi y). \end{cases}$$

For the comparison, we choose the following parameters : a tolerance of the GMRES method of  $10^{-9}$  and 5 time iterations to compute the implicit problem. The results are gathered in Table 3. For each time iteration (*i.e.* iteration to compute variables at  $t^{n+1}$  from variables at time  $t^n$ ), we compare the average number of sub-iterations (grey columns of Table 3). Other columns correspond to the average time for each time iteration. Those results are computed for high order polynomial method in space. Indeed, for small order polynomial method in space (and so, few Gauss-Lobatto points), the classical preconditioners of the GMRES method (e.g Jacobi) are sufficient. Contrariwise, the method developed here is suitable for large and not-sparse jacobian matrices, which result from high order polynomial methods. For this reason, we consider the 4-order in the following test (Table 3).

Firstly, we remark that the GMRES method is not able to solve this problem: an additional preconditioning

	Cells	$c \Delta t = 0.1$		$c \Delta t = 1$		$c \Delta t = 10$		$c \Delta t = 100$	
		Time	Sub-iter	Time	Sub-iter	Time	Sub-iter	Time	Sub-iter
GMRES	$24 \times 24$	-	nc	-	nc	-	nc	-	nc
	$32 \times 32$	-	nc	-	nc	-	nc	-	nc
	$48 \times 48$	-	nc	-	nc	-	nc	-	nc
GMRES+ Jacobi	$24 \times 24$	1.0E+0	27	1.2E+1	324	1.8E+2	4130	-	nc
	$32 \times 32$	3.7E+0	45	5.8E+1	530	3.8E+2	6070	-	nc
	$48 \times 48$	1.5E+1	38	3.8E+2	1100	-	nc	-	nc
GMRES+ pressure- PC	$24 \times 24$	4.6E+0	2	1.1E+1	4	1.9E+1	6	2.3E+1	7
	$32 \times 32$	1.3E+1	2	1.6E+1	3	3.1E+1	6	3.5E+1	8
	$48 \times 48$	2.6E+1	2	3.4E+1	3	5.3E+1	6	7.1E+1	8
GMRES+ velocity- PC	$24 \times 24$	9.5E+0	2	3.2 E+1	2	7.5E+1	2	1.4 E+2	3
	$32 \times 32$	2.5E+1	2	5.6E+1	2	2.4E+2	2	7.0E+2	3
	$48 \times 48$	1.2E+2	2	1.5E+2	2	7.7E+2	2	-	nc

TABLE 3. Average sub-iterations number to converge, in one time iteration, and average time for one time iteration. When the method fails to converge, we note "nc" in the iteration column.

is clearly needed. The bad-conditionin can be explained by the high order discretization and the hyperbolic structure. Secondly, for the most complex problems (which correspond to large  $c \Delta t$  and fine spatial grids), the physics-based preconditioners are able to treat the test with a better efficiency than the classical Jacobi preconditioning.

Let us now compare both of the two physics-based preconditioners (PC). We observe that both physics-based PC converge quickly (the velocity-PC begin a little bit quicker than the pressure-PC) but the pressure-PC is more efficient in time than the velocity-PC. There are two reasons to explain this difference :

- For small  $c \Delta t$ , the prediction and correction matrices are diagonal and so directly inverted : only the matrix associated to the Schur complement needs to be inverted. The size of this matrix is more important for the velocity-PC than for the pressure-PC, which explains the additional cost of the velocity-PC. However, for shallow water equations (Section 2), the prediction-correction matrices are

both advection matrices, and consequently, the cost between the two methods will be the same in this case.

- For large  $c \Delta t$ , the big additional cost and the non-convergence in one test for the velocity-PC come from the conditioning of the Schur complement. For the pressure-PC, the Schur complement writes  $P_{p,\text{schur}} = I_1 - c^2 \Delta t^2 \Delta I_1$  and this operator is coercive. When  $c \Delta t \gg 1$ , the limit operator is also coercive and well-conditioned. For the velocity-PC, the Schur complement is  $P_{u,\text{schur}} = I_2 - c^2 \Delta t^2 \nabla(\nabla \cdot I_2)$  and this operator is also coercive. But when  $c \Delta t \gg 1$ , the limit operator is not coercive and it is ill-conditioned. Indeed, the kernel of  $-\nabla(\nabla \cdot I_2)$  contains all the curl-free vector fields and the plane wave analysis (Subsection 1.4) shows that the limit operator is a multi-scale operator with two propagation velocities  $0 \ll |c \Delta t|$ , which leads to an ill-conditioning. In short, the additional cost or non-convergence come from the large time necessary to invert the Schur complement inside the velocity-PC. Nevertheless, this problem can be solved by finding a good algebraic PC for this operator (multi-grid methods, approximation based on the scale separation etc.)

**Remark 1.12. Remarks about the implementation and the method**

- The method depends on the compatibility of the boundary conditions between the Schur complement and the model. It is necessary to write correctly the boundary conditions on the sub model. The implementation is also important : if the boundary conditions are imposed weakly, the method is less efficient.
- The high-order Lagrangian polynomial matrices with Gauss-Lobatto points are not that ill-conditioned. The method has already been extensively tested, using a different code and discretization. It shows that the results with classical preconditioners like Jacobi are still worth.
- The implementation is clearly not optimal (Matrix storage, linkage with the libraries "Paralution" used etc) and after optimizing the implementation and finding the good solver for each sub-step, better results can be obtained.

Secondly, we propose to compare explicit and implicit methods for long time computation. We consider the same test case with the final time  $T_f = 10$  and the results are detailed in Table 4. The sound-speed  $c$  is set to 1, the mesh to  $32 \times 32$ , the polynomials are degree three and the explicit scheme is time-order two.

Scheme / results	$\Delta t$	nb time iter	time
Explicit	0.0005	10000	740 sec
Implicit I	1	10	90 sec
Implicit II	2	5	55 sec

TABLE 4. We use GMRES+velocity-PC for the implicit method. The time in the third column is the total time of the simulation. The explicit time step is maximal to have stability.

## 2. SHALLOW WATER EQUATIONS

The second studied system (5) is a simplification of shallow water+Exner system (1) which describes the movement of the flow between a free surface of an incompressible fluid and the topography of the ground, in an area where vertical dimension is neglected with respect to the horizontal scale. Equations (5) becomes in non-conservation form

$$\begin{cases} \partial_t h + \nabla \cdot (h \mathbf{u}) = 0, & (21a) \\ h \partial_t \mathbf{u} + h (\mathbf{u} \cdot \nabla) \mathbf{u} + \nabla \left( \frac{gh^2}{2} \right) = 0. & (21b) \end{cases}$$

In this section, we adapt for shallow water equations (21), the method developed previously.

## 2.1. Standard equations

**Boundary conditions** We propose to add some admissible boundary conditions to shallow water model (21). We consider the simple case  $\mathbf{u} \cdot \mathbf{n} \equiv \mathbf{0}$  on  $\partial\Omega$  with  $\mathbf{n}$  the unit normal vector at the boundary.

**Dissipation of the total energy** Shallow water equations (5) (without topography source term) together with those homogeneous boundary conditions conserve energy, as epitomized in the following proposition.

**Proposition 2.1.** *The total energy, defined by  $\mathcal{E} = \int_{\Omega} \left[ h \frac{\|\mathbf{u}\|^2}{2} + p \right] d\mathbf{x}$ , is conserved by homogeneous boundary conditions.*

*Proof.* The proof of D.D.Schnack in [17] about Ideal MHD is adapted here to non-conservative form of shallow water equations (21). Applying a dot product by  $\mathbf{u}$  to Equation (21b), multiplying Equation (21a) by  $\frac{\|\mathbf{u}\|^2}{2}$ , we obtain

$$\begin{aligned} \partial_t \left( h \frac{\|\mathbf{u}\|^2}{2} \right) + h \mathbf{u} \cdot (\mathbf{u} \cdot \nabla) \mathbf{u} + \mathbf{u} \cdot \nabla p + \frac{\|\mathbf{u}\|^2}{2} \nabla \cdot (h \mathbf{u}) &= 0, \\ \partial_t \left( h \frac{\|\mathbf{u}\|^2}{2} \right) + \nabla \cdot \left( h \mathbf{u} \frac{\|\mathbf{u}\|^2}{2} \right) + \mathbf{u} \cdot \nabla p &= 0. \end{aligned}$$

Afterwards, let us deal with the term  $\mathbf{u} \cdot \nabla p$ . We multiply Equation (21a) by  $gh$  and noticing that the pressure is  $p = \frac{gh^2}{2}$  appear

$$\partial_t \left( \frac{gh^2}{2} \right) + \mathbf{u} \cdot \nabla \left( \frac{gh^2}{2} \right) + gh^2 \nabla \cdot \mathbf{u} = 0,$$

thus

$$\partial_t p + \mathbf{u} \cdot \nabla p + 2p \nabla \cdot \mathbf{u} = 0.$$

Adding this equation to the equation on the kinetic energy, we obtain

$$\partial_t \left( h \frac{\|\mathbf{u}\|^2}{2} + p \right) + \nabla \cdot \left( h \mathbf{u} \frac{\|\mathbf{u}\|^2}{2} \right) + 2 \nabla \cdot (p \mathbf{u}) = 0.$$

Now we integrate and apply the flux-divergence theorem to obtain

$$\partial_t \int_{\Omega} \left( h \frac{\|\mathbf{u}\|^2}{2} + p \right) d\mathbf{x} = - \int_{\partial\Omega} \left( h \frac{\|\mathbf{u}\|^2}{2} + 2p \right) (\mathbf{u} \cdot \mathbf{n}) d\boldsymbol{\sigma} = - \int_{\partial\Omega} \left( \sqrt{\frac{p}{2g}} \|\mathbf{u}\|^2 + 2p \right) (\mathbf{u} \cdot \mathbf{n}) d\boldsymbol{\sigma}.$$

The considered boundary conditions allow to obtain the energy conservation.  $\square$

## 2.2. Riemann invariants and propagation

We perform the same weft as the one followed in Section 1 and compute the wave velocities of the original problem (5) to compare it with wave velocities of preconditioner in Subsections 2.5 and 2.6. We first linearize non conservative shallow water equations (21) around constant state  $(h^n, \mathbf{u}^n)$  to obtain

$$\begin{cases} \partial_t \delta h + h^n \nabla \cdot \delta \mathbf{u} + \mathbf{u}^n \cdot \nabla \delta h = 0, \\ h^n \partial_t \delta \mathbf{u} + h^n (\mathbf{u}^n \cdot \nabla) \delta \mathbf{u} + gh^n \nabla \delta h = 0. \end{cases} \quad (22)$$

**Proposition 2.2.** *The Riemann invariants for linearized shallow water equations are  $(0, -k_2, k_1)^t$ ,  $(\sqrt{gh^n}, g \frac{k_1}{\|\mathbf{k}\|}, g \frac{k_2}{\|\mathbf{k}\|})^t$  and  $(\sqrt{gh^n}, -g \frac{k_1}{\|\mathbf{k}\|}, -g \frac{k_2}{\|\mathbf{k}\|})^t$ , traveling with respectively the velocity  $\mathbf{k} \cdot \mathbf{u}^n$ , the velocity  $\mathbf{k} \cdot \mathbf{u}^n + c \|\mathbf{k}\|$  and the velocity  $\mathbf{k} \cdot \mathbf{u}^n - c \|\mathbf{k}\|$  with the gravity waves speed  $c = \sqrt{gh^n}$ .*

Dispersion relations (canceling the eigenvalues)	Riemann invariants (basis vectors of the eigenspaces)
$\omega = \mathbf{k} \cdot \mathbf{u}^n$	$\text{Ker}(A_{\text{sh}}) = \text{Span} \left\{ \begin{pmatrix} 0 \\ -k_2 \\ k_1 \end{pmatrix} \right\}$
$\omega = \mathbf{k} \cdot \mathbf{u}^n + \sqrt{gh^n} \ \mathbf{k}\ $	$\text{Ker}(A_{\text{sh}}) = \text{Span} \left\{ \begin{pmatrix} \sqrt{gh^n} \\ g \frac{k_1}{\ \mathbf{k}\ } \\ g \frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$
$\omega = \mathbf{k} \cdot \mathbf{u}^n - \sqrt{gh^n} \ \mathbf{k}\ $	$\text{Ker}(A_{\text{sh}}) = \text{Span} \left\{ \begin{pmatrix} \sqrt{gh^n} \\ -g \frac{k_1}{\ \mathbf{k}\ } \\ -g \frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$

TABLE 5. [Shallow water equations] Dispersion relations and Riemann invariants for the shallow water model

**Remark 2.3.** The sound velocity  $c$  for the acoustic wave equations is replaced by the surface wave velocity  $\sqrt{gh^n}$ .

*Proof.* Applying the space-and time-Fourier transform to the linearized shallow water equations (22), we obtain the equation (12) with  $\mathbf{V}^0 = (\mathcal{F}(\delta h), \mathcal{F}(\delta u_1), \mathcal{F}(\delta u_2))^t$  and

$$A_{\text{sh}} = \begin{pmatrix} i\omega - i\mathbf{k} \cdot \mathbf{u}^n & -ih^n k_1 & -ih^n k_2 \\ -igh^n k_1 & ih^n(\omega - \mathbf{k} \cdot \mathbf{u}^n) & 0 \\ -igh^n k_2 & 0 & ih^n(\omega - \mathbf{k} \cdot \mathbf{u}^n) \end{pmatrix}.$$

After calculating the eigenvalues and the corresponding eigenvectors, we summarize the results in Table 5.  $\square$

### 2.3. Discretization and linearization

**Time discretization** The same notation as in Subsection 1.2 are kept :  $\Delta t$  for time step,  $\Delta x$  for size mesh,  $\mathbf{u}^n$  and  $h^n$  always stand for velocity and height at time  $t^n = n \Delta t$ . As explained in Introduction, time multi-scales cohabit for shallow water equations (5) and create several difficulties. For instance, the simulation time  $T$ , given by the sedimentation behavior, must be much greater than the time step  $\Delta t$ . Thus, to take into account problems due to those time multi-scales, the following semi-discrete  $\theta$ -scheme is computed with the boundary conditions and  $\theta \in [0, 1]$

$$\begin{cases} \frac{h^{n+1} - h^n}{\Delta t} + \theta \nabla \cdot (h^{n+1} \mathbf{u}^{n+1}) + (1 - \theta) \nabla \cdot (h^n \mathbf{u}^n) = 0, & \text{in } \Omega, \\ h^n \frac{\mathbf{u}^{n+1} - \mathbf{u}^n}{\Delta t} + \theta h^{n+1} (\mathbf{u}^{n+1} \cdot \nabla) \mathbf{u}^{n+1} + (1 - \theta) h^n (\mathbf{u}^n \cdot \nabla) \mathbf{u}^n + \theta gh^{n+1} \nabla h^{n+1} + (1 - \theta) gh^n \nabla h^n = 0, & \text{in } \Omega, \\ \mathbf{u}^{n+1} \cdot \mathbf{n} \equiv 0, & \text{on } \partial\Omega, \end{cases}$$

which can be rewritten in the form

$$\begin{cases} G \begin{pmatrix} h^{n+1} \\ \mathbf{u}^{n+1} \end{pmatrix} = E \begin{pmatrix} h^n \\ \mathbf{u}^n \end{pmatrix}, & \text{in } \Omega, \\ \mathbb{B} \begin{pmatrix} h^{n+1} \\ \mathbf{u}^{n+1} \end{pmatrix} = 0 & \text{on } \partial\Omega. \end{cases}$$

with

$$G : \begin{pmatrix} h \\ \mathbf{u} \end{pmatrix} \mapsto \begin{pmatrix} h + \theta \Delta t \nabla \cdot (h \mathbf{u}) \\ h^n \mathbf{u} + \theta \Delta t h (\mathbf{u} \cdot \nabla) \mathbf{u} + \theta \Delta t h g \nabla h \end{pmatrix}$$

and

$$E : \begin{pmatrix} h \\ \mathbf{u} \end{pmatrix} \mapsto \begin{pmatrix} h - (1 - \theta) \Delta t \nabla \cdot (h \mathbf{u}) \\ h \mathbf{u} - (1 - \theta) \Delta t h (\mathbf{u} \cdot \nabla) \mathbf{u} - (1 - \theta) \Delta t h g \nabla h \end{pmatrix}.$$

In order to design a preconditioner, we would like to follow the same method as in Section 1 *i.e.* using the Schur theory. Hence, the studied system needs to be a matrix system for Schur complement to be computed. A linearization of equations is then necessary to introduce the associated linearized matrix system.

**Linearization** Neglecting the second order terms yields the linearized system

$$G \begin{pmatrix} h^n \\ \mathbf{u}^n \end{pmatrix} + J_{\text{ac}_G}^n \begin{pmatrix} \delta h^{n+1} \\ \delta \mathbf{u}^{n+1} \end{pmatrix} = E \begin{pmatrix} h^n \\ \mathbf{u}^n \end{pmatrix}$$

where  $J_{\text{ac}_G}^n$  is the Jacobian matrix of  $G$  at time  $n$  and  $\delta h^{n+1}$  (resp.  $\delta \mathbf{u}^{n+1}$ ) stands for the difference  $h^{n+1} - h^n$  (resp.  $\mathbf{u}^{n+1} - \mathbf{u}^n$ ). The Jacobian matrix is decomposed in four blocs

$$J_{\text{ac}_G}^n = \begin{pmatrix} D_1 & U \\ L & \mathbb{D}_2 \end{pmatrix}$$

with an advection term for  $h^n$

$$D_1 = I_1 + \theta \Delta t \nabla \cdot (\mathbf{u}^n I_1),$$

an advection-convection term for  $\mathbf{u}^n$

$$\mathbb{D}_2 = h^n I_2 + \theta \Delta t h^n (\mathbf{u}^n \cdot \nabla) I_2 + \theta \Delta t h^n \begin{pmatrix} \partial_x u_1^n & \partial_y u_1^n \\ \partial_x u_2^n & \partial_y u_2^n \end{pmatrix} I_2$$

and some coupling terms

$$U = \theta \Delta t \nabla \cdot (h^n I_2), \quad L = \theta g \Delta t \nabla (h^n I_1) + \theta \Delta t I_1 (\mathbf{u}^n \cdot \nabla) \mathbf{u}^n.$$

In the previous relations  $I_d$  stands for the  $d \times d$ -Identity matrix. Here, the operators  $U$  and  $L$  are a generalization of those defined in the acoustic wave equation (Equation (15)).

Hence, the linearized semi-discrete system becomes

$$\begin{pmatrix} D_1 & U \\ L & \mathbb{D}_2 \end{pmatrix} \begin{pmatrix} \delta h^{n+1} \\ \delta \mathbf{u}^{n+1} \end{pmatrix} = \begin{pmatrix} -\Delta t \nabla \cdot (h^n \mathbf{u}^n) \\ -\Delta t h^n (\mathbf{u}^n \cdot \nabla) \mathbf{u}^n - \Delta t g h^n \nabla h^n \end{pmatrix}. \quad (23)$$

#### 2.4. Preconditioning for the shallow water equations

As for the wave equations, we propose in this subsection to write the preconditioner based on the Schur complement on the velocity. We focalize our study on the case where the Schur complement is computed on the velocity because this is the most interesting case where we want to extend the method to more complicate fluid model (Euler equation, MHD). Indeed for more complex models, the coupling between all the equations is ensured by the velocity equation. The Schur complement allows to simplify some terms. When we apply the Schur complement to the velocity we obtain a "parabolization" of the coupling terms (in Chacon's terminology) which are the stiff terms in the equation. The algorithm which is deduced writes

$$\text{(Velocity preconditioner)} \begin{cases} \text{Predictor : } D_1 \delta h^* = -\Delta t \nabla \cdot (h^n \mathbf{u}^n), \text{ in } \Omega, \\ \text{Propagation : } P_{\mathbf{u}} \delta \mathbf{u}^{n+1} = -L \delta h^* - \Delta t h^n (\mathbf{u}^n \cdot \nabla) \mathbf{u}^n - \Delta t g h^n \nabla h^n, \text{ in } \Omega, \\ \text{Corrector : } D_1 \delta h^{n+1} = D_1 \delta h^* - U \delta \mathbf{u}^{n+1}, \text{ in } \Omega, \end{cases}$$

where  $D_1$  and  $\mathbb{D}_2$  are introduced in (23) and with

$$P_{\mathbf{u}} = \mathbb{D}_2 - LD_1^{-1}U,$$

the Schur complement on the velocity. We introduce the boundary conditions for each operator : no boundary condition for  $D_1$  and Dirichlet boundary condition  $\mathbf{u} \cdot \mathbf{n} = 0$  for  $P_{\mathbf{u}}$ .

**Remark 2.4.** As denoted for the acoustic wave equations, using a splitting operator instead of the Schur theory produces the same result. In deed, splitting Equations (23) in the following form

$$\begin{cases} D_1 \delta h^* = -\Delta t \nabla (h^n \mathbf{u}^n), \\ L \delta h^{n+1} + \mathbb{D}_2 \delta \mathbf{u}^{n+1} = -\Delta t h^n (\mathbf{u}^n \cdot \nabla) \mathbf{u}^n - \Delta t g h^n \nabla h^n, \\ D_1 \delta h^{n+1} - D_1 \delta h^* + U \delta \mathbf{u}^{n+1} = 0. \end{cases}$$

and injecting the third equation in the second one gives the same algorithm.

**Remark 2.5.** The Schur complement  $P_{\mathbf{u}}$  requests the inversion of the advection term  $D_1$ . Some approximations have to be found to compute easily this invert while keeping physics-based characteristics such as wave propagation speeds close to those of the initial operator. In the two following subsections, we suggest two approximations to compute  $P_{\mathbf{u}}$ , adapted from magnetohydrodynamics in [7] : the first approximation is especially for the low Froude number whereas the second one is more general.

The following hypothesis is assumed throughout the two next subsections

$$(H_1) \quad (h^n, \mathbf{u}^n) \text{ are constants.}$$

**Remark 2.6.** The hypothesis  $(H_1)$  is a classical hypothesis used frequently in studies of physic waves. In deed, the overall model would be too difficult to analyse. However, our preconditioners can operate without  $(H_1)$  and that restriction is only needed in Fourier analysis (Subsections 2.5 and 2.6).

## 2.5. Schur complement for small flow approximation (low Froude number)

In this section, we propose to design a Schur complement in the low Froude number case by adapting the results of [7]. The small flow regime corresponds to the low Froude number  $F_r \ll 1$  where the Froude number corresponds to the Mach number (ratio between acoustic and advection phenomena) and is equal to

$$F_r = \frac{\mathbf{u}^n \cdot \mathbf{k}}{\sqrt{g h^n} \|\mathbf{k}\|}.$$

Assuming the flow is small means that the characteristic speed of the fluid is smaller than the gravity waves' speed. Moreover, we impose the time step to be of the same order than the gravity waves' speed, which leads to (with a rescaling with respect to gravity waves)  $\Delta t \sim \sqrt{g h^n} = O(1)$  and to the approximation

$$(H_2) \quad D_1 \approx I_1.$$

**Remark 2.7.** The aim of our study is to design a preconditioner for any Froude number. We begin with the small flow approximation because of the simplicity of the Schur inversion in this case. However, we do not apply all simplifications linked to this approximation : only those indispensable for the Schur inversion are made in order to anticipate the study of any Froude number case.

The previous hypothesis  $(H_2)$  makes the operator  $LD_1^{-1}U$  which appears in  $P_{\mathbf{u}}$  to be easier to compute. Thanks to the relation

$$L \delta h = (\theta \Delta t (\mathbf{u}^n \cdot \nabla) \mathbf{u}^n) \delta h + \theta \Delta t g \nabla (h^n \delta h),$$



and the choice  $\delta h = U\delta\mathbf{u} = \theta \Delta t \nabla \cdot (h^n \delta\mathbf{u})$ , this operator becomes

$$LU(\delta\mathbf{u}) = (\theta \Delta t (\mathbf{u}^n \cdot \nabla) \mathbf{u}^n) (\theta \Delta t \nabla \cdot (h^n \delta\mathbf{u})) + \theta^2 \Delta t^2 g \nabla (h^n \nabla \cdot (h^n \delta\mathbf{u})). \quad (24)$$

Eventually,  $P_{\mathbf{u}}$  is approximated by

$$P_{\mathbf{u}}^{\text{low}} \delta\mathbf{u} = h^n \delta\mathbf{u} + \theta \Delta t h^n (\mathbf{u}^n \cdot \nabla) \delta\mathbf{u} + \theta \Delta t h^n (\delta\mathbf{u} \cdot \nabla) \mathbf{u}^n - \underbrace{(\theta \Delta t (\mathbf{u}^n \cdot \nabla) \mathbf{u}^n) (\theta \Delta t \nabla \cdot (h^n \delta\mathbf{u}))}_{\text{advection term (A)}} - \theta^2 \Delta t^2 g \nabla (h^n \nabla \cdot (h^n \delta\mathbf{u})). \quad (25)$$

**Proposition 2.8.** *Assume that  $\mathbf{u}^n$  is constant (hypothesis  $(H_1)$ ), then the low flow approximate Schur complement has propagation properties very close to those of the initial shallow water operator. For a low Froude number, they both are spectrally equivalent.*

*Proof.* When the time step tends to zero, the low Froude number preconditioner (25) is consistent with the following model

$$\begin{cases} \partial_t \delta h + h^n \nabla \cdot \delta\mathbf{u} + \mathbf{u}^n \cdot \nabla \delta h = 0, \\ h^n \partial_{tt} \delta\mathbf{u}^n + h^n (\mathbf{u}^n \cdot \nabla) \partial_t \delta\mathbf{u} - g (h^n)^2 \nabla (\nabla \cdot \delta\mathbf{u}) = 0. \end{cases}$$

Thanks to space-and time-Fourier transform, the matrix of Equation (12) is

$$A_{\mathbf{u}}^{\text{low}}(\mathbf{k}, \omega) = \begin{pmatrix} i(\omega - \mathbf{k} \cdot \mathbf{u}^n) & -ih^n k_1 & -ih^n k_2 \\ 0 & -h^n \omega (\omega - \mathbf{u}^n \cdot \mathbf{k}) + g (h^n)^2 k_1^2 & g (h^n)^2 k_1 k_2 \\ 0 & g (h^n)^2 k_1 k_2 & -h^n \omega (\omega - \mathbf{k} \cdot \mathbf{u}^n) + g (h^n)^2 k_2^2 \end{pmatrix}.$$

Imposing one of the eigenvalues to be equal to zero and thus the corresponding eigenspace to be equal to the kernel enables us to compute the propagation properties (cf Table 6, where the dispersion relations and the Riemann invariants are computed with respect to the Froude number  $F_r$ .)  $\square$

To conclude, the propagation speeds for the models corresponding to the low Froude number preconditioner are close to the propagation speeds of the shallow water if and only if the Froude number is small. Moreover, the kernel associated with the dispersion relations are close to the results given by the shallow water model only in the small Froude regime. Therefore, as expected, the models corresponding to the preconditioners have the same propagation properties as the full model only in the small Froude regime.

**Remark 2.9.** For a Froude number asymptotically small ( $\mathbf{u} = \mathbf{0}$ ), we recognise the acoustic wave equation with  $c = \sqrt{gh}$ .

## 2.6. Schur complement for arbitrary flow approximation (any Froude number)

In this subsection, we compute an approximation of the Schur complement  $P_{\mathbf{u}}$  proposed in [8] valid for an arbitrary Froude number *i.e.* without any hypothesis on the flow. Since hypothesis  $(H_2)$  is not required, an other method to compute easily the invert of  $D_1$  has to be found.

**Computation of the Schur complement for any flow** The idea is to construct an operator  $\mathcal{M}$  such that  $U\mathcal{M} \approx D_1 U$ , to obtain accordingly the following approximation

$$P_{\mathbf{u}}^{\text{any flow}} = (\mathbb{D}_2 \mathcal{M} - LU) \mathcal{M}^{-1}.$$

Therefore,  $\delta\mathbf{u}$  is in the kernel of  $P_{\mathbf{u}}^{\text{any flow}}$  if and only if there exists  $\delta\mathbf{u}^*$  such that

$$\begin{cases} (\mathbb{D}_2 \mathcal{M} - LU) \delta\mathbf{u}^* = 0, \\ \delta\mathbf{u} = \mathcal{M} \delta\mathbf{u}^*. \end{cases}$$

Dispersion relations	Riemann invariants
$\omega = \mathbf{u}^n \cdot \mathbf{k}$	$\text{Ker}(A_{\mathbf{u}}^{\text{low}}) = \text{Span} \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ -k_2 \\ k_1 \end{pmatrix} \right\}$
$\omega = \frac{\mathbf{u}^n \cdot \mathbf{k}}{2} + \sqrt{\frac{(\mathbf{u}^n \cdot \mathbf{k})^2}{4} + gh^n \ \mathbf{k}\ ^2}$ $= \mathbf{u}^n \cdot \mathbf{k} + \sqrt{gh^n \ \mathbf{k}\ } \left( 1 + \sqrt{\frac{F_r^2}{4} + 1} - \left( 1 + \frac{F_r}{2} \right) \right)$	$\text{Ker}(A_{\mathbf{u}}^{\text{low}}) = \text{Span} \left\{ \begin{pmatrix} \frac{\sqrt{gh^n}}{\sqrt{gh^n}} \\ \frac{\sqrt{g} \left( -\frac{\mathbf{u}^n \cdot \mathbf{k}}{2} + \sqrt{\frac{(\mathbf{u}^n \cdot \mathbf{k})^2}{4} + gh^n \ \mathbf{k}\ ^2} \right)}{\sqrt{h^n}} \frac{k_1}{\ \mathbf{k}\ ^2} \\ \frac{\sqrt{g} \left( -\frac{\mathbf{u}^n \cdot \mathbf{k}}{2} + \sqrt{\frac{(\mathbf{u}^n \cdot \mathbf{k})^2}{4} + gh^n \ \mathbf{k}\ ^2} \right)}{\sqrt{h^n}} \frac{k_2}{\ \mathbf{k}\ ^2} \end{pmatrix} \right\}$ $= \text{Span} \left\{ \begin{pmatrix} \frac{\sqrt{gh^n}}{g \left( \sqrt{\frac{F_r^2}{4} + 1} - \frac{F_r}{2} \right)} \frac{k_1}{\ \mathbf{k}\ } \\ g \left( \sqrt{\frac{F_r^2}{4} + 1} - \frac{F_r}{2} \right) \frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$
$\omega = \frac{\mathbf{u}^n \cdot \mathbf{k}}{2} - \sqrt{\frac{(\mathbf{u}^n \cdot \mathbf{k})^2}{4} + gh^n \ \mathbf{k}\ ^2}$ $= \mathbf{u}^n \cdot \mathbf{k} - \sqrt{gh^n \ \mathbf{k}\ } \left( 1 + \sqrt{\frac{F_r^2}{4} + 1} - \left( 1 - \frac{F_r}{2} \right) \right)$	$\text{Ker}(A_{\mathbf{u}}^{\text{low}}) = \text{Span} \left\{ \begin{pmatrix} \frac{\sqrt{gh^n}}{-\sqrt{g} \left( \frac{\mathbf{u}^n \cdot \mathbf{k}}{2} + \sqrt{\frac{(\mathbf{u}^n \cdot \mathbf{k})^2}{4} + gh^n \ \mathbf{k}\ ^2} \right)} \frac{k_1}{\ \mathbf{k}\ ^2} \\ \frac{\sqrt{gh^n}}{-\sqrt{g} \left( \frac{\mathbf{u}^n \cdot \mathbf{k}}{2} + \sqrt{\frac{(\mathbf{u}^n \cdot \mathbf{k})^2}{4} + gh^n \ \mathbf{k}\ ^2} \right)} \frac{k_2}{\ \mathbf{k}\ ^2} \end{pmatrix} \right\}$ $= \text{Span} \left\{ \begin{pmatrix} \frac{\sqrt{gh^n}}{-g \left( \sqrt{\frac{F_r^2}{4} + 1} + \frac{F_r}{2} \right)} \frac{k_1}{\ \mathbf{k}\ } \\ -g \left( \sqrt{\frac{F_r^2}{4} + 1} + \frac{F_r}{2} \right) \frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$

TABLE 6. [Shallow water equations] Dispersion relations and Riemann invariants for low Froude number velocity-preconditioner without advection terms

The previous system prevents us from computing the invert of  $\mathcal{M}$ . Let us define

$$\mathcal{M} = I_2 + \theta \Delta t \frac{\mathbf{u}^n}{h^n} (\nabla \cdot (h^n I_2)), \quad (26)$$

to obtain  $D_1 U \delta \mathbf{u} = U \mathcal{M} \delta \mathbf{u}$ . Definition (26) of the operator  $\mathcal{M}$  leads to

$$\begin{aligned} \mathbb{D}_2 \mathcal{M} \delta \mathbf{u} &= h^n \delta \mathbf{u} + \theta \Delta t [h^n \mathbf{u}^n \cdot \nabla \delta \mathbf{u} + h^n (\delta \mathbf{u} \cdot \nabla \mathbf{u}^n) + \mathbf{u}^n \nabla \cdot (h^n \delta \mathbf{u})] \\ &\quad + \theta^2 \Delta t^2 \left[ (\mathbf{u}^n \cdot \nabla \mathbf{u}^n) \nabla \cdot (h^n \delta \mathbf{u}) + h^n \mathbf{u}^n \cdot \nabla \left( \frac{\mathbf{u}^n}{h^n} \nabla \cdot (h^n \delta \mathbf{u}) \right) \right] \end{aligned}$$

The approximate Schur complement  $P_{\mathbf{u}}^{\text{any flow}}$  is given by the operator  $\mathbb{D}_2 \mathcal{M}$ , the  $LU$  low Froude operator (24)

:

$$\begin{cases} P_{\mathbf{u}}^{\text{any flow}} \delta \mathbf{u}^* = h^n \delta \mathbf{u}^* + \theta \Delta t [h^n \mathbf{u}^n \cdot \nabla \delta \mathbf{u}^* + h^n (\delta \mathbf{u}^* \cdot \nabla \mathbf{u}^n) + \mathbf{u}^n \nabla \cdot (h^n \delta \mathbf{u}^*)] \\ \quad + \theta^2 \Delta t^2 \left[ h^n \mathbf{u}^n \cdot \nabla \left( \frac{\mathbf{u}^n}{h^n} \nabla \cdot (h^n \delta \mathbf{u}^*) \right) - g \nabla (h^n \nabla \cdot (h^n \delta \mathbf{u}^*)) \right], \\ \delta \mathbf{u} = \delta \mathbf{u}^* + \theta \Delta t \frac{\mathbf{u}^n}{h^n} (\nabla \cdot (h^n \delta \mathbf{u}^*)). \end{cases}$$

### Study of the wave propagation associated to $P_{\mathbf{u}}^{\text{any flow}}$

**Proposition 2.10.** *Assuming  $\mathbf{u}^n$  and  $h^n$  constant (Hypothesis  $(H_1)$ ) then the Schur operator and the shallow water model have the same wave propagation properties.*

*Proof.* The arbitrary Froude number preconditioner when the time step tends to zero is consistent with the following model

$$\begin{cases} \partial_t(\delta h) + h^n \nabla \cdot \delta \mathbf{u} + \mathbf{u}^n \cdot \nabla \delta h = 0, \\ h^n \partial_{tt} \delta \mathbf{u}^* + h^n (\nabla \cdot \partial_t \delta \mathbf{u}^*) \mathbf{u}^n + h^n (\mathbf{u}^n \cdot \nabla) \partial_t \delta \mathbf{u}^* + h^n \mathbf{u}^n (\mathbf{u}^n \cdot \nabla) (\nabla \cdot \delta \mathbf{u}^*) - g (h^n)^2 \nabla \nabla \cdot \delta \mathbf{u}^* = 0, \end{cases} \quad (27)$$

with the following relation, which results from the operator  $\mathcal{M}$

$$\partial_{tt} \delta \mathbf{u} = \partial_{tt} \delta \mathbf{u}^* + \mathbf{u}^n \nabla \cdot \partial_t \delta \mathbf{u}^*.$$

System (27) leads to Equation (12) with the matrix

$$A_{\delta \mathbf{u}^*} = \begin{pmatrix} i(\omega - \mathbf{k} \cdot \mathbf{u}^n) & -ik_1 h^n & -ik_2 h^n \\ 0 & h^n (u_1^n k_1 - \omega) (\omega - \mathbf{k} \cdot \mathbf{u}^n) + g (h^n)^2 k_1^2 & h^n u_1^n k_2 (\omega - \mathbf{k} \cdot \mathbf{u}^n) + g (h^n)^2 k_1 k_2 \\ 0 & h^n u_2^n k_1 (\omega - \mathbf{k} \cdot \mathbf{u}^n) + g (h^n)^2 k_1 k_2 & h^n (u_2^n k_2 - \omega) (\omega - \mathbf{k} \cdot \mathbf{u}^n) + g (h^n)^2 k_2^2 \end{pmatrix},$$

and the operator  $\mathcal{M}$  becomes due to the Fourier transform,  $-\omega^2 \mathcal{F}(\delta \mathbf{u}) = \left[ -\omega^2 \mathbb{I}_2 + \omega \begin{pmatrix} u_1^n k_1 & u_1^n k_2 \\ u_2^n k_1 & u_2^n k_2 \end{pmatrix} \right] \mathcal{F}(\delta \mathbf{u}^*)$ ,

thus  $\mathcal{F}(\delta \mathbf{u}^*) = \begin{pmatrix} \frac{\omega - k_2 u_2^n}{\omega - \mathbf{k} \cdot \mathbf{u}^n} & \frac{k_2 u_1^n}{\omega - \mathbf{k} \cdot \mathbf{u}^n} \\ \frac{k_1 u_2^n}{\omega - \mathbf{k} \cdot \mathbf{u}^n} & \frac{\omega - k_1 u_1^n}{\omega - \mathbf{k} \cdot \mathbf{u}^n} \end{pmatrix} \mathcal{F}(\delta \mathbf{u})$ . Hence, it becomes with the three unknowns  $(\mathcal{F}(\delta h), \mathcal{F}(\delta u_1), \mathcal{F}(\delta u_2))^t$

$$\begin{pmatrix} \mathcal{F}(\delta h) \\ \mathcal{F}(\delta u_1^*) \\ \mathcal{F}(\delta u_2^*) \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 0 & \frac{\omega - k_2 u_2^n}{\omega - \mathbf{k} \cdot \mathbf{u}^n} & \frac{k_2 u_1^n}{\omega - \mathbf{k} \cdot \mathbf{u}^n} \\ 0 & \frac{k_1 u_2^n}{\omega - \mathbf{k} \cdot \mathbf{u}^n} & \frac{\omega - k_1 u_1^n}{\omega - \mathbf{k} \cdot \mathbf{u}^n} \end{pmatrix}}_{=A_{\mathcal{M}}} \begin{pmatrix} \mathcal{F}(\delta h) \\ \mathcal{F}(\delta u_1) \\ \mathcal{F}(\delta u_2) \end{pmatrix}.$$

Then, Equation (12) for the unknowns  $\delta \mathbf{u}$  writes with the matrix

$$A_{\mathbf{u}}^{\text{any flow}} = A_{\delta \mathbf{u}^*} A_{\mathcal{M}} = \begin{pmatrix} i(\omega - \mathbf{k} \cdot \mathbf{u}^n) & -\frac{ik_1 h^n \omega}{\omega - \mathbf{k} \cdot \mathbf{u}^n} & -\frac{ik_2 h^n \omega}{\omega - \mathbf{k} \cdot \mathbf{u}^n} \\ 0 & -h^n \omega (\omega - \mathbf{u}^n \cdot \mathbf{k}) + \frac{g (h^n)^2 k_1^2 \omega}{\omega - \mathbf{u}^n \cdot \mathbf{k}} & \frac{g (h^n)^2 k_1 k_2 \omega}{\omega - \mathbf{u}^n \cdot \mathbf{k}} \\ 0 & \frac{g (h^n)^2 k_1 k_2 \omega}{\omega - \mathbf{u}^n \cdot \mathbf{k}} & -h^n \omega (\omega - \mathbf{k} \cdot \mathbf{u}^n) + \frac{g (h^n)^2 k_2^2 \omega}{\omega - \mathbf{u}^n \cdot \mathbf{k}} \end{pmatrix}$$

and the associated eigenvalues and eigenspaces are synthesized in Table 7, which concludes the proof.  $\square$

## 2.7. Numerical results

In this subsection, hypothesis  $(H_1)$  is no longer applied. As previously for the linear case (Subsection 1.5), our first numerical test is a comparison for different time steps between our method and the classical GMRES method with and without Jacobi preconditioner. In addition, we propose to study the effect of the Froude number on the results. We consider the following solution on  $\Omega = [0, 1]^2$

$$\begin{cases} h(t, x, y) = 1.0 + x^2(1 - x^2)y^2(1 - y^2), \\ u_1(t, x, y) = \alpha x(1 - x)(1 - 2y) + \epsilon \sin(2\pi x) \sin(2\pi y), \\ u_2(t, x, y) = -\alpha y(1 - y)(1 - 2x) + \epsilon \sin(2\pi x) \sin(2\pi y), \end{cases}$$

Dispersion relations	Riemann invariants
$\omega = \mathbf{u}^n \cdot \mathbf{k}$	$\text{Ker}(A_{\mathbf{u}}^{\text{any flow}}) = \text{Span} \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ -k_2 \\ k_1 \end{pmatrix} \right\}$
$\omega = \mathbf{u}^n \cdot \mathbf{k} + \sqrt{gh^n} \ \mathbf{k}\ $	$\text{Ker}(A_{\mathbf{u}}^{\text{any flow}}) = \text{Span} \left\{ \begin{pmatrix} \frac{\mathbf{k} \cdot \mathbf{u}^n}{\ \mathbf{k}\ } + \sqrt{gh^n} \\ g \frac{k_1}{\ \mathbf{k}\ } \\ g \frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$
$\omega = \mathbf{u}^n \cdot \mathbf{k} - \sqrt{gh^n} \ \mathbf{k}\ $	$\text{Ker}(A_{\mathbf{u}}^{\text{any flow}}) = \text{Span} \left\{ \begin{pmatrix} -\frac{\mathbf{k} \cdot \mathbf{u}^n}{\ \mathbf{k}\ } + \sqrt{gh^n} \\ -g \frac{k_1}{\ \mathbf{k}\ } \\ -g \frac{k_2}{\ \mathbf{k}\ } \end{pmatrix} \right\}$

TABLE 7. [Shallow water equations] Dispersion relations and Riemann invariants for arbitrary Froude number velocity-preconditioner

which satisfies  $\mathbf{u} \cdot \mathbf{n} = 0$  at the boundary and corresponds to a steady state when  $\epsilon = 0$  of shallow water equations with source terms. The coefficient  $\alpha$  allows to reduce the Froude number and  $\epsilon$  enables us to study the perturbation ( $\epsilon = 10^{-5}$ ) of a steady state in low Froude regime ( $\alpha = 10^{-5}$ ). The parameters used in Subsection 1.5 are chosen : a tolerance of the GMRES method of  $10^{-9}$  and 5 time iterations to compute the implicit problem. Eventually, the gravity waves' speed  $c = \sqrt{gh}$  verifies  $c = O(1)$ . We compare two time steps  $\Delta t = 0.1$  and  $\Delta t = 1$ . The results are gathered in Table 8 and correspond to the expected outcomes. Indeed,

	Cells	$\Delta t = 0.1$		$\Delta t = 1$	
		Time	Sub-iter	Time	Sub-iter
GMRES	$24 \times 24$	-	nc	-	nc
	$32 \times 32$	-	nc	-	nc
	$48 \times 48$	-	nc	-	nc
GMRES Jacobi	$24 \times 24$	1.2E+1	255	-	nc
	$32 \times 32$	6.8E+1	720	-	nc
	$48 \times 48$	-	nc	-	nc
GMRES Low Froude	$24 \times 24$	3.1E+1	3	3.9E+2	5
	$32 \times 32$	9.1E+1	3	1.0E+3	5
	$48 \times 48$	4.4E+2	3	4.5E+3	5

TABLE 8. Number of sub-iterations to converge in one time iteration and CPU time by time iteration for different solvers and different time steps.

the method is efficient in the low Froude regime and converges with few iterations even when the classical Jacobi preconditioner is ineffective. As for the linear case, a large part of the CPU cost comes from solving the Schur complement. Indeed, when  $\Delta t$  is large and the acoustic outweighs (*i.e.* in low Froude regime), the dominant operator in the Schur complement is  $-\theta^2 \Delta t^2 g \nabla (h^n \nabla \cdot (h^n \delta \mathbf{u}))$  which is not coercive and admit two scales associated with 0 and  $c = \sqrt{h^n g}$ . Finding a method to treat efficiently this operator is clearly an important question.

As second test case, we propose to compare the low-Froude preconditioner for different low-Froude numbers when the time step ( $\Delta t = 1.0$ ) and the spatial mesh size ( $32 \times 32$ ) are fixed. Each step is solved thanks to an exact solver in order to avoid the problem linked to the inversion of the Schur complement and we focus the study on the global convergence of the method. According to the results in Table 9, the convergence of the GMRES + low-Froude preconditioner is very effective in the low-Froude regime and slower when the Froude

Froude Number	$F_r = O(10^{-5})$	$F_r = O(10^{-3})$	$F_r = O(10^{-2})$	$F_r = O(10^{-1})$	$F_r = O(1)$	$F_r = O(10)$
Low Froude PC	5	5	5	6	23	nc
Any Froude PC	5	5	5	6	20	nc

TABLE 9. Number of iterations for the convergence of the GMRES + low-Froude preconditioner.

number is close to one. These results correspond to assumptions used to construct the Schur complement and match the conclusions obtained with the plane waves' study. They justify also the new method proposed in [8] if the Froude number is around one or larger (this method does not converge for a Froude number around 10). A less expected result is that the Any-Froude preconditioner is not better for Froude number close to one. These results are compatible with the previous result on the propagation property of the any Froude preconditioner. Indeed the wave velocities are good but not the eigenvectors.

## CONCLUSION

In this paper, we have studied the "physics-based" preconditioning methods for hyperbolic systems. This method approximates the Jacobian matrix by a suitable succession of simpler problems and as expected, this technique is efficient if we have a relevant algorithm to solve each sub-systems (which are purely advection-diffusion operators). The approximation between the preconditioner and the model comes from the approximation of the physics (in the nonlinear case) and a time splitting.

We have aimed to analyze the physical approximation by a study based on plane waves. This study allows to verify formally if the preconditioning has the same propagation properties that the model and if it is not the case, we have computed an estimation of the error of the wave velocities. This analysis would be used to find how to correct or to simplify the operators of the preconditioning method for more complex physical models.

Additionally, we have shown three difficulties. Firstly, as expected the classical method is valid only in the Low-Mach (low Froude) regime. The plane wave study shows that the correction proposed in [8] allows to solve partially the problem. The second problem comes from the velocity Schur complement which is non coercive in the high time step limit. Consequently, it will be interesting to find a way to split the two scale linked to the two velocity waves 0 and  $\pm c$  and to solve each scale. This problem is not detailed in the previous works. Thirdly, we have noticed that the implementation of the boundary conditions is important since when we impose weakly the boundary conditions the results are worse. We assume that imposing by penalization the boundary conditions allows to reduce the time splitting errors.

In the future, it will be interesting to extend the method to the shallow water + Exner model and also use the preconditioner with compatible finite element spaces. The compatible spaces given the Inf-Sup theory allows to avoid the spurious pressure mods and allow to design good preconditioning for the Schur operator (the good space to work and design preconditioning is the  $H(\text{div})$  space).

## ACKNOWLEDGEMENTS

The authors would like to thank all the supervisors for their continuous support, especially E. Franck. They also acknowledge all the organizers of the CEMRACS'15.

## REFERENCES

- [1] E. Audusse, and M-O. Bristeau. Transport of pollutant in shallow water, a two time steps kinetic method. *ESAIM: M2AN*, 37(2):389–416, 2003.
- [2] I. Babuska, M. Griebel, and J. Pitkäranta. The problem of selecting the shape functions for a p-type finite element. Research Report BN-1090, November 1988.
- [3] M. Bilanceri, F. Beux, I. Elmahi, H. Guillard, and M.V. Salvetti. Linearised implicit time-advancing applied to sediment transport simulations. Research Report RR-7492, INRIA, December 2010.

- [4] F. Bourdel, P.A. Mazet, and P. Helluy. Resolution of the non-stationary or harmonic Maxwell equations by a discontinuous finite element method. Application to an E.M.I. (electromagnetic impulse) case. In *Proceedings of the 10th international conference on computing methods in applied sciences and engineering*, pages 405–422, 1992.
- [5] H. Brezis. *Functional Analysis, Sobolev Spaces and Partial Differential Equations*. Universitext. Springer New York, 2011.
- [6] F. Brezzi and M. Fortin. *Mixed and Hybrid Finite Element Methods*. Springer Verlag, New York, 1991.
- [7] L. Chacón. An optimal, parallel, fully implicit Newton-Krylov solver for three dimensional viscoresistive magnetohydrodynamics. *Phys. Plasmas*, 15:056103, 2008.
- [8] L. Chacón. Scalable parallel implicit solvers for 3D magnetohydrodynamics. *Journal of Physics: Conference Series*, 125:012041, 2008.
- [9] L. Chacón, and D.A. Knoll. A 2D high- $\beta$  Hall MHD implicit nonlinear solver. *J. Comput. Phys.*, 188(2):573–592, 2003.
- [10] A. Crestetto. *Optimisation de méthodes numériques pour la physique des plasmas. Application aux faisceaux de particules chargées*. PhD thesis, University of Strasbourg, 2012.
- [11] A.-J.-C. de Saint-Venant. Théorie du mouvement non-permanent des eaux, avec application aux crues des rivières et à l'introduction des marées dans leur lit. *C.R. Acad. Sci. Paris, Section Mécanique*, 73:147–154, 1871.
- [12] M. Duruflé, P. Grob, and P. Joly. Influence of Gauss and Gauss-Lobatto quadrature rules on the accuracy of a quadrilateral finite element method in the time domain. *Num. Methods Part. Diff. Eq.*, 25(3):526–551, 2009.
- [13] A. Ern, and J.-L. Guermond. *Theory and Practice of Finite Elements*. Applied Mathematical Sciences. Springer New York, 2004.
- [14] A. Ern, and J.-L. Guermond. Discontinuous Galerkin methods for Friedrichs' systems. I. General theory. *SIAM J. Numer. Anal.*, 44(2):753–778, 2006.
- [15] J.F. Gerbeau, and B. Perthame. Derivation of viscous Saint-Venant system for laminar shallow water; Numerical validation. *Discrete Contin. Dynam. Systems*, 1:89–102, 2001.
- [16] A.J. Grass. Sediments transport by waves and currents. Technical Report No. FL29, SERC London Cent. Mar. Technol., 1981.
- [17] D.D. Schnack. *Lectures in Magnetohydrodynamics : With an Appendix on Extended MHD*. Lect. Notes Phys. 780 (Springer, Berlin Heidelberg), 2009.
- [18] L.C. Van Rijn. Sediment transport. Part I : Bed load transport. *J. Hydraul. Eng.*, 110(10):1431–1456, 1984.