

## NUMERICAL SOLUTION OF THE VISCOUS SURFACE WAVE WITH DISCONTINUOUS GALERKIN METHOD \*

LEI WU<sup>1</sup> AND CHI-WANG SHU<sup>1</sup>

**Abstract.** We consider an incompressible viscous flow without surface tension in a finite-depth domain of two dimensions, with free top boundary and fixed bottom boundary. This system is governed by the Navier–Stokes equations in this moving domain and the transport equation on the moving boundary. In this paper, we construct a stable numerical scheme to simulate the evolution of this system by discontinuous Galerkin method, and discuss the error analysis of the fluid under certain assumptions. Our formulation is mainly based on the geometric structure introduced in [Y. Guo and Ian Tice, *Anal. PDE* **6** (2013) 287–369; Y. Guo and Ian Tice, *Arch. Ration. Mech. Anal.* **207** (2013) 459–531; L. Wu, *SIAM J. Math. Anal.* **46** (2014) 2084–2135], and the natural energy estimate, which is rarely used in the numerical study of this system before.

**Mathematics Subject Classification.** 35Q30, 35R35, 74S05.

Received June 1, 2014.

Published online 25 June 2015.

### 1. INTRODUCTION

We consider an incompressible viscous flow in the moving domain

$$\Lambda(t) = \{y = (y_1, y_2) \in \mathbb{T} \times \mathbb{R} : -1 < y_2 < \eta(y_1, t)\}, \quad (1.1)$$

where  $\mathbb{T}$  denotes the 1-torus. We denote the initial domain  $\Lambda(0) = \Lambda_0$ . For each  $t$ , the flow is described by its velocity and pressure  $(u, p) : \Lambda(t) \mapsto \mathbb{R}^2 \times \mathbb{R}$  which satisfies the incompressible Navier–Stokes equations

$$\begin{cases} \partial_t u + u \cdot \nabla u + \nabla p = \nu \Delta u + S & \text{in } \Lambda(t), \\ \nabla \cdot u = 0 & \text{in } \Lambda(t), \\ (pI - \nu \mathbb{D}(u))\mu = g\eta\mu & \text{on } \{y_2 = \eta(y_1, t)\}, \\ u = 0 & \text{on } \{y_2 = -1\}, \\ \partial_t \eta = u_2 - u_1 \partial_{y_1} \eta & \text{on } \{y_2 = \eta(y_1, t)\}, \\ u(t=0) = u_0 & \text{in } \Lambda_0, \\ \eta(t=0) = \eta_0 & \text{on } \mathbb{T}, \end{cases} \quad (1.2)$$

---

*Keywords and phrases.* Stability, free boundary, Navier–Stokes equation.

\* *Research supported by DOE grant DE-FG02-08ER25863 and NSF grants DMS-1112700 and DMS-1418750.*

<sup>1</sup> Division of Applied Mathematics, Brown University, Providence, RI 02912, USA. [Lei\\_Wu@brown.edu](mailto:Lei_Wu@brown.edu); [shu@dam.brown.edu](mailto:shu@dam.brown.edu)

for  $\mu$  the outward-pointing unit normal vector on  $\{y_2 = \eta\}$ ,  $I$  the  $2 \times 2$  identity matrix,  $(\mathbb{D}u)_{ij} = \partial_i u_j + \partial_j u_i$  the symmetric gradient of  $u$ ,  $g$  the gravitational constant,  $\nu > 0$  the viscosity and  $S = S(t, y_1, y_2)$  an external source term. The fifth equation in (1.2) implies the free surface is convected with the fluid. Note in (1.2), we have shifted the actual pressure  $\bar{p}$  by the constant atmosphere pressure  $p_{atm}$  according to  $p = \bar{p} + gy_2 - p_{atm}$ .

We always assume the natural condition that there exists a positive number  $\delta$  such that  $\eta_0 + 1 \geq \delta > 0$  on  $\Sigma$ , which means the initial free surface is strictly separated from the bottom.

Traditionally, based on the handling of the free surface, this type of problems can be solved via moving-grid technique as in [17], marker-and-cell method as in [13], volume-of-fluid method as in [14] and level-set method as in [9]. In each case, finite difference method, finite volume method and finite element method can be applied to solve the Navier–Stokes equation in  $\Lambda(t)$ . However, to the best of authors’ knowledge, there is very little in the literature on solving this problem with discontinuous Galerkin method other than [9]. On the other hand, in spite of many computational tests presented in the literature, there is very little discussion on the stability and convergence of the numerical scheme. In this paper, we employ the idea from [10, 11, 20], to construct a stable numerical scheme and give detailed analysis.

Our central idea is to flatten the free surface *via* a coordinates transform. First, we define a fixed domain

$$\Omega = \{x = (x_1, x_2) \in \Sigma \times \mathbb{R} \mid -1 < x_2 < 0\}, \tag{1.3}$$

for which we write the coordinates  $x \in \Omega$ . In this slab, we take  $\Sigma : \{x_2 = 0\}$  as the upper boundary and  $\Sigma_b : \{x_2 = -1\}$  as the lower boundary.

Consider the geometric transform from  $\Omega$  to  $\Lambda(t)$ , which is first introduced in [2] and further extended in [10, 20]:

$$\Phi : (x_1, x_2) \mapsto (x_1, x_2 + \eta(1 + x_2)) = (y_1, y_2). \tag{1.4}$$

We may directly verify this transform maps  $\Omega$  into  $\Lambda(t)$  with the Jacobi matrix

$$\nabla\Phi = \begin{pmatrix} 1 & 0 \\ A & J \end{pmatrix}, \tag{1.5}$$

and the transform matrix

$$\mathcal{A} = ((\nabla\Phi)^{-1})^T = \begin{pmatrix} 1 & -AK \\ 0 & K \end{pmatrix}, \tag{1.6}$$

where

$$\begin{aligned} \tilde{b} &= 1 + x_2, & A &= \partial_1 \eta \tilde{b}, \\ J &= 1 + \eta, & K &= 1/J. \end{aligned} \tag{1.7}$$

Here we denote the derivative with respect to  $x_1$  as  $\partial_1$  and with respect to  $x_2$  as  $\partial_2$ . Define the transformed operators as follows:

$$\begin{aligned} (\nabla_{\mathcal{A}} f)_i &= \mathcal{A}_{ij} \partial_j f, \\ \nabla_{\mathcal{A}} \cdot \mathbf{g} &= \mathcal{A}_{ij} \partial_j g_i, \\ \Delta_{\mathcal{A}} f &= \nabla_{\mathcal{A}} \cdot \nabla_{\mathcal{A}} f, \\ \mathcal{N} &= (-\partial_1 \eta, 1), \\ \chi &= \partial_1 \eta, \\ (\mathbb{D}_{\mathcal{A}} u)_{ij} &= \mathcal{A}_{ik} \partial_k u_j + \mathcal{A}_{jk} \partial_k u_i, \\ S_{\mathcal{A}}(p, u) &= pI - \mathbb{D}_{\mathcal{A}} u, \end{aligned} \tag{1.8}$$

where the summation index should be understood in the Einstein convention. If we extend the divergence  $\nabla_{\mathcal{A}} \cdot$  to act on symmetric tensors in the natural way, then a straightforward computation reveals  $\nabla_{\mathcal{A}} \cdot S_{\mathcal{A}}(p, u) = \nabla_{\mathcal{A}} p - \Delta_{\mathcal{A}} u$  for vector fields satisfying  $\nabla_{\mathcal{A}} \cdot u = 0$ .

In our new coordinates, the original system (1.2) becomes

$$\begin{cases} \partial_t u - \partial_t \eta \tilde{b} K \partial_2 u + u \cdot \nabla_{\mathcal{A}} u - \nu \Delta_{\mathcal{A}} u + \nabla_{\mathcal{A}} p = S & \text{in } \Omega, \\ \nabla_{\mathcal{A}} \cdot u = 0 & \text{in } \Omega, \\ S_{\mathcal{A}}(p, u) \mathcal{N} = g \eta \mathcal{N} & \text{on } \Sigma, \\ u = 0 & \text{on } \Sigma_b, \\ u(x, 0) = u_0(x) & \text{in } \Omega, \\ \partial_t \eta = u \cdot \mathcal{N} & \text{on } \Sigma, \\ \eta(x', 0) = \eta_0(x') & \text{on } \Sigma. \end{cases} \tag{1.9}$$

Since  $\mathcal{A}$  depends on  $\eta$  through the transform, most of the operators in the Navier–Stokes equations are related to the free surface  $\eta$ . Hence, the Navier–Stokes equations and the transport equation are essentially coupled.

Based on [10], equation (1.9) with  $S = 0$  possesses a natural energy equality as follows:

$$\int_{\Omega} J(t) |u(t)|^2 + g \int_{\Sigma} |\eta(t)|^2 + \nu \int_0^t \int_{\Omega} J(s) |\mathbb{D}_{\mathcal{A}} u(s)|^2 ds = \int_{\Omega} J(0) |u(0)|^2 + g \int_{\Sigma} |\eta(0)|^2. \tag{1.10}$$

Hence, we try to construct a numerical scheme to recover this energy stability for the numerical solution.

In the following, we refer to the term “continuous case” when we consider the exact solution triple  $(u, p, \eta)$  which is sufficiently smooth. On the other hand, we refer to the term “discrete case” when we consider the numerical solution triple  $(u_h, p_h, \eta_h)$ .

Throughout this paper,  $C > 0$  denotes a positive constant that only depends on the parameters  $\Omega$ ,  $g$  and  $\nu$  of the problem, and the exact solution  $(u, p, \eta)$ . It is referred as universal and can change from one inequality to another. When we write  $C(z)$ , it means a positive constant depending on the quantity  $z$ .  $a \lesssim b$  denotes  $a \leq Cb$ , where  $C$  is a universal constant as defined above.

The method we discuss in this paper belongs to the class of discontinuous Galerkin (DG) methods. DG methods are finite element methods which use completely discontinuous piecewise polynomial solution and test spaces. They are particularly useful for convection dominated wave problems, and have the advantage of flexibility in adaptivity and efficient parallel implementation. We refer to the references [3–5, 19] for more details.

## 2. NUMERICAL SCHEME

### 2.1. Fundamental settings

In our construction and analysis of the numerical scheme, we focus on the semi-discrete form of the system. The time discretization is based on the Runge–Kutta method for differential-algebraic equations of index 2, as presented in [12]. We do not discuss the fully-discrete scheme in this paper.

We choose the bulk domain as  $\Omega : (x_1, x_2) \in [0, 1] \times [-1, 0]$  and the surface domain as  $\Sigma : x_1 \in [0, 1]$ , which are 1-periodic in the  $x_1$  direction. The surface mesh is defined by dividing  $\Sigma$  into  $N$  uniform elements with length  $h = 1/N$ . The bulk mesh construction can be divided into two steps: first we divide  $\Omega$  into  $N \times N$  uniform squares, where each element is sized  $h \times h$ . Then, we divide each square into two right triangles by cutting along the diagonal from the left-up corner to the right-down corner. The bulk mesh is shown in Figure 1.

**Remark 2.1.** Since the discretization of the weighted differential operators  $\nabla_{\mathcal{A}}$ ,  $\nabla_{\mathcal{A}} \cdot$  and  $\Delta_{\mathcal{A}}$  depends on  $\eta$ , *i.e.* the surface variable, we need the bulk discretization to match the surface discretization in the vertical direction. Hence, we need to first choose the strip-shape mesh as in Figure 2, whose projection on the upper

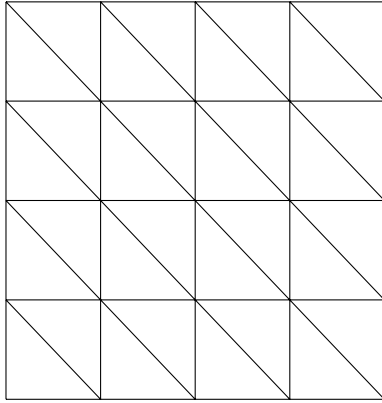


FIGURE 1. Mesh distribution in  $N = 4$ .

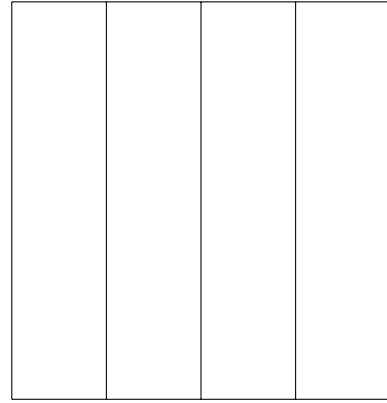


FIGURE 2. Strip-shape mesh distribution in  $N = 4$ .

boundary is exactly the surface mesh. Then in each strip, any triangulation is permitted. For convenience, we made the choice as in Figure 1.

Let  $\mathcal{E}_h$  be the set consisting of all the triangular elements in the bulk mesh and  $\partial\mathcal{E}_h$  be the set of all the sides of the triangular cells. Let  $\mathcal{F}_h$  be the set consisting of all the interval elements in the surface mesh and  $\partial\mathcal{F}_h$  be the set of all the boundary points of intervals.

Define the usual Sobolev space in  $\Omega$  as  $H^k(\Omega)$  and in  $\Sigma$  as  $H^k(\Sigma)$ . We define the space  $P_h^k(\Omega)$  where  $f \in P_h^k(\Omega)$  if and only if  $f|_E \in P^k(E)$  for all  $E \in \mathcal{E}_h$  ( $P^k(E)$  denotes the set of polynomials of degree at most  $k$  defined on  $E$ ). Similarly, we can define the space  $P_h^k(\Sigma)$ .

We define the solution spaces as

$$X_h^k = \left\{ u_h \in (L^2(\Omega))^2 : u_h \in (P_h^k(\Omega))^2 \right\}, \tag{2.1}$$

$$M_h^k = \{ p_h \in L^2(\Omega) : p_h \in P_h^k(\Omega) \}, \tag{2.2}$$

$$S_h^k = \{ \eta_h \in L^2(\Sigma) : \eta_h \in P_h^k(\Sigma) \}, \tag{2.3}$$

for  $k \geq 1$ .

For discrete functions  $f_h \in X_h^k$  and  $g_h \in S_h^k$ , we may define the  $H_h$  norms as follows:

$$\|f_h\|_{H_h}^2 = \sum_{E \in \mathcal{E}_h} \|f_h\|_{H^1(E)}^2, \tag{2.4}$$

$$\|g_h\|_{H_h}^2 = \sum_{F \in \mathcal{F}_h} \|g_h\|_{H^1(F)}^2. \tag{2.5}$$

Define the  $X_h$  norm for  $v_h \in X_h^k$  as

$$\|v_h\|_{X_h}^2 = \sum_{E \in \mathcal{E}_h} \left( \|\nabla v_h\|_{L^2(E)}^2 + \int_{\partial E} \frac{1}{h} [v_h]^2 \right). \tag{2.6}$$

Before stating the scheme, we announce the notations used below. For a given boundary belonging to a reference cell  $E \in \mathcal{E}_h$  and a function  $f$  defined in  $\Omega$ ,  $f^{\text{ext}(E)}$  denotes the value of  $f$  read from the exterior direction on the

boundary and  $f^{\text{int}(E)}$  denotes that read from the interior direction. Also for the boundary shared by two cells, we define

$$\{f\} = \frac{1}{2} \left( f^{\text{ext}(E)} + f^{\text{int}(E)} \right), \quad (2.7)$$

$$[f] = f^{\text{int}(E)} - f^{\text{ext}(E)}. \quad (2.8)$$

For the side on the boundary of  $\Omega$ , we make a special definition that on  $\Sigma_b$ ,  $f^{\text{int}(E)} = \{f\} = [f]$  and  $f^{\text{ext}(E)} = 0$ , while on  $\Sigma$ ,  $f^{\text{int}(E)} = f^{\text{ext}(E)} = \{f\}$  and  $[f] = 0$ . The definitions are reasonable, since the system (1.9) for  $u$  gives Dirichlet-type condition on  $\Sigma_b$  and Neumann-type condition on  $\Sigma$ .

The numerical scheme is based on the weak formulation of the system (1.9). Taking test functions  $v_h \in X_h^k$  and  $q_h \in M_h^{k-1}$ . We multiply  $Jv_h$  and  $Jq_h$  respectively on both sides of the Navier–Stokes equations and integrate over  $E \in \mathcal{E}_h$  to obtain

$$\begin{aligned} & \int_E J \partial_t u \cdot v_h - \int_E J \left( \partial_t \tilde{\eta} b K \partial_2 u \right) \cdot v_h \\ & + \int_E J (u \cdot \nabla_{\mathcal{A}} u) \cdot v_h - \nu \int_E J \Delta_{\mathcal{A}} u \cdot v_h + \int_E J \nabla_{\mathcal{A}} p \cdot v_h = \int_E J (S \cdot v_h), \end{aligned} \quad (2.9)$$

$$\int_E J (\nabla_{\mathcal{A}} \cdot u) q_h = 0. \quad (2.10)$$

Considering the transport equation in (1.9), we make an addition and subtraction to modify two terms in above formulation as

$$\begin{aligned} & \left( \int_E J \partial_t u \cdot v_h - \int_E J \left( \partial_t \tilde{\eta} b K \partial_2 u \right) \cdot v_h + \frac{1}{2} \int_{\partial E \cap \Sigma} \partial_t \eta (u \cdot v_h) \right) \\ & + \left( \int_E J (u \cdot \nabla_{\mathcal{A}} u) \cdot v_h - \frac{1}{2} \int_{\partial E \cap \Sigma} (u \cdot \mathcal{N})(u \cdot v_h) \right) - \nu \int_E J \Delta_{\mathcal{A}} u \cdot v_h + \int_E J \nabla_{\mathcal{A}} p \cdot v_h = \int_E J (S \cdot v_h), \end{aligned} \quad (2.11)$$

$$\int_E J (\nabla_{\mathcal{A}} \cdot u) q_h = 0, \quad (2.12)$$

where  $\partial E$  denotes the boundary of the triangular element  $E$ . This trick is to enforce the transport relation and shows its power in the stability proof.

Multiplying a test functions  $\phi_h \in S_h^k$  on both sides of the transport equation and integrating over  $F \in \mathcal{F}_h$ , we obtain

$$\int_F \partial_t \eta \phi_h + \int_F \bar{u}_1 \partial_1 \eta \phi_h - \int_F \bar{u}_2 \phi_h = 0, \quad (2.13)$$

where  $\bar{u}_1$  and  $\bar{u}_2$  denote the traces of  $u_1$  and  $u_2$  on  $\Sigma$ .

The weak formulations of the system (1.9) are based on the integration by parts of (2.11), (2.12) and (2.13). In the following, we define and analyze the numerical scheme term by term.

## 2.2. Discretization

We define the semi-discrete scheme as follows:

$$\left\{ \begin{array}{l} \kappa_F (\eta_h, u_h, \phi_h) = 0, \\ \lambda_F (\chi_h, \psi_h) = 0, \\ \zeta_E (\eta_h, u_h, v_h) + \gamma_E (u_h, \eta_h, u_h, v_h) + \nu \alpha_E (\eta_h, u_h, \eta_h, v_h) \\ \quad + \beta_E (p_h, \eta_h, v_h) + g \mu_E (\eta_h, \eta_h, v_h) = \omega_E (\eta_h, v_h), \\ \rho_E (u_h, \eta_h, q_h) = 0, \end{array} \right. \quad (2.14)$$

for any test functions  $\phi_h \in S_h^k$ ,  $\psi_h \in S_h^k$ ,  $v_h \in X_h^k$  and  $q_h \in M_h^{k-1}$ , where the first two equations denote the transport discretization and the last two equations denote the fluid discretization. This forms a complete differential-algebraic system for  $\eta_h(t) \in S_h^k$ ,  $\chi_h(t) \in S_h^k$ ,  $u_h(t) \in X_h^k$  and  $p_h(t) \in M_h^{k-1}$ . The detailed definitions of above multi-linear terms  $\kappa$ ,  $\lambda$ ,  $\zeta$ ,  $\gamma$ ,  $\alpha$ ,  $\beta$ ,  $\mu$ ,  $\omega$ , and  $\rho$  are in the following.

2.2.1. Discretization of transport terms as  $\kappa_F$  and  $\lambda_F$

We define

$$\begin{aligned} \kappa_F(\eta_h, u_h, \phi_h) &= \int_F \partial_t \eta_h \phi_h - \int_F (\bar{u}_2)_h \phi_h - \int_F \partial_1 (\bar{u}_1)_h \eta_h \phi_h - \int_F (\bar{u}_1)_h \eta_h \partial_1 \phi_h \\ &\quad + (\hat{u}_1)_h \hat{\eta}_h \phi_h^-|_{F+1/2} - (\hat{u}_1)_h \hat{\eta}_h \phi_h^+|_{F-1/2} + \frac{1}{2} (\eta_h^- \phi_h^- [u_h]|_{F+1/2} - \eta_h^+ \phi_h^+ [u_h]|_{F-1/2}) \end{aligned} \tag{2.15}$$

$$\lambda_F(\chi_h, \psi_h) = \int_F \chi_h \psi_h + \int_F \eta_h \partial_1 \psi_h - \eta_h^+ \psi_h^-|_{F+1/2} + \eta_h^+ \psi_h^+|_{F-1/2}, \tag{2.16}$$

where  $F - 1/2$  and  $F + 1/2$  denote the left and right boundary points of  $F$ . Here, for each boundary point, “-” means the value read from the left and “+” means the value read from the right. Also,  $(\hat{u}_1)_h \hat{\eta}_h$  denotes the numerical flux on the boundary. We always take

$$(\hat{u}_1)_h = \{(\bar{u}_1)_h\}, \tag{2.17}$$

and  $\hat{\eta}_h$  should be determined from the sign of  $\{(\bar{u}_1)_h\}$  following the upwinding rule as

$$\hat{\eta}_h = \begin{cases} \eta_h^- & \text{if } \{(\bar{u}_1)_h\} \geq 0, \\ \eta_h^+ & \text{if } \{(\bar{u}_1)_h\} < 0, \end{cases} \tag{2.18}$$

where  $(\bar{u}_1)_h$  is the trace of  $(u_1)_h$  on  $\Sigma$ . Note  $\frac{1}{2} (\eta_h^- \phi_h^- [u_h]|_{F+1/2} - \eta_h^+ \phi_h^+ [u_h]|_{F-1/2})$  are extra penalty terms, which helps to show the stability. Since  $u_h$  is also discontinuous across the boundary point, these penalty terms are used to transform  $u_h$  into  $\{u_h\}$  on the boundary after integrating by parts in (2.15).

In all the applications below, we use  $\chi_h$  to discretize  $\partial_1 \eta$ , and  $\eta_h$  to discretize  $\eta$ . Hence, we denote  $\mathcal{N}_h$  for the vector  $(-\chi_h, 1)$ . Also,  $A_h, J_h, K_h$  and  $\mathcal{A}_h$  can be defined in the same convention.

2.2.2. Spacial discretization of temporal terms as  $\zeta_E$

We define

$$\begin{aligned} \zeta_E(\eta_h, u_h, v_h) &= \int_E \partial_t (J_h u_h) \cdot v_h + \int_E \partial_t \eta_h \tilde{b} u_h \cdot \partial_2 v_h - \int_{\partial E} \partial_t \eta_h \tilde{b} (\hat{u}_h \cdot v_h^{\text{int}(E)}) n_2 \\ &\quad + \frac{1}{2} \int_{\partial E \cap \Sigma} \partial_t \eta_h (u_h \cdot v_h), \end{aligned} \tag{2.19}$$

where  $n_E = (n_1, n_2)$  denotes the outward normal vector on  $\partial E$ . The last term in (2.19) is the newly-added penalty term in the formulation (2.11). The argument  $\eta_h$  in  $\zeta_E(\eta_h, u_h, v_h)$  denotes  $\eta_h$  and its derivatives. We utilize the upwinding flux in this scheme, *i.e.*

$$\hat{u}_h = \begin{cases} u_h^+ & \text{if } \partial_t \eta_h \tilde{b} \geq 0, \\ u_h^- & \text{if } \partial_t \eta_h \tilde{b} < 0, \end{cases} \tag{2.20}$$

where “+” denotes the value read from the up direction and “-” denote that from the down direction. If a side is vertical, then its contribution is zero.

### 2.2.3. Discretization of convection terms as $\gamma_E$

We define

$$\begin{aligned} \gamma_E(u_h, \eta_h, u_h, v_h) &= -\frac{1}{2} \int_{\partial E \cap \Sigma} (u_h \cdot \mathcal{A}_h) (u_h \cdot v_h) + \int_E J_h (u_h \cdot \nabla_{\mathcal{A}_h} u_h) \cdot v_h \\ &\quad + \frac{1}{2} \int_E J_h (\nabla_{\mathcal{A}_h} \cdot u_h) (u_h \cdot v_h) - \frac{1}{2} \int_{\partial E \setminus \partial \Omega} [u_h \cdot J_h \mathcal{A}_h] \cdot n_E \frac{u_h^{\text{int}(E)} \cdot v_h^{\text{int}(E)}}{2} \\ &\quad + \int_{\partial E^-} |\{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E| \left( u_h^{\text{int}(E)} - u_h^{\text{ext}(E)} \right) \cdot v_h^{\text{int}(E)}, \end{aligned} \quad (2.21)$$

where  $\nabla_{\mathcal{A}_h}$  is understood as in (1.8) with  $\mathcal{A}$  replaced by  $\mathcal{A}_h$  and

$$\partial E^- = \{x \in \partial E : \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E < 0\}. \quad (2.22)$$

The first term in (2.21) is the newly added penalty term in the formulation (2.11). Since

$$J_h \mathcal{A}_h = \begin{pmatrix} J_h - A_h \\ 0 & 1 \end{pmatrix}, \quad (2.23)$$

we only need one  $\eta_h$  argument in  $\gamma_E(u_h, \eta_h, u_h, v_h)$ .

### 2.2.4. Discretization of diffusion terms as $\alpha_E$

We utilize the Symmetric Internal Penalty Galerkin Method (SIPG) or the Nonsymmetric Internal Penalty Galerkin Method (NIPG) to define

$$\begin{aligned} \alpha_E(\eta_h, u_h, \eta_h, v_h) &= \frac{1}{2} \int_E J_h \mathbb{D}_{\mathcal{A}_h} u_h : \mathbb{D}_{\mathcal{A}_h} v_h - \int_{\partial E \setminus \Sigma} \{\mathbb{D}_{\mathcal{A}_h} u_h J_h \mathcal{A}_h\} \cdot n_E \cdot v_h^{\text{int}(E)} \\ &\quad \pm \frac{1}{2} \int_{\partial E \setminus \Sigma} \mathbb{D}_{\mathcal{A}_h^{\text{int}(E)}} v_h^{\text{int}(E)} J_h^{\text{int}(E)} \mathcal{A}_h^{\text{int}(E)} \cdot n_E \cdot [u_h] + \frac{\sigma}{h} \int_{\partial E} [u_h] \cdot v_h^{\text{int}(E)}. \end{aligned} \quad (2.24)$$

For the third term, if we take +, it is NIPG and if we take -, it is SIPG. The penalty constant  $\sigma$  can always be taken as 1 for NIPG and a sufficiently large number for SIPG. Note that we need two  $\eta_h$  arguments in  $\alpha_E(\eta_h, u_h, \eta_h, v_h)$  to denote the dependence of  $\mathcal{A}_h$ .

### 2.2.5. Discretization of pressure terms as $\beta_E$

We define

$$\beta_E(p_h, \eta_h, v_h) = - \int_E J_h p_h \nabla_{\mathcal{A}_h} \cdot v_h + \int_{\partial E \setminus \Sigma} v_h^{\text{int}(E)} \cdot \{p_h J_h \mathcal{A}_h\} \cdot n_E. \quad (2.25)$$

### 2.2.6. Discretization of forcing terms as $\mu_E$

We define

$$\mu_E(\eta_h, \eta_h, v_h) = \int_{\partial E \cap \Sigma} \eta_h (v_2)_h - \int_{\partial E \cap \Sigma} \eta_h \partial_1 \eta_h (\bar{v}_1)_h + \frac{1}{4} [\eta_h^2] \bar{v}_h^- |_{\partial E \cap \Sigma + 1/2} - \frac{1}{4} [\eta_h^2] \bar{v}_h^+ |_{\partial E \cap \Sigma - 1/2}, \quad (2.26)$$

where  $\partial E \cap \Sigma + 1/2$  and  $\partial E \cap \Sigma - 1/2$  denote the left and right boundary points of  $\partial E \cap \Sigma$  respectively, *i.e.* similar to  $F - 1/2$  and  $F + 1/2$ . Note that the forcing term is nontrivial only for the top cells and for all other cells it is zero,  $\mu_E(\eta_h, \eta_h, v_h) = 0$ .  $\frac{1}{4} [\eta_h^2] \bar{v}_h^- |_{\partial E \cap \Sigma + 1/2} - \frac{1}{4} [\eta_h^2] \bar{v}_h^+ |_{\partial E \cap \Sigma - 1/2}$  are extra penalty terms, which help to show the stability. Since  $\eta_h$  is discontinuous across the boundary, these penalty terms are used to transform  $\eta_h^2$  into  $\{\eta_h^2\}$  on the boundary after integrating by parts in (2.26).

2.2.7. Discretization of source terms as  $\omega_E$

We define

$$\omega_E(\eta_h, v_h) = \int_E J_h S \cdot v_h. \tag{2.27}$$

2.2.8. Discretization of divergence terms as  $\rho_E$

We define

$$\rho_E(u_h, \eta_h, q_h) = - \int_E J_h q_h \nabla_{\mathcal{A}_h} \cdot u_h + \frac{1}{2} \int_{\partial E \setminus \Sigma} [u_h] \cdot q_h^{\text{int}(E)} J_h^{\text{int}(E)} \mathcal{A}_h^{\text{int}(E)} \cdot n_E. \tag{2.28}$$

2.3. Properties and estimates of discretization

2.3.1. Estimate of temporal terms

We consider the temporal terms, *i.e.*

$$\left( \int_E J \partial_t u \cdot v_h - \int_E J \left( \partial_t \tilde{\eta} \tilde{b} K \partial_2 u \right) \cdot v_h + \frac{1}{2} \int_{\partial E \cap \Sigma} \partial_t \eta (u \cdot v_h) \right). \tag{2.29}$$

In the continuous case, a direct integration by parts reveals

$$\begin{aligned} \int_E J \partial_t u \cdot v_h &= \int_E \partial_t (J u) \cdot v_h - \int_E \partial_t J (u \cdot v_h), \\ - \int_E J \left( \partial_t \tilde{\eta} \tilde{b} K \partial_2 u \right) \cdot v_h &= - \int_E \partial_t \tilde{\eta} \tilde{b} \partial_2 u \cdot v_h \\ &= \int_E \partial_t \tilde{\eta} \tilde{b} u \cdot \partial_2 v_h + \int_E \partial_2 \left( \partial_t \tilde{\eta} \tilde{b} \right) (u \cdot v_h) - \int_{\partial E} \partial_t \tilde{\eta} \tilde{b} (u \cdot v_h) n_2. \end{aligned} \tag{2.31}$$

Since  $\partial_t J = \partial_2(\partial_t \tilde{\eta} \tilde{b})$ , we can simplify above terms into

$$\int_E J \partial_t u \cdot v_h - \int_E J \left( \partial_t \tilde{\eta} \tilde{b} K \partial_2 u \right) \cdot v_h = \int_E \partial_t (J u) \cdot v_h + \int_E \partial_t \tilde{\eta} \tilde{b} u \cdot \partial_2 v_h - \int_{\partial E} \partial_t \tilde{\eta} \tilde{b} (u \cdot v_h) n_2. \tag{2.32}$$

Hence, we may define the discretization of the temporal term as (2.19). Since in our bulk mesh, there are only two effective boundary integrals for each triangle, we may call them the upper boundary  $\partial E + 1/2$  and lower boundary  $\partial E - 1/2$  without ambiguity. If we denote  $f(u_h) = -\partial_t \tilde{\eta} \tilde{b} u_h$ , then (2.19) is actually

$$\begin{aligned} \zeta_E(\eta_h, u_h, v_h) &- \frac{1}{2} \int_{\partial E \cap \Sigma} \partial_t \eta_h (u_h \cdot v_h) \\ &= \int_E \partial_t (J_h u_h) \cdot v_h - \int_E f(u_h) \cdot \partial_2 v_h + \int_{\partial E + 1/2} \hat{f}(u_h) \cdot v_h^- - \int_{\partial E - 1/2} \hat{f}(u_h) \cdot v_h^+, \end{aligned} \tag{2.33}$$

where the flux reduces to

$$\hat{f}(u_h) = f(\hat{u}_h) = \begin{cases} f(u_h^+) & \text{if } \partial_t \tilde{\eta} \tilde{b} \geq 0, \\ f(u_h^-) & \text{if } \partial_t \tilde{\eta} \tilde{b} < 0. \end{cases} \tag{2.34}$$

**Lemma 2.2.** *If we take  $v_h = u_h$ , the discretized temporal term satisfies*

$$\sum_{E \in \mathcal{E}_h} \zeta_E(\eta_h, u_h, u_h) \geq \frac{1}{2} \partial_t \int_{\Omega} J_h |u_h|^2. \tag{2.35}$$



*Proof.* Note the fact that  $\partial_t J_h = \partial_2(\partial_t \eta_h \tilde{b})$ . We can directly compute

$$\begin{aligned}
 & \zeta_E(\eta_h, u_h, u_h) - \frac{1}{2} \int_{\partial E \cap \Sigma} \partial_t \eta_h |u_h|^2 \\
 &= \int_E \partial_t (J_h u_h) \cdot u_h - \int_E f(u_h) \cdot \partial_2 u_h + \int_{\partial E+1/2} \hat{f}(u_h) \cdot u_h^- - \int_{\partial E-1/2} \hat{f}(u_h) \cdot u_h^+ \\
 &= \frac{1}{2} \partial_t \int_E J_h |u_h|^2 + \frac{1}{2} \int_E \partial_t J_h |u_h|^2 - \frac{1}{2} \int_E \partial_2 (\partial_t \eta_h \tilde{b}) |u_h|^2 \\
 &\quad - \frac{1}{2} \int_{\partial E+1/2} f(u_h^-) \cdot u_h^- + \frac{1}{2} \int_{\partial E-1/2} f(u_h^+) \cdot u_h^+ + \int_{\partial E+1/2} \hat{f}(u_h) \cdot u_h^- - \int_{\partial E-1/2} \hat{f}(u_h) \cdot u_h^+ \\
 &= \frac{1}{2} \partial_t \int_E J_h |u_h|^2 \\
 &\quad - \frac{1}{2} \int_{\partial E+1/2} f(u_h^-) \cdot u_h^- + \frac{1}{2} \int_{\partial E-1/2} f(u_h^+) \cdot u_h^+ + \int_{\partial E+1/2} \hat{f}(u_h) \cdot u_h^- - \int_{\partial E-1/2} \hat{f}(u_h) \cdot u_h^+ \\
 &= \frac{1}{2} \partial_t \int_E J_h |u_h|^2 + \hat{F}_{\partial E+1/2} - \hat{F}_{\partial E-1/2} + \Theta_E, \tag{2.36}
 \end{aligned}$$

where

$$\hat{F}_{\partial E+1/2} = -\frac{1}{2} \int_{\partial E+1/2} f(u_h^-) \cdot u_h^- + \int_{\partial E+1/2} \hat{f}(u_h) \cdot u_h^-, \tag{2.37}$$

$$\hat{F}_{\partial E-1/2} = -\frac{1}{2} \int_{\partial E-1/2} f(u_h^+) \cdot u_h^+ + \int_{\partial E-1/2} \hat{f}(u_h) \cdot u_h^+, \tag{2.38}$$

and

$$\Theta_E = \frac{1}{2} \int_{\partial E-1/2} (f(u_h^+) \cdot u_h^+ - f(u_h^-) \cdot u_h^-) + \int_{\partial E-1/2} (\hat{f}(u_h) \cdot u_h^- - \hat{f}(u_h) \cdot u_h^+). \tag{2.39}$$

Based on the flux definition (2.34), we have the estimate

$$\Theta_E = \int_{\partial E-1/2} \left( \partial_t \eta_h \tilde{b} (u_h^+ - u_h^-) \cdot \left( \hat{u}_h - \frac{u_h^+ + u_h^-}{2} \right) \right) \geq 0. \tag{2.40}$$

Hence, combining with (2.36), we obtain

$$\zeta_E(\eta_h, u_h, u_h) - \frac{1}{2} \int_{\partial E \cap \Sigma} \partial_t \eta_h |u_h|^2 \geq \frac{1}{2} \partial_t \int_E J_h |u_h|^2 + \hat{F}_{\partial E+1/2} - \hat{F}_{\partial E-1/2}. \tag{2.41}$$

When summing up over all  $E \in \mathcal{E}_h$ , we can easily see when  $\partial E \subset \Omega$ , all the terms involving  $\hat{F}$  are canceled out. Therefore, only the terms on  $\partial\Omega$  remain, *i.e.*

$$\sum_{E \in \mathcal{E}_h} \zeta_E(\eta_h, u_h, u_h) - \frac{1}{2} \int_{\Sigma} \partial_t \eta_h |u_h|^2 \geq \sum_{E \in \mathcal{E}_h} \frac{1}{2} \partial_t \int_E J_h |u_h|^2 + \sum_{\partial E \cap \Sigma \neq \emptyset} \hat{F}_{\partial E+1/2} - \sum_{\partial E \cap \Sigma_b \neq \emptyset} \hat{F}_{\partial E-1/2}. \tag{2.42}$$

On  $\Sigma_b$  we have  $u_h^- = 0$ , so  $-\hat{F}_{\partial E-1/2} = 0$ , *i.e.*

$$\sum_{\partial E \cap \Sigma_b \neq \emptyset} \hat{F}_{\partial E-1/2} = 0. \tag{2.43}$$

On  $\Sigma$  we have  $u_h^+ = u_h^- = u_h|_\Sigma$  and we always take  $\hat{u}_h = u_h$ , so it implies

$$\hat{F}_{\partial E+1/2} = -\frac{1}{2} \int_{\partial E+1/2} f(u_h^-) u_h^- + \int_{\partial E+1/2} \hat{f}(u_h) u_h^- = -\frac{1}{2} \int_{\partial E \cap \Sigma} \partial_t \eta_h |u_h|^2. \tag{2.44}$$

Then summing up over  $E \in \mathcal{E}_h$  gives a full integration over  $\Sigma$ , *i.e.*

$$\sum_{\partial E \cap \Sigma \neq \emptyset} \hat{F}_{\partial E+1/2} = -\frac{1}{2} \int_{\Sigma} \partial_t \eta_h |u_h|^2. \tag{2.45}$$

Therefore, combining (2.42), (2.43) and (2.45), we deduce

$$\sum_{E \in \mathcal{E}_h} \zeta_E(\eta_h, u_h, u_h) - \frac{1}{2} \int_{\Sigma} \partial_t \eta_h |u_h|^2 \geq \frac{1}{2} \partial_t \int_E J_h |u_h|^2 - \frac{1}{2} \int_{\Sigma} \partial_t \eta_h |u_h|^2. \tag{2.46}$$

Hence, our result easily follows. □

**Remark 2.3.** This proof is based on the mesh as in Figure 1. For general triangulation, it is possible to have three effective boundaries for each element. However, it is easy to see Lemma 2.2 still holds.

2.3.2. Estimate of convection terms

We consider the convection term

$$\int_E J(u \cdot \nabla_{\mathcal{A}} u) \cdot v_h - \frac{1}{2} \int_F (u \cdot \mathcal{N})(u \cdot v_h). \tag{2.47}$$

This is the key nonlinear term in the Navier–Stokes equations. Our discretization is inspired by the idea in [7]. In the continuous case, when summing up over all  $E \in \mathcal{E}_h$ , we can see

$$\sum_{E \in \mathcal{E}_h} \gamma_E(u, \eta, u, v_h) = \int_{\Omega} J(u \cdot \nabla_{\mathcal{A}} u) \cdot v_h + \frac{1}{2} \int_{\Omega} J(\nabla_{\mathcal{A}} \cdot u)(u \cdot v_h) - \frac{1}{2} \int_{\Sigma} (u \cdot \mathcal{N})(u \cdot v_h), \tag{2.48}$$

where all the other boundary terms vanish. Hence, considering the continuous  $\mathcal{A}$ -divergence-free condition for  $u$ , this discretization is consistent.

**Lemma 2.4.** *If we take  $v_h = u_h$ , the discretized convection term satisfies*

$$\sum_{E \in \mathcal{E}_h} \gamma_E(u_h, \eta_h, u_h, u_h) \geq 0. \tag{2.49}$$

*Proof.* We divide the proof into several steps:

**Step 1:** Direct integration by parts.

We plug the test function  $v_h = u_h$  into (2.21) to obtain

$$\begin{aligned} \gamma_E(u_h, \eta_h, u_h, u_h) + \frac{1}{2} \int_F (u_h \cdot \mathcal{N}_h) |u_h|^2 &= \int_E J_h(u_h \cdot \nabla_{\mathcal{A}_h} u_h) \cdot u_h + \frac{1}{2} \int_E J_h(\nabla_{\mathcal{A}_h} \cdot u_h)(u_h \cdot u_h) \\ &- \frac{1}{2} \int_{\partial E \setminus \partial \Omega} [u_h \cdot J_h \mathcal{A}_h] \cdot n_E \frac{u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)}}{2} + \int_{\partial E^-} |\{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E| \left( u_h^{\text{int}(E)} - u_h^{\text{ext}(E)} \right) \cdot u_h^{\text{int}(E)}. \end{aligned} \tag{2.50}$$

A direct integration by parts yields

$$\begin{aligned} \int_E J_h (u_h \cdot \nabla_{\mathcal{A}_h} u_h) \cdot u_h &= - \int_E J_h (u_h \cdot \nabla_{\mathcal{A}_h} u_h) \cdot u_h - \int_E J_h (\nabla_{\mathcal{A}_h} \cdot u_h) (u_h \cdot u_h) \\ &\quad + \int_{\partial E} \left( u_h^{\text{int}(E)} \cdot J_h^{\text{int}(E)} \mathcal{A}_h^{\text{int}(E)} \cdot n_E \right) \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right). \end{aligned} \tag{2.51}$$

Therefore, (2.50) can be simplified as

$$\begin{aligned} &\gamma_E (u_h, \eta_h, u_h, u_h) + \frac{1}{2} \int_{\partial E \cap \Sigma} (u_h \cdot \mathcal{N}_h) |u_h|^2 \\ &= - \left( \int_E J_h (u_h \cdot \nabla_{\mathcal{A}_h} u_h) \cdot u_h + \frac{1}{2} \int_E J_h (\nabla_{\mathcal{A}_h} \cdot u_h) (u_h \cdot u_h) \right) \\ &\quad - \frac{1}{2} \int_{\partial E \setminus \partial \Omega} [u_h \cdot J_h \mathcal{A}_h] \cdot n_E \frac{u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)}}{2} \\ &\quad + \int_{\partial E^-} |\{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E| \left( u_h^{\text{int}(E)} - u_h^{\text{ext}(E)} \right) \cdot u_h^{\text{int}(E)} \\ &\quad + \int_{\partial E} \left( u_h^{\text{int}(E)} \cdot J_h^{\text{int}(E)} \mathcal{A}_h^{\text{int}(E)} \cdot n_E \right) \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) \\ &= I + II + III + IV, \end{aligned} \tag{2.52}$$

where

$$I = - \left( \int_E J_h (u_h \cdot \nabla_{\mathcal{A}_h} u_h) \cdot u_h + \frac{1}{2} \int_E J_h (\nabla_{\mathcal{A}_h} \cdot u_h) (u_h \cdot u_h) \right), \tag{2.53}$$

and *II*, *III* and *IV* can be understood respectively.

**Step 2:** Estimates of *II* and *IV*.

In *IV*, for  $e \in \partial E \setminus \partial \Omega$ , we have the decomposition

$$u_h^{\text{int}(E)} \cdot J_h^{\text{int}(E)} \mathcal{A}_h^{\text{int}(E)} = \{u_h \cdot J_h \mathcal{A}_h\} + \frac{1}{2} \left( u_h^{\text{int}(E)} \cdot J_h^{\text{int}(E)} \mathcal{A}_h^{\text{int}(E)} - u_h^{\text{ext}(E)} \cdot J_h^{\text{ext}(E)} \mathcal{A}_h^{\text{ext}(E)} \right). \tag{2.54}$$

For  $e \subset \partial E \setminus \partial \Omega$ , we can sum up the second term on the right-hand side of (2.54) over  $E$  and its neighboring cells, which have the opposition outward normal vectors, to show

$$\begin{aligned} &\int_e \frac{1}{2} \left( u_h^{\text{int}(E)} \cdot J_h^{\text{int}(E)} \mathcal{A}_h^{\text{int}(E)} - u_h^{\text{ext}(E)} \cdot J_h^{\text{ext}(E)} \mathcal{A}_h^{\text{ext}(E)} \right) \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) \\ &\quad + \int_e \frac{1}{2} \left( u_h^{\text{ext}(E)} \cdot J_h^{\text{ext}(E)} \mathcal{A}_h^{\text{ext}(E)} - u_h^{\text{int}(E)} \cdot J_h^{\text{int}(E)} \mathcal{A}_h^{\text{int}(E)} \right) \cdot (-n_E) \left( u_h^{\text{ext}(E)} \cdot u_h^{\text{ext}(E)} \right) \\ &= \int_e [u_h \cdot J_h \mathcal{A}_h] \cdot n_e \{u_h \cdot u_h\}, \end{aligned} \tag{2.55}$$

where  $n_e$  denotes the normal vector read from the same reference cell as  $[u_h \cdot J_h \mathcal{A}_h]$ . Note for  $e \subset \partial E \cap \partial \Omega$ ,  $u_h^{\text{int}(E)} \cdot J_h^{\text{int}(E)} \mathcal{A}_h^{\text{int}(E)} = \{u_h \cdot J_h \mathcal{A}_h\}$ . Therefore, we obtain

$$\sum_{E \in \mathcal{E}_h} IV = \sum_{e \in \partial \mathcal{E}_h} \int_e \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) + \sum_{e \in \partial \mathcal{E}_h \setminus \partial \Omega} \int_e [u_h \cdot J_h \mathcal{A}_h] \cdot n_e \{u_h \cdot u_h\}. \tag{2.56}$$

Also, in  $II$ , for  $e \subset \partial E \setminus \partial\Omega$ , we can sum up over  $E$  and its neighboring cells to achieve

$$\begin{aligned}
 -\frac{1}{2} \int_e [u_h \cdot J_h \mathcal{A}_h] \cdot n_E \frac{u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)}}{2} - \frac{1}{2} \int_e (-[u_h \cdot J_h \mathcal{A}_h]) \cdot (-n_E) \frac{u_h^{\text{ext}(E)} \cdot u_h^{\text{ext}(E)}}{2} \\
 = -\frac{1}{2} \int_e [u_h \cdot J_h \mathcal{A}_h] \cdot n_e \{u_h \cdot u_h\}, \tag{2.57}
 \end{aligned}$$

which further leads to

$$\sum_{E \in \mathcal{E}_h} II = - \sum_{e \in \partial \mathcal{E}_h \setminus \partial\Omega} \frac{1}{2} \int_e [u_h \cdot J_h \mathcal{A}_h] \cdot n_e \{u_h \cdot u_h\}. \tag{2.58}$$

**Step 3:** Further estimates in (2.52).

Hence, summing over all  $E \in \mathcal{E}_h$  in (2.52) and combining (2.56) and (2.58), we deduce

$$\begin{aligned}
 & \sum_{E \in \mathcal{E}_h} \gamma_E(u_h, \eta_h, u_h, u_h) + \frac{1}{2} \int_{\Sigma} (u_h \cdot \mathcal{N}_h) |u_h|^2 \\
 &= - \left( \int_{\Omega} J_h(u_h \cdot \nabla_{\mathcal{A}_h} u_h) \cdot u_h + \frac{1}{2} \int_{\Omega} J_h(\nabla_{\mathcal{A}_h} \cdot u_h)(u_h \cdot u_h) - \sum_{e \in \partial \mathcal{E}_h \setminus \partial\Omega} \frac{1}{2} \int_e [u_h \cdot J_h \mathcal{A}_h] \cdot n_E \{u_h \cdot u_h\} \right) \\
 &+ \sum_{E \in \mathcal{E}_h} \left( \int_{\partial E^-} |\{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E| \left( u_h^{\text{int}(E)} - u_h^{\text{ext}(E)} \right) \cdot u_h^{\text{int}(E)} + \int_{\partial E} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) \right). \tag{2.59}
 \end{aligned}$$

We can further estimate the last two terms in (2.59) as follows:

$$\begin{aligned}
 & \sum_{E \in \mathcal{E}_h} \left( \int_{\partial E^-} |\{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E| \left( u_h^{\text{int}(E)} - u_h^{\text{ext}(E)} \right) \cdot u_h^{\text{int}(E)} + \int_{\partial E} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) \right) \\
 &= \sum_{E \in \mathcal{E}_h} \int_{\partial E^-} |\{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E| u_h^{\text{ext}(E)} \cdot \left( u_h^{\text{ext}(E)} - u_h^{\text{int}(E)} \right) \\
 &+ \sum_{E \in \mathcal{E}_h} \int_{\partial E^+ \cap \partial\Omega} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) + Z \sum_{E \in \mathcal{E}_h} \int_{\partial E^- \cap \Sigma} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right). \tag{2.60}
 \end{aligned}$$

Then combining (2.59) and (2.60), we have the complete form

$$\begin{aligned}
 & \sum_{E \in \mathcal{E}_h} \gamma_E(u_h, \eta_h, u_h, u_h) + \frac{1}{2} \int_{\Sigma} (u_h \cdot \mathcal{N}_h) |u_h|^2 \\
 &= - \left( \int_{\Omega} J_h(u_h \cdot \nabla_{\mathcal{A}_h} u_h) \cdot u_h + \frac{1}{2} \int_{\Omega} J_h(\nabla_{\mathcal{A}_h} \cdot u_h)(u_h \cdot u_h) \right. \\
 &- \sum_{e \in \partial \mathcal{E}_h} \frac{1}{2} \int_{e \not\subset \partial\Omega} [u_h \cdot J_h \mathcal{A}_h] \cdot n_E \{u_h \cdot u_h\} \left. \right) + \sum_{E \in \mathcal{E}_h} \int_{\partial E^-} |\{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E| u_h^{\text{ext}(E)} \cdot \left( u_h^{\text{ext}(E)} - u_h^{\text{int}(E)} \right) \\
 &+ \sum_{E \in \mathcal{E}_h} \int_{\partial E^+ \cap \partial\Omega} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) + \sum_{E \in \mathcal{E}_h} \int_{\partial E^- \cap \Sigma} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right). \tag{2.61}
 \end{aligned}$$

**Step 4:** Synthesis.

Summing up (2.50) over  $E \in \mathcal{E}_h$  and adding it to (2.61) imply

$$\begin{aligned}
 & \sum_{E \in \mathcal{E}_h} \gamma_E(u_h, \eta_h, u_h, u_h) + \frac{1}{2} \int_{\Sigma} (u_h \cdot \mathcal{N}_h) |u_h|^2 \\
 = & \frac{1}{2} \sum_{E \in \mathcal{E}_h} \int_{\partial E^-} |\{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E| \left| u_h^{\text{ext}(E)} - u_h^{\text{int}(E)} \right|^2 + \sum_{E \in \mathcal{E}_h} \frac{1}{2} \int_{\partial E^+ \cap \partial \Omega} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) \\
 & + \sum_{E \in \mathcal{E}_h} \frac{1}{2} \int_{\partial E^- \cap \Sigma} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) \\
 \geq & \sum_{E \in \mathcal{E}_h} \frac{1}{2} \int_{\partial E^+ \cap \Sigma} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) + \sum_{E \in \mathcal{E}_h} \frac{1}{2} \int_{\partial E^- \cap \Sigma} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) \\
 = & \frac{1}{2} \int_{\Sigma} \{u_h \cdot J_h \mathcal{A}_h\} \cdot n_E \left( u_h^{\text{int}(E)} \cdot u_h^{\text{int}(E)} \right) \\
 = & \frac{1}{2} \int_{\Sigma} (u_h \cdot \mathcal{N}_h) |u_h|^2, \tag{2.62}
 \end{aligned}$$

where the last equality can be directly verified by the definitions of  $J_h$ ,  $\mathcal{A}_h$  and  $\mathcal{N}_h$ . Then our result naturally follows.  $\square$

2.3.3. Estimate of diffusion terms

In order to show the coercivity of the diffusion term, we need the discrete form of Korn’s inequality. The proof here is based on [18].

**Lemma 2.5.** Assume  $f_h \in X_h^k$  satisfies for  $e \subset \partial \mathcal{E}_h \cap \Sigma_b$ , it holds that  $f_h^{\text{ext}}|_e = 0$ . Then for sufficiently large  $\sigma_0 > 0$ , we have

$$\sum_{E \in \mathcal{E}_h} \int_E |\mathbb{D}f_h|^2 + \sum_{e \in \partial \mathcal{E}_h} \frac{\sigma_0}{h} \int_e [f_h]^2 \gtrsim \|f_h\|_{X_h}^2 \gtrsim \|f_h\|_{H_h}^2. \tag{2.63}$$

*Proof.* The second inequality has been shown in [7], so we turn to the first one. It is easy to see the key part is to show the derivatives of  $f_h$  can be controlled. Hence, we only need to show

$$\sum_{E \in \mathcal{E}_h} \int_E |\mathbb{D}f_h|^2 + \sum_{e \in \partial \mathcal{E}_h} \frac{\sigma_0}{h} \int_e [f_h]^2 \gtrsim \|f_h\|_{H_h}^2. \tag{2.64}$$

If this is not true, then we can construct a sequence  $f_h^n \in X_h^k$  satisfying

$$\|f_h^n\|_{H_h} = 1, \tag{2.65}$$

and

$$\sum_{E \in \mathcal{E}_h} \|\mathbb{D}f_h^n\|_{L^2(E)} + \sum_{e \in \partial \mathcal{E}_h} \frac{\sigma_0}{h} \int_e [f_h^n]^2 \leq \frac{1}{n}. \tag{2.66}$$

To abuse the notations, we can extract the weakly convergent subsequence

$$f_h^n \rightharpoonup f_h \text{ in } H^1(E) \text{ for } \forall E \in \mathcal{E}_h. \tag{2.67}$$

By the compact embedding theorem in each cell  $E$  and the weak lower semi-continuity of  $H^1(E)$  norm, we have the strongly convergent subsequence

$$f_h^n \rightarrow f_h \text{ in } L^2(\Omega). \tag{2.68}$$

By (2.66), we also have

$$\sum_{E \in \mathcal{E}_h} \|\mathbb{D}f_h^n\|_{L^2(\Omega)} \rightarrow 0. \tag{2.69}$$

Notice the fact that in each cell  $E$ , we still have the continuous Korn’s inequality

$$\|f_h^m - f_h^n\|_{H^1(E)} \lesssim \|\mathbb{D}f_h^m - \mathbb{D}f_h^n\|_{L^2(E)} + \|f_h^m - f_h^n\|_{L^2(E)} \text{ for } \forall m, n \in \mathbb{N}. \tag{2.70}$$

Hence, combining all above, we know  $f_h^n$  is a Cauchy’s sequence under the norm  $H^1(E)$ , which means it is a Cauchy’s sequence under the norm  $H_h$ . Then

$$f_h^n \rightarrow f_h \text{ in } H_h. \tag{2.71}$$

Thus this means

$$\|f_h\|_{H_h} = 1, \tag{2.72}$$

which further implies

$$\mathbb{D}f_h = 0 \text{ in } \forall E \in \mathcal{E}_h, \tag{2.73}$$

and

$$\sum_{e \in \partial \mathcal{E}_h} \frac{\sigma_0}{h} \int_e [f_h]^2 = 0. \tag{2.74}$$

This means

$$f_h = (a, b) + c(-x_1, x_2), \tag{2.75}$$

for some constant  $a, b, c$  and  $f_h$  is continuous in  $\bar{\Omega}$ . Certainly, the zero bottom implies  $f_h = 0$ , which contradicts  $\|f_h\|_{H_h} = 1$ . Therefore, our hypothesis is invalid and (2.64) holds.  $\square$

**Lemma 2.6.** *Assume  $f_h$  satisfies for  $e \subset \partial \mathcal{E}_h \cap \Sigma_b$ , it holds that  $f_h^{\text{ext}}|_e = 0$ . Also,  $\eta_h$  satisfies  $\eta_h + 1 \geq \delta > 0$  and*

$$Q = \sup_{F \in \mathcal{F}_h} \|\eta_h\|_{W^{1,\infty}(F)} < \infty. \tag{2.76}$$

Then for sufficiently large  $\sigma_0 > 0$ , we have

$$\sum_{E \in \mathcal{E}_h} \int_E J_h |\mathbb{D}_{\mathcal{A}_h} f_h|^2 + \sum_{e \in \partial \mathcal{E}_h} \frac{\sigma_0}{h} \int_e [f_h]^2 \gtrsim C(Q) \|f_h\|_{X_h}^2 \gtrsim C(Q) \|f_h\|_{H_h}^2. \tag{2.77}$$

*Proof.* Note that in each cell  $E$ , the free surface  $\eta_h$  is smooth. Hence, we can change it back by the transform  $\Phi_h^{-1}$  to a curved cell  $\Phi_h^{-1}E$  which satisfies

$$\int_E J_h(x) |\mathbb{D}_{\mathcal{A}_h(x)} f_h(x)|^2 dx = \int_{\Phi_h^{-1}E} |\mathbb{D}f_h(y)|^2 dy. \tag{2.78}$$

Hence, in this curved cell, we have the Korn's inequality

$$\|f_h\|_{H^1(\Phi_h^{-1}E)} \lesssim \|\mathbb{D}f_h\|_{L^2(\Phi_h^{-1}E)} + \|f_h\|_{L^2(\Phi_h^{-1}E)}. \quad (2.79)$$

Since the transform  $\Phi_h^{-1}$  is a diffeomorphism between  $E$  and  $\Phi_h^{-1}(E)$ , then the  $H^1$  norms in these two spaces are comparable. Hence, we have

$$\|f_h\|_{H^1(E)} \lesssim C(Q) \left( \int_E J_h |\mathbb{D}_{\mathcal{A}_h} f_h|^2 + \|f_h\|_{L^2(E)} \right). \quad (2.80)$$

Then a similar proof as that of Lemma 2.5 naturally yields the desired result.  $\square$

**Lemma 2.7.** *Assume  $f_h$  satisfies for  $e \subset \partial\mathcal{E}_h \cap \Sigma_b$ , it holds that  $f_h^{\text{ext}}|_e = 0$ . Also,  $\eta_h$  satisfies  $\eta_h + 1 \geq \delta > 0$  and*

$$Q = \sup_{F \in \mathcal{F}_h} \|\eta_h\|_{W^{1,\infty}(F)} < \infty. \quad (2.81)$$

*Then there exists a sufficiently small constant  $\delta' > 0$  such that for  $\eta'$  satisfying*

$$\sup_{F \in \mathcal{F}_h} \|\eta'_h - \eta_h\|_{W^{1,\infty}(F)} \leq \delta', \quad (2.82)$$

*and for sufficiently large  $\sigma'_0 > 0$ , we have*

$$\sum_{E \in \mathcal{E}_h} \int_E J'_h |\mathbb{D}_{\mathcal{A}'_h} f_h|^2 + \sum_{e \in \partial\mathcal{E}_h} \frac{\sigma'_0}{h} \int_e [f_h]^2 \gtrsim C(Q) \|f_h\|_{X_h}^2 \gtrsim C(Q) \|f_h\|_{H_h}^2. \quad (2.83)$$

*Proof.* By Lemma 2.6, we know

$$\sum_{E \in \mathcal{E}_h} \int_E J_h |\mathbb{D}_{\mathcal{A}_h} f_h|^2 + \sum_{e \in \partial\mathcal{E}_h} \frac{\sigma_0}{h} \int_e [f_h]^2 \gtrsim C(Q) \|f_h\|_{X_h}^2 \gtrsim C(Q) \|f_h\|_{H_h}^2. \quad (2.84)$$

We can rewrite our formula in a perturbed form as

$$\begin{aligned} & \sum_{E \in \mathcal{E}_h} \int_E J'_h |\mathbb{D}_{\mathcal{A}'_h} f_h|^2 - \sum_{E \in \mathcal{E}_h} \int_E J_h |\mathbb{D}_{\mathcal{A}_h} f_h|^2 \\ &= \sum_{E \in \mathcal{E}_h} \int_E J'_h \left( \mathbb{D}_{\mathcal{A}'_h - \mathcal{A}_h} f_h \right) \left( 2\mathbb{D}_{\mathcal{A}_h} f_h + \mathbb{D}_{\mathcal{A}'_h - \mathcal{A}_h} f_h \right) + \sum_{E \in \mathcal{E}_h} \int_E (J'_h - J_h) |\mathbb{D}_{\mathcal{A}_h} f_h|^2 \\ &\lesssim \delta' \|f_h\|_{H_h}^2. \end{aligned} \quad (2.85)$$

Naturally, when taking  $\delta'$  sufficiently small, we can absorb the perturbation into the principle part. Then our result easily follows.  $\square$

For the diffusion term

$$- \int_E J \Delta_{\mathcal{A}} u \cdot v_h, \quad (2.86)$$

a direct integration by parts implies

$$- \sum_{E \in \mathcal{E}_h} \int_E J \Delta_{\mathcal{A}} u \cdot v_h = \frac{1}{2} \int_{\Omega} J \mathbb{D}_{\mathcal{A}} u : \mathbb{D}_{\mathcal{A}} v_h - \sum_{e \in \partial\mathcal{E}_h \setminus \Sigma} \int_e \{ \mathbb{D}_{\mathcal{A}} u J_{\mathcal{A}} \cdot n_E \} \cdot [v_h] - \int_{\Sigma} (\mathbb{D}_{\mathcal{A}} u) \mathcal{N} \cdot v_h^{\text{int}(E)}. \quad (2.87)$$

Also in the continuous case, our discretization yields

$$\sum_{E \in \mathcal{E}_h} \alpha_E (\eta, u, \eta, v_h) = \frac{1}{2} \int_{\Omega} J \mathbb{D}_{\mathcal{A}} u : \mathbb{D}_{\mathcal{A}} v_h - \sum_{e \in \partial \mathcal{E}_h \setminus \Sigma} \int_e \{ \mathbb{D}_{\mathcal{A}} u J_{\mathcal{A}} \cdot n_E \} \cdot [v_h]. \tag{2.88}$$

We can notice the difference is the extra physical boundary term  $-\int_{\Sigma} (\mathbb{D}_{\mathcal{A}} u) \cdot \mathcal{N} \cdot v_h^{\text{int}(E)}$ , which contributes to the boundary condition on  $\Sigma$ . This is separately added to the scheme as a forcing term later combined with the pressure contribution. Hence, our scheme is consistent.

**Lemma 2.8.** *Assume  $\eta_h$  satisfies  $\eta_h + 1 \geq \delta > 0$ . If we take  $v_h = u_h$  and choose the penalty constant  $\sigma$  properly, the discretized diffusion term with NIPG satisfies*

$$\sum_{E \in \mathcal{E}_h} \alpha_E (\eta_h, u_h, \eta_h, u_h) \geq \int_{\Omega} J_h |\mathbb{D}_{\mathcal{A}_h} u_h|^2. \tag{2.89}$$

If we further assume  $\eta_h$  satisfies

$$Q = \sup_{F \in \mathcal{F}_h} \|\eta_h\|_{W^{1,\infty}(F)} < \infty, \tag{2.90}$$

then the discretized diffusion term with SIPG or NIPG satisfies

$$\sum_{E \in \mathcal{E}_h} \alpha_E (\eta_h, u_h, \eta_h, u_h) \gtrsim C(Q) \|u_h\|_{X_h}^2. \tag{2.91}$$

*Proof.* We can directly compute

$$\begin{aligned} \sum_{E \in \mathcal{E}_h} \alpha_E (\eta_h, u_h, \eta_h, v_h) &= \frac{1}{2} \sum_{E \in \mathcal{E}_h} \int_E J_h \mathbb{D}_{\mathcal{A}_h} u_h : \mathbb{D}_{\mathcal{A}_h} v_h - \sum_{e \in \partial \mathcal{E}_h \setminus \Sigma} \int_e \{ \mathbb{D}_{\mathcal{A}_h} u_h J_h \mathcal{A}_h \cdot n_E \} \cdot [v_h] \\ &\pm \sum_{e \in \partial \mathcal{E}_h \setminus \Sigma} \int_e \{ \mathbb{D}_{\mathcal{A}_h} v_h J_h \mathcal{A}_h \cdot n_E \} \cdot [u_h] + \frac{\sigma}{h} \sum_{e \in \partial \mathcal{E}_h} [u_h] [v_h]. \end{aligned} \tag{2.92}$$

For NIPG, when  $v_h = u_h$ , (2.92) reduces to

$$\sum_{E \in \mathcal{E}_h} \alpha_E (\eta_h, u_h, \eta_h, u_h) = \frac{1}{2} \int_{\Omega} J_h |\mathbb{D}_{\mathcal{A}_h} u_h|^2 + \frac{\sigma}{h} \sum_e [u_h]^2. \tag{2.93}$$

Hence, we may simply take the penalty  $\sigma = 1$  and by the discrete Korn’s inequality in Lemma 2.6, our result naturally follows.

For SIPG, when  $v_h = u_h$ , (2.92) reduces to

$$\sum_{E \in \mathcal{E}_h} \alpha_E (\eta_h, u_h, \eta_h, u_h) = \frac{1}{2} \int_{\Omega} J_h |\mathbb{D}_{\mathcal{A}_h} u_h|^2 - 2 \sum_{e \in \partial \mathcal{E}_h \setminus \Sigma} \int_e \{ \mathbb{D}_{\mathcal{A}_h} u_h J_h \mathcal{A}_h \cdot n_E \} \cdot [u_h] + \frac{\sigma}{h} \sum_{e \in \partial \mathcal{E}_h} [u_h]^2. \tag{2.94}$$

By the discrete Korn’s inequality in Lemma 2.6, we have

$$\sum_{E \in \mathcal{E}_h} \int_E J_h |\mathbb{D}_{\mathcal{A}_h} u_h|^2 + \sum_{e \in \partial \mathcal{E}_h} \frac{\sigma_0}{h} \int_e [u_h]^2 \gtrsim C(Q) \|u_h\|_{X_h}^2 \gtrsim C(Q) \|u_h\|_{H_h}^2. \tag{2.95}$$



Then we utilize Hölder’s inequality, the trace theorem and Cauchy’s inequality to estimate

$$\begin{aligned}
 \sum_{e \in \partial \mathcal{E}_h \setminus \partial \Omega} \int_e \{ \mathbb{D}_{\mathcal{A}_h} u_h J_h \mathcal{A}_h \cdot n_E \} \cdot [u_h] &\leq C(Q) \sum_{e \in \partial \mathcal{E}_h \setminus \partial \Omega} \|\nabla u_h\|_{L^2(e)} \|[u_h]\|_{L^2(e)} \\
 &\leq C(Q) \sum_{e \in \partial \mathcal{E}_h \setminus \partial \Omega} \|u_h\|_{H^1(e)} \|[u_h]\|_{L^2(e)} \\
 &\leq CC(Q) \sum_{e \in \partial \mathcal{E}_h \setminus \partial \Omega} h \|u_h\|_{H^1(e)}^2 + \frac{C(Q)}{4C} \frac{1}{h} \sum_{e \in \partial \mathcal{E}_h \setminus \partial \Omega} \|[u_h]\|_{L^2(e)}^2 \\
 &\leq CC(Q) \sum_{E \in \mathcal{E}_h} h \|u_h\|_{H^{3/2}(E)}^2 + \frac{C(Q)}{4C} \frac{1}{h} \sum_{e \in \partial \mathcal{E}_h \setminus \partial \Omega} \|[u_h]\|_{L^2(e)}^2 \\
 &\leq CC(Q) \|u_h\|_{H_h}^2 + \frac{C(Q)}{4C} \frac{1}{h} \sum_{e \in \partial \mathcal{E}_h \setminus \partial \Omega} \|[u_h]\|_{L^2(e)}^2. \tag{2.96}
 \end{aligned}$$

The last inequality is valid since  $u_h \in P_h^k$ . Then we take  $C$  sufficiently small, and  $\sigma \geq \sigma + C(Q)/(4C)$  to absorb (2.96) into (2.95). Hence, our result easily follows.  $\square$

2.3.4. Estimate of pressure terms

For the pressure term

$$\int_E J \nabla_{\mathcal{A}} p \cdot v_h, \tag{2.97}$$

a direct integration by parts reveals the following equality in the continuous case

$$\sum_{E \in \mathcal{E}_h} \int_E J \nabla_{\mathcal{A}} p \cdot v_h = - \int_{\Omega} J p \nabla_{\mathcal{A}} \cdot v_h + \sum_{e \in \partial \mathcal{E}_h \setminus \Sigma} \int_e [v_h] \cdot (p J \mathcal{A}) \cdot n_e + \int_{\Sigma} p \mathcal{N} \cdot v_h^{\text{int}(E)}. \tag{2.98}$$

It is easy to see in the continuous case, our discretization reduces to

$$\sum_{E \in \mathcal{E}_h} \beta_E(p, \eta, v_h) = - \int_{\Omega} J p \nabla_{\mathcal{A}} \cdot v_h + \sum_{e \in \partial \mathcal{E}_h \setminus \Sigma} \int_e [v_h] \cdot (p J \mathcal{A}) \cdot n_e. \tag{2.99}$$

We can notice the difference is the extra physical boundary term  $\int_{\Sigma} p \mathcal{N} \cdot v^{\text{int}(E)}$ , which contributes to the boundary condition on  $\Sigma$ . This is separately added to the scheme as a forcing term later combined with the diffusion contribution. Hence, our scheme is consistent.

**Lemma 2.9.** *When we take  $v_h = u_h$ , the discretized pressure term satisfies*

$$\sum_{E \in \mathcal{E}_h} \beta_E(p_h, \eta_h, u_h) = - \int_{\Omega} J_h p_h \nabla_{\mathcal{A}_h} \cdot u_h + \sum_{e \in \partial \mathcal{E}_h \setminus \Sigma} \int_e [u_h] \cdot \{ p_h J_h \mathcal{A}_h \} \cdot n_e, \tag{2.100}$$

where  $n_e$  denotes the outward normal vector read from the same direction as  $[u_h]$ .

*Proof.* This is a direct corollary of the definition, so we omit the proof here.  $\square$

2.3.5. Estimate of forcing terms

Considering the physical boundary condition on  $\Sigma$  of the system (1.9), we need to take the upper boundary condition as a forcing term, *i.e.*

$$(pI - \mathbb{D}_{\mathcal{A}}u) \mathcal{N} = \eta \mathcal{N} \quad \text{on } \Sigma. \tag{2.101}$$

Multiplying the test function  $v_h$  and integrating over  $F$  yield

$$\int_F \eta \mathcal{N} v_h = \int_F \eta (\bar{v}_2)_h - \int_F \eta \partial_1 \eta (\bar{v}_1)_h. \tag{2.102}$$

In the continuous case, our discretization reduces to

$$\sum_{E \in \mathcal{E}_h} \mu_E(\eta, \eta, v_h) = \int_{\Sigma} \eta (\bar{v}_2)_h - \int_{\Sigma} \eta \partial_1 \eta (\bar{v}_1)_h, \tag{2.103}$$

and the penalty terms vanish. Hence, this discretization is consistent.

**Lemma 2.10.** *When we take  $v_h = u_h$ , the discretized forcing term satisfies*

$$\sum_{E \in \mathcal{E}_h} \mu_E(\eta_h, \eta_h, u_h) \geq \frac{1}{2} \int_{\Sigma} \partial_t |\eta_h|^2. \tag{2.104}$$

*Proof.* We may sum up over  $E \in \mathcal{E}_h$  to obtain

$$\begin{aligned} & \sum_{E \in \mathcal{E}_h} \mu_E(\eta_h, \eta_h, v_h) \\ &= \int_{\Sigma} \eta_h (\bar{v}_2)_h + \int_{\Sigma} \eta_h \partial_1 \eta_h (\bar{v}_1)_h + \frac{1}{2} \sum_{F \in \mathcal{F}_h} \left( \{(\bar{u}_1)_h\} (\eta_h^-)^2 \Big|_{F+1/2} - \{(\bar{u}_1)_h\} (\eta_h^+)^2 \Big|_{F-1/2} \right). \end{aligned} \tag{2.105}$$

Taking  $\phi_h = \eta_h$  in the transport discretization (2.15) and integrating by parts imply

$$\begin{aligned} & \frac{1}{2} \int_F \partial_t |\eta_h|^2 - \int_F (\bar{u}_2)_h \eta_h + \int_F (\bar{u}_1)_h \eta_h \partial_1 \eta_h - (\bar{u}_1)_h^- (\eta_h^-)^2 \Big|_{F+1/2} + (\bar{u}_1)_h^+ (\eta_h^+)^2 \Big|_{F-1/2} \\ & + (\hat{u}_1)_h \hat{\eta}_h \eta_h^- \Big|_{F+1/2} - (\hat{u}_1)_h \hat{\eta}_h \eta_h^+ \Big|_{F-1/2} + \frac{1}{2} \left( (\eta_h^-)^2 [u_h] \Big|_{F+1/2} - (\eta_h^+)^2 [u_h] \Big|_{F-1/2} \right) = 0, \end{aligned} \tag{2.106}$$

which further leads to

$$\begin{aligned} & \frac{1}{2} \int_F \partial_t |\eta_h|^2 - \int_F (\bar{u}_2)_h \eta_h + \int_F (\bar{u}_1)_h \eta_h \partial_1 \eta_h \\ & - \{(\bar{u}_1)_h\} (\eta_h^-)^2 \Big|_{F+1/2} + \{(\bar{u}_1)_h\} (\eta_h^+)^2 \Big|_{F-1/2} + (\hat{u}_1)_h \hat{\eta}_h \eta_h^- \Big|_{F+1/2} - (\hat{u}_1)_h \hat{\eta}_h \eta_h^+ \Big|_{F-1/2} = 0. \end{aligned} \tag{2.107}$$

Summing over  $F \in \mathcal{F}_h$  implies

$$\begin{aligned} & \frac{1}{2} \int_{\Sigma} \partial_t |\eta_h|^2 - \int_{\Sigma} (\bar{u}_2)_h \eta_h + \int_{\Sigma} (\bar{u}_1)_h \eta_h \partial_1 \eta_h - \sum_{F \in \mathcal{F}_h} \{(\bar{u}_1)_h\} (\eta_h^-)^2 \Big|_{F+1/2} \\ & + \sum_{F \in \mathcal{F}_h} \{(\bar{u}_1)_h\} (\eta_h^+)^2 \Big|_{F-1/2} - \sum_{F \in \mathcal{F}_h} (\hat{u}_1)_h \hat{\eta}_h \eta_h^- \Big|_{F+1/2} + \sum_{F \in \mathcal{F}_h} (\hat{u}_1)_h \hat{\eta}_h \eta_h^+ \Big|_{F-1/2} = 0. \end{aligned} \tag{2.108}$$

Since we always take  $\hat{u}_1 = \{(\bar{u}_1)_h\}$ , combining (2.105) and (2.108) yields

$$\begin{aligned} \sum_{E \in \mathcal{E}_h} \mu_E(\eta_h, \eta_h, v_h) &= \frac{1}{2} \int_{\Sigma} \partial_t |\eta_h|^2 + \sum_{F \in \mathcal{F}_h} \left( -\frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^-)^2 |_{F+1/2} + \frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^+)^2 |_{F-1/2} \right. \\ &\quad \left. + \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^- |_{F+1/2} - \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^+ |_{F-1/2} \right). \end{aligned} \quad (2.109)$$

We can decompose the boundary terms in (2.109) to obtain

$$\begin{aligned} -\frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^-)^2 |_{F+1/2} + \frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^+)^2 |_{F-1/2} + \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^- |_{F+1/2} - \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^+ |_{F-1/2} \\ = \Psi_{F+1/2} - \Psi_{F-1/2} + \Theta_F, \end{aligned} \quad (2.110)$$

where

$$\Psi_{F+1/2} = -\frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^-)^2 |_{F+1/2} + \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^- |_{F+1/2}, \quad (2.111)$$

$$\Psi_{F-1/2} = -\frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^-)^2 |_{F-1/2} + \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^- |_{F-1/2}, \quad (2.112)$$

and

$$\Theta_F = \frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^+)^2 |_{F-1/2} - \frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^-)^2 |_{F-1/2} - \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^+ |_{F-1/2} + \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^- |_{F-1/2}. \quad (2.113)$$

By the flux definition (2.18), we can derive

$$\Theta_F = \{(\bar{u}_1)_h\} (\eta_h^+ - \eta_h^-) \left( \frac{\eta_h^+ + \eta_h^-}{2} - \hat{\eta}_h \right) \geq 0. \quad (2.114)$$

When we sum up (2.110) over  $F \in \mathcal{F}_h$ , all  $\Psi_{F+1/2}$  are canceled out due to the periodicity. Then we have

$$\begin{aligned} \sum_{F \in \mathcal{F}_h} \left( -\frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^-)^2 |_{F+1/2} + \frac{1}{2} \{(\bar{u}_1)_h\} (\eta_h^+)^2 |_{F-1/2} \right. \\ \left. + \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^- |_{F+1/2} - \{(\bar{u}_1)_h\} \hat{\eta}_h \eta_h^+ |_{F-1/2} \right) \geq 0. \end{aligned} \quad (2.115)$$

Hence, combining (2.109) and (2.115), we can obtain the desired result.  $\square$

### 2.3.6. Estimate of divergence terms

In the continuous case, our discretization reduces to

$$\sum_{E \in \mathcal{E}_h} \rho_E(u, \eta, q_h) = - \int_{\Omega} J(\nabla_{\mathcal{A}} \cdot u) q_h. \quad (2.116)$$

Hence, this discretization is consistent.

**Lemma 2.11.** *When we take  $q_h = p_h$ , the discretized divergence term satisfies*

$$\sum_{E \in \mathcal{E}_h} \rho_E(u_h, \eta_h, p_h) = - \int_{\Omega} J_h p_h \nabla_{\mathcal{A}_h} \cdot u_h + \sum_{e \in \partial \mathcal{E}_h \setminus \Sigma} \int_e [u_h] \cdot \{p_h J_h \mathcal{A}_h\} \cdot n_e. \quad (2.117)$$

*Proof.* This is a natural corollary of the definition, so we omit the proof here.  $\square$

### 3. STABILITY ANALYSIS

**Condition 3.1.** The free surface  $\eta_h(t)$  satisfies  $1 + \eta_h(t) \geq \delta > 0$  for some  $\delta > 0$  independent of  $h$  and  $t$ .

**Condition 3.2.** The free surface  $\eta_h(t)$  satisfies  $\sup_{F \in \mathcal{F}_h} \|\eta_h(t)\|_{W^{1,\infty}(F)} \leq Q$  for some  $Q > 0$  independent of  $h$  and  $t$ .

**Theorem 3.3.** If  $S = 0$ , suppose Condition 3.1 is valid in  $t \in [0, T]$  for some  $T > 0$ . Then there exists a unique numerical solution triple  $(u_h, p_h, \eta_h) \in X_h^k \times M_h^{k-1} \times S_h^k$  to the scheme (2.14) with NIPG, which satisfies the estimate

$$\begin{aligned} \left\| \sqrt{J_h(t)} u_h(t) \right\|_{L^2(\Omega)}^2 + g \|\eta_h(t)\|_{L^2(\Sigma)}^2 + \nu \int_0^t \int_{\Omega} J_h(s) |\mathbb{D}_{\mathcal{A}_h(s)} u_h(s)|^2 ds \\ \leq \left\| \sqrt{J_h(0)} u_h(0) \right\|_{L^2(\Omega)}^2 + g \|\eta_h(0)\|_{L^2(\Sigma)}^2, \end{aligned} \tag{3.1}$$

for any  $t \in [0, T]$ .

For general  $S$ , suppose Conditions 3.1 and 3.2 are valid in  $t \in [0, T]$  for some  $T > 0$ . Then there exists a unique numerical solution triple  $(u_h, p_h, \eta_h) \in X_h^k \times M_h^{k-1} \times S_h^k$  to the scheme (2.14) with SIPG or NIPG, which satisfies the estimate

$$\begin{aligned} \left\| \sqrt{J_h(t)} u_h(t) \right\|_{L^2(\Omega)}^2 + \|\eta_h(t)\|_{L^2(\Sigma)}^2 + \int_0^t \|u_h(s)\|_{X_h}^2 ds \\ \lesssim C(Q) \left( \left\| \sqrt{J_h(0)} u_h(0) \right\|_{L^2(\Omega)}^2 + \|\eta_h(0)\|_{L^2(\Sigma)}^2 + \int_0^t \int_{\Omega} |S(r)|^2 dr \right), \end{aligned} \tag{3.2}$$

for any  $t \in [0, T]$ .

*Proof.* The existence and uniqueness follow from a standard argument for the differential-algebraic equations, so we omit it here and focus on the energy estimate. In the system (2.14), we take the test function  $v_h = u_h$  and  $q_h = p_h$ , and sum up over  $E \in \mathcal{E}_h$ . Then it yields

$$\begin{aligned} \sum_{E \in \mathcal{E}_h} \zeta_E(\eta_h, u_h, u_h) + \sum_{E \in \mathcal{E}_h} \gamma_E(u_h, \eta_h, u_h, u_h) \\ + \nu \sum_{E \in \mathcal{E}_h} \alpha_E(\eta_h, u_h, \eta_h, u_h) + \sum_{E \in \mathcal{E}_h} \beta_E(p_h, \eta_h, u_h) + g \sum_{E \in \mathcal{E}_h} \mu_E(\eta_h, \eta_h, u_h) = \sum_{E \in \mathcal{E}_h} \omega_E(\eta_h, u_h), \end{aligned} \tag{3.3}$$

$$\sum_{E \in \mathcal{E}_h} \rho_E(u_h, \eta_h, p_h) = 0. \tag{3.4}$$

By Lemmas 2.9 and 2.11, we have

$$\sum_{E \in \mathcal{E}_h} \beta_E(p_h, \eta_h, u_h) = \sum_{E \in \mathcal{E}_h} \rho_E(u_h, \eta_h, p_h) = 0. \tag{3.5}$$

Hence, we can simplify (3.3) into

$$\sum_{E \in \mathcal{E}_h} \zeta_E(\eta_h, u_h, u_h) + \sum_{E \in \mathcal{E}_h} \gamma_E(u_h, \eta_h, u_h, u_h) + \nu \sum_{E \in \mathcal{E}_h} \alpha_E(\eta_h, u_h, \eta_h, u_h) + g \sum_{E \in \mathcal{E}_h} \mu_E(\eta_h, \eta_h, u_h) = \sum_{E \in \mathcal{E}_h} \omega_E(\eta_h, u_h). \tag{3.6}$$

Then by Lemmas 2.2, 2.4, 2.8 and 2.10, we have

$$\frac{1}{2} \partial_t \int_{\Omega} J_h |u_h|^2 + \frac{g}{2} \partial_t \int_{\Sigma} |\eta_h|^2 + CC(Q) \|u_h\|_{X_h}^2 \leq \sum_{E \in \mathcal{E}_h} \omega_E(u_h). \tag{3.7}$$

Note Conditions 3.1 and 3.2 guarantee the existence of  $C(Q) > 0$  which is independent of  $t$ . An application of Cauchy’s inequality implies

$$\sum_{E \in \mathcal{E}_h} \omega_E(u_h) \leq C' \int_{\Omega} J_h^2 |u_h|^2 + \frac{1}{4C'} \int_{\Omega} |S|^2 \leq Q^2 C' \|u_h\|_{X_h}^2 + \frac{1}{4C'} \int_{\Omega} |S|^2. \tag{3.8}$$

In (3.8), when taking  $C'$  sufficiently small, we can always absorb it into  $CC(Q) \|u_h\|_{X_h}^2$  in (3.7). Hence, we have

$$\frac{1}{2} \partial_t \int_{\Omega} J_h |u_h|^2 + \frac{g}{2} \partial_t \int_{\Sigma} |\eta_h|^2 + \frac{CC(Q)}{2} \|u_h\|_{X_h}^2 \leq C(Q) \int_{\Omega} |S|^2. \tag{3.9}$$

Then integrating over  $[0, t]$  leads to the desired result. The  $S = 0$  case is easily derived from the general case without discussion of the source term.  $\square$

**Remark 3.4.** Conditions 3.1 and 3.2 are not always satisfied *a priori*. In [1, 15, 16], this type of assumptions were also introduced in the numerical analysis.

#### 4. DISCUSSION ON THE ERROR ANALYSIS

For the continuous solution  $(u, \eta)$ , the analysis in [20] reveals in the Navier–Stokes equations of (1.9), we need  $H^1$  norm of  $\eta(t)$  to bound  $L^2$  norm of  $u(t)$ . However, the result in [6] implies in the transport equation of (1.9), we need  $H^2$  norm of  $u(t)$  to control  $H^1$  norm of  $\eta(t)$ . This type of inconsistent coupling cannot be improved even if we go to higher order derivatives. Hence, this implies the coupled system in (1.9) is not closed in the usual Sobolev norms, which means we cannot expect to obtain the error estimates for the whole system (1.9).

In the following, we mainly analyze the error in the Navier–Stokes equations provided we have the error estimates of the free surface.

**Condition 4.1.** The free surface  $\eta_h$  satisfies the error estimates

$$\|\eta_h - \eta\|_{L^2} \lesssim h^{k+1}, \tag{4.1}$$

$$\|\eta_h - \eta\|_{L^\infty} \lesssim h^k, \tag{4.2}$$

$$\|\chi_h - \chi\|_{L^2} \lesssim h^k, \tag{4.3}$$

$$\|\chi_h - \chi\|_{L^\infty} \lesssim h^{k-1}, \tag{4.4}$$

$$\|\partial_t(\eta_h - \eta)\|_{L^2} \lesssim h^k, \tag{4.5}$$

$$\|\partial_t(\eta_h - \eta)\|_{L^\infty} \lesssim h^{k-1}. \tag{4.6}$$

##### 4.1. Velocity error analysis

We decompose the velocity error and the pressure error as follows:

$$u_h - u = (u_h - \mathcal{P}u) + (\mathcal{P}u - u) = \epsilon_u + \delta_u, \tag{4.7}$$

$$p_h - p = (p_h - \mathcal{P}p) + (\mathcal{P}p - p) = \epsilon_p + \delta_p, \tag{4.8}$$

where  $\mathcal{P}$  is some projection such that  $\mathcal{P}u \in X_h^k$  and  $\mathcal{P}p \in M_h^{k-1}$  achieving the optimal accuracy, *i.e.*

$$\|\delta_u\|_{L^2} \lesssim h^{k+1}, \tag{4.9}$$

$$\|\delta_p\|_{L^2} \lesssim h^k. \tag{4.10}$$

Hence, the key part of the error estimates is  $\epsilon_u$  and  $\epsilon_p$ .

**Lemma 4.2.** *Suppose Conditions 3.1, 3.2 and 4.1 are valid. Assume the exact solution satisfies  $u \in C^2(\Omega)$ ,  $p \in C^1(\Omega)$  and  $\eta \in C^1(\Sigma)$ . Then for  $h$  sufficiently small, the numerical solution to the scheme (2.14) satisfies the estimate*

$$\partial_t \int_{\Omega} J_h |\epsilon_u|^2 + \|\epsilon_u\|_{H_h}^2 \lesssim C(Q) \left( h^{k-1-r} \|\epsilon_u\|_{H_h} + h^{k-r} \|\epsilon_p\|_{L^2} + \|\epsilon_u\|_{L^2}^2 \right) \text{ for } 1/2 < r < 1. \tag{4.11}$$

*Proof.* The discretization of the Navier–Stokes equations is

$$\zeta_E(\eta_h, u_h, v_h) + \gamma_E(u_h, \eta_h, u_h, v_h) + \nu \alpha_E(\eta_h, u_h, \eta_h, v_h) + \beta_E(p_h, \eta_h, v_h) + g \mu_E(\eta_h, \eta_h, v_h) = \omega_E(\eta_h, v_h), \tag{4.12}$$

$$\rho_E(u_h, \eta_h, q_h) = 0. \tag{4.13}$$

The consistency of the scheme (2.14) implies the exact solution  $(u, p, \eta)$  also satisfies this scheme, *i.e.*

$$\zeta_E(\eta, u, v_h) + \gamma_E(u, \eta, u, v_h) + \nu \alpha_E(\eta, u, \eta, v_h) + \beta_E(p, \eta, v_h) + g \mu_E(\eta, \eta, v_h) = \omega_E(\eta, v_h), \tag{4.14}$$

$$\rho_E(u, \eta, q_h) = 0. \tag{4.15}$$

Therefore, taking the difference of above two sets of equations, we obtain the error equations

$$\begin{aligned} & (\zeta_E(\eta_h, u_h, v_h) - \zeta_E(\eta, u, v_h)) + (\gamma_E(u_h, \eta_h, u_h, v_h) - \gamma_E(u, \eta, u, v_h)) \\ & \quad + \nu (\alpha_E(\eta_h, u_h, \eta_h, v_h) - \alpha_E(\eta, u, \eta, v_h)) \\ & + (\beta_E(p_h, \eta_h, v_h) - \beta_E(p, \eta, v_h)) + g (\mu_E(\eta_h, \eta_h, v_h) - \mu_E(\eta, \eta, v_h)) = (\omega_E(\eta_h, v_h) - \omega_E(\eta, v_h)), \end{aligned} \tag{4.16}$$

$$(\rho_E(u_h, \eta_h, q_h) - \rho_E(u, \eta, q_h)) = 0. \tag{4.17}$$

Now we need to analyze each term in the error equations (4.16) and (4.17). We always decompose the difference of the bilinear forms into two parts: the energy part  $W_*$  and the remaining part  $R_*$ , such that

$$\text{Difference} = W_* + R_*, \tag{4.18}$$

where  $W_*$  builds the main body of the error equations and is put in the left-hand side (LHS), and  $R_*$  is moved to the right-hand side (RHS) and can be taken as the perturbed source term. The  $*$  can be  $\zeta$ ,  $\gamma$  or  $\alpha$ , *etc.* For example, we can decompose the convection difference  $\gamma(u_h, \eta_h, u_h, v_h) = \sum_{E \in \mathcal{E}_h} \gamma_E(u_h, \eta_h, u_h, v_h)$  as follows:

$$\begin{aligned} & \gamma(u_h, \eta_h, u_h, v_h) - \gamma(u, \eta, u, v_h) \\ & = \gamma(u_h, \eta_h, u_h - u, v_h) + \gamma(u_h, \eta_h - \eta, u, v_h) + \gamma(u_h - u, \eta, u, v_h) \\ & = \gamma(u_h, \eta_h, \epsilon_u, v_h) + \gamma(u_h, \eta_h, \delta_u, v_h) + \gamma(u_h, \eta_h - \eta, u, v_h) + \gamma(u_h - u, \eta, u, v_h) \\ & = \gamma(u_h, \eta_h, \epsilon_u, \epsilon_u) + (\gamma(u_h, \eta_h, \delta_u, \epsilon_u) + \gamma(u_h, \eta_h - \eta, u, \epsilon_u) + \gamma(u_h - u, \eta, u, \epsilon_u)) \\ & = W_\gamma + R_\gamma, \end{aligned} \tag{4.19}$$

where

$$W_\gamma = \gamma(u_h, \eta_h, \epsilon_u, \epsilon_u), \tag{4.20}$$

$$R_\gamma = \gamma(u_h, \eta_h, \delta_u, \epsilon_u) + \gamma(u_h, \eta_h - \eta, u, \epsilon_u) + \gamma(u_h - u, \eta, u, \epsilon_u). \tag{4.21}$$

$W_*$  can be bounded as in Theorem 3.3 and  $R_*$  can be estimated based on the following rules:

**Rule 1:**

For the estimates of the product terms as

$$\int F_1 F_2 \dots F_k, \tag{4.22}$$

we apply Hölder’s inequality to obtain two terms in the  $L^2$  norm and other terms in the  $L^\infty$  norm. The priority of taking the  $L^\infty$  norm is as follows:

- (1)  $F_i$  only containing the exact solution  $(u, p, \eta)$  has the highest priority.
- (2)  $F_i$  containing the free surface error  $\eta_h - \eta, \chi_h - \chi$  or  $\partial_t(\eta_h - \eta)$ , or the numerical solution  $\eta_h, \chi_h$  or  $\partial_t \eta_h$  has the second priority.
- (3)  $F_i$  containing the projection error  $u - \mathcal{P}u$  or  $p - \mathcal{P}p$  for some  $\mathcal{P}$  has the third priority.
- (4)  $F_i$  only containing the velocity error  $u_h - \mathcal{P}u$  or the pressure error  $p_h - \mathcal{P}p$  for some projection  $\mathcal{P}$ , or the numerical solution  $u_h$  or  $p_h$ , has the lowest priority.

In this fashion, we can make the best use of the regularity of the exact solutions and the known error estimates in the free surface to avoid the estimates of  $(u_h, p_h)$  in undesired norms due to the embedding theorem.

**Rule 2:**

For the boundary term, we apply the trace theorem

$$\|F\|_{L^2(\partial E)} \lesssim \|F\|_{H^r(E)}, \tag{4.23}$$

for  $r > 1/2$ . This introduces more regularity in the estimates, so we should always try to avoid to use it directly and apply certain projection property to eliminate the boundary terms.

**Rule 3:**

Note the simple fact that for  $r > 0$ , we have

$$h^r \|f_h\|_{H^r(E)} \lesssim \|f_h\|_{L^2(E)}, \tag{4.24}$$

where  $h$  is the mesh size and  $f_h \in P_h^k$ . This can bound the higher order Sobolev norm by the lower norm at the price of some order of  $h$ .

In the error equations (4.16) and (4.17), we always take the test functions  $v_h = \epsilon_u$  and  $q_h = \epsilon_p$ . Since the estimates are standard, we omit the details here and just present the results. We can simplify the error equations (4.16) and (4.17) as:

$$W_\zeta + W_\gamma + W_\alpha + W_\beta = -R_\zeta - R_\gamma - R_\alpha - R_\beta - R_\mu - R_\omega \tag{4.25}$$

$$W_\rho = -R_\rho. \tag{4.26}$$

Then we can utilize a similar argument as in the proof of Theorem 3.3 to bound the energy part as

$$\frac{1}{2} \partial_t \int_\Omega J_h |\epsilon_u|^2 + C(Q) \|\epsilon_u\|_{H_h}^2 \lesssim R_\rho - R_\zeta - R_\gamma - R_\alpha - R_\beta - R_\mu - R_\omega. \tag{4.27}$$

Utilizing the estimates about the remaining part, we can further obtain

$$\begin{aligned} & \frac{1}{2} \partial_t \int_{\Omega} J_h |\epsilon_u|^2 + C(Q) \|\epsilon_u\|_{H_h}^2 \lesssim \\ & h^{k-1-r} \|\epsilon_u\|_{H_h} + h^{k-r} \|\epsilon_p\|_{L^2} + h^{k-1} (1 + \|\epsilon_u\|_{L^2}) \|\epsilon_u\|_{H_h} + \|\epsilon_u\|_{L^2} \|\epsilon_u\|_{H_h} \quad \text{for } 1/2 < r < 1. \end{aligned} \quad (4.28)$$

The stability shows both  $u_h$  and  $u$  are bounded in  $L^2$ , so  $\epsilon_u$  is also bounded in  $L^2$ . Then we get

$$\frac{1}{2} \partial_t \int_{\Omega} J_h |\epsilon_u|^2 + C(Q) \|\epsilon_u\|_{H_h}^2 \lesssim h^{k-1-r} \|\epsilon_u\|_{H_h} + h^{k-r} \|\epsilon_p\|_{L^2} + \|\epsilon_u\|_{L^2} \|\epsilon_u\|_{H_h} \quad \text{for } 1/2 < r < 1. \quad (4.29)$$

Applying Cauchy's inequality to the last term of (4.29) implies

$$\begin{aligned} & \frac{1}{2} \partial_t \int_{\Omega} J_h |\epsilon_u|^2 + C(Q) \|\epsilon_u\|_{H_h}^2 \lesssim h^{k-1-r} \|\epsilon_u\|_{H_h} + h^{k-r} \|\epsilon_p\|_{L^2} + \frac{1}{4C'} \|\epsilon_u\|_{L^2}^2 + C' \|\epsilon_u\|_{H_h}^2 \\ & \text{for } 1/2 < r < 1. \end{aligned} \quad (4.30)$$

Taking  $C'$  sufficiently small, we can absorb it into the left-hand side to achieve

$$\frac{1}{2} \partial_t \int_{\Omega} J_h |\epsilon_u|^2 + \frac{C(Q)}{2} \|\epsilon_u\|_{H_h}^2 \lesssim h^{k-1-r} \|\epsilon_u\|_{H_h} + h^{k-r} \|\epsilon_p\|_{L^2} + C(Q) \|\epsilon_u\|_{L^2}^2 \quad \text{for } 1/2 < r < 1. \quad (4.31)$$

This is the desired result. □

### 4.2. Pressure error analysis

In order for further analysis, we first need to show the inf-sup condition of our pressure discretization.

**Lemma 4.3.** *In the pressure discretization  $\beta(p_h, \eta_h, v_h) = \sum_{E \in \mathcal{E}_h} \beta_E(p_h, \eta_h, v_h)$ , suppose Conditions 3.1, 3.2 and 4.1 are valid. Assume the exact solution satisfies  $\eta(t) \in C^1(\Sigma)$ . Then for sufficiently small time  $T$  which is independent of  $h$ , there exists a constant  $\Xi > 0$  such that*

$$\inf_{p_h \in M_h^{k-1}} \sup_{v_h \in X_h^k} \frac{\beta(p_h, \eta_h, v_h)}{\|p_h\|_{L^2} \|v_h\|_{H_h}} \geq \Xi, \quad (4.32)$$

for  $t \in [0, T]$ , where  $\Xi$  is independent of  $h$ .

*Proof.* This result is similar to ([8], Thm. 4.5), so we omit the details here. Note the fact that for  $J_0 \mathcal{A}_0$  which is the value of  $J \mathcal{A}$  at  $t = 0$ , and  $\bar{J}_0 \bar{\mathcal{A}}_0$  the average of  $J_0 \mathcal{A}_0$  in each cell  $E$  which is piecewise constants, we naturally have

$$\|J_0 \mathcal{A}_0 - \bar{J}_0 \bar{\mathcal{A}}_0\|_{L^\infty} \lesssim h. \quad (4.33)$$

Hence, the weight functions in the differential operators is only a small perturbation, which can be absorbed into the main part. □

**Lemma 4.4.** *Suppose Conditions 3.1, 3.2 and 4.1 are valid. Assume the exact solution satisfies  $u \in C^2(\Omega)$ ,  $p \in C^1(\Omega)$  and  $\eta \in C^1(\Sigma)$ . Then for sufficiently small  $T$  and  $h$ , the numerical solution to the scheme (2.14) satisfies the estimate*

$$\int_0^t \|\epsilon_p\|_{L^2} \lesssim C(Q) \left( h^{k-1-r} + \int_0^t \|\epsilon_u\|_{L^2} + \|\epsilon_u(t)\|_{L^2} + \frac{1}{h} \int_0^t \|\epsilon_u\|_{H_h} \right) \quad \text{for } 1/2 < r < 1, \quad (4.34)$$

within  $t \in [0, T]$ .



*Proof.* We can take any test function  $v_h \in X_h^k$  in the error equation (4.16) and integrate over time  $[0, t]$ . Here we do not distinguish between the energy part and remaining part, but only concentrate on

$$\int_0^t \beta(p_h, \eta_h, v_h). \quad (4.35)$$

All the other terms can be moved into the right-hand side and estimated in terms of  $\|v_h\|_{H_h}$ . Basically, the estimates are similar to the velocity error estimates in the proof of Lemma 4.2, so we omit the details here. Based on the inf-sup condition in Lemma 4.3, we can deduce the desired result.  $\square$

### 4.3. Error analysis of the fluid

**Theorem 4.5.** *Suppose Conditions 3.1, 3.2 and 4.1 are valid. Assume the exact solution satisfies  $u \in C^2(\Omega)$ ,  $p \in C^1(\Omega)$  and  $\eta \in C^1(\Sigma)$ . Then for sufficiently small  $T$  and  $h$ , the numerical solution to the scheme (2.14) satisfies the estimates*

$$\|u_h - u\|_{L^2} \lesssim C(Q)h^{k-1-r}, \quad (4.36)$$

$$\left( \int_0^t \|u_h - u\|_{H_h}^2 \right)^{1/2} \lesssim C(Q)h^{k-1-r}, \quad (4.37)$$

$$\int_0^t \|p_h - p\|_{L^2} \lesssim C(Q)h^{k-1-r} \quad \text{for } 1/2 < r < 1, \quad (4.38)$$

in  $t \in [0, T]$ .

*Proof.* In Lemmas 4.2 and 4.4, we have shown

$$\partial_t \int_{\Omega} J_h |\epsilon_u|^2 + \|\epsilon_u\|_{H_h}^2 \lesssim C(Q) \left( h^{k-1-r} \|\epsilon_u\|_{H_h} + h^{k-r} \|\epsilon_p\|_{L^2} + \|\epsilon_u\|_{L^2}^2 \right) \quad \text{for } 1/2 < r < 1, \quad (4.39)$$

and

$$\int_0^t \|\epsilon_p\|_{L^2} \lesssim C(Q) \left( h^{k-1-r} + \int_0^t \|\epsilon_u\|_{L^2} + \|\epsilon_u(t)\|_{L^2} + \frac{1}{h} \int_0^t \|\epsilon_u\|_{H_h} \right) \quad \text{for } 1/2 < r < 1. \quad (4.40)$$

Integrating over  $[0, t]$  for any  $t \in [0, T]$  in (4.39), we obtain

$$\int_{\Omega} J_h(t) |\epsilon_u(t)|^2 + \int_0^t \|\epsilon_u\|_{H_h}^2 \lesssim C(Q) \left( h^{k+2} + h^{k-1-r} \int_0^t \|\epsilon_u\|_{H_h} + h^{k-r} \int_0^t \|\epsilon_p\|_{L^2} + \int_0^t \|\epsilon_u\|_{L^2}^2 \right). \quad (4.41)$$

Then we plug (4.40) into (4.41) to eliminate  $\epsilon_p$  and obtain

$$\begin{aligned} & \int_{\Omega} J_h(t) |\epsilon_u(t)|^2 + \int_0^t \|\epsilon_u\|_{H_h}^2 \\ & \lesssim C(Q) \left( h^{2k-1-2r} + h^{k-1-r} \int_0^t \|\epsilon_u\|_{H_h} + h^{k-r} \|\epsilon_u(t)\|_{L^2} + \int_0^t \|\epsilon_u\|_{L^2}^2 \right) \\ & \lesssim C(Q) \left( h^{2k-1-2r} + h^{2k-2-2r} + \left( \int_0^t \|\epsilon_u\|_{H_h} \right)^2 + h^{k-r} \|\epsilon_u(t)\|_{L^2} + \int_0^t \|\epsilon_u\|_{L^2}^2 \right) \\ & \lesssim C(Q) \left( h^{2k-2-2r} + \sqrt{t} \int_0^t \|\epsilon_u\|_{H_h}^2 + h^{k-r} \|\epsilon_u(t)\|_{L^2} + \int_0^t \|\epsilon_u\|_{L^2}^2 \right). \end{aligned} \quad (4.42)$$

When  $T$  is sufficiently small, we can absorb  $C(Q)\sqrt{t} \int_0^t \|\epsilon_u\|_{H_h}^2$  into the left-hand side of (4.42) to achieve

$$\int_{\Omega} J_h(t) |\epsilon_u(t)|^2 + \int_0^t \|\epsilon_u\|_{H_h}^2 \lesssim C(Q) \left( h^{2k-2-2r} + h^{k-r} \|\epsilon_u(t)\|_{L^2} + \int_0^t \|\epsilon_u\|_{L^2}^2 \right). \quad (4.43)$$

Since (4.43) holds for any  $t \in [0, T]$ , we can take the maximum to get

$$\max_{t \in [0, T]} \|\epsilon_u(t)\|_{L^2}^2 \lesssim C(Q) \left( h^{2k-2-2r} + h^{k-r} \max_{t \in [0, T]} \|\epsilon_u(t)\|_{L^2} + T \max_{t \in [0, T]} \|\epsilon_u(t)\|_{L^2}^2 \right). \tag{4.44}$$

When  $T$  and  $h$  are sufficiently small, we can absorb  $C(Q)(h^{k-r} + T) \max_{t \in [0, T]} \|\epsilon_u\|_{L^2}^2$  into the left-hand side of (4.44) to obtain

$$\max_{t \in [0, T]} \|\epsilon_u(t)\|_{L^2}^2 \lesssim C(Q)h^{2k-2-2r}. \tag{4.45}$$

Then this leads to

$$\max_{t \in [0, T]} \|\epsilon_u(t)\|_{L^2} \lesssim C(Q)h^{k-1-r}. \tag{4.46}$$

Hence, we plug (4.46) into (4.43) and (4.40) to achieve

$$\int_0^t \|\epsilon_u\|_{H_h}^2 \lesssim C(Q)h^{2k-2-2r}, \tag{4.47}$$

$$\int_0^t \|\epsilon_p\|_{L^2} \lesssim C(Q)h^{k-1-r}. \tag{4.48}$$

Combining with the projection errors  $\delta_u$  and  $\delta_p$ , we show the desired result. □

**Remark 4.6.** The convergent rate  $k - 1 - r$  is not optimal. If we can obtain a nicer error estimate of the temporal term, then this rate can be further improved.

### 5. NUMERICAL TESTS

We perform the accuracy tests for our numerical scheme. Our tests are based on a set of exact solutions

$$\eta = Z \sin(2\pi(x_1 - t)), \tag{5.1}$$

$$u_1 = Z(6x_2^5 + 15x_2^4 + 12x_2^3 + 3x_2^2) \sin(2\pi(x_1 - t)) + (6x_2^5 + 5x_2^4 + 1), \tag{5.2}$$

$$u_2 = \pi Z^2(4x_2^6 + 15x_2^5 + 21x_2^4 + 13x_2^3 + 3x_2^2) \sin(4\pi(x_1 - t)) \tag{5.3}$$

$$+ 2\pi Z(4x_2^6 + 7x_2^5 + 2x_2^4 - x_2^3) \cos(2\pi(x_1 - t)),$$

$$p = Z \sin(2\pi(x_1 - t)), \tag{5.4}$$

with the source term

$$\begin{aligned} S_1 = & (u_1)_t - \eta_t(1 + x_2)K\partial_2 u_1 + (u_1(\partial_1 u_1 - AK\partial_2 u_1) + Ku_2\partial_2 u_1) \\ & - \nu(\partial_{11} u_1 + (1 + A^2)K^2\partial_{22} u_1 - 2AK\partial_{12} u_1 + (AK^2\partial_2 A - A\partial_1 K - K\partial_1 A)\partial_2 u_1) \\ & + (\partial_1 p - AK\partial_2 p), \end{aligned} \tag{5.5}$$

$$\begin{aligned} S_2 = & (u_2)_t - \eta_t(1 + x_2)K\partial_2 u_2 + (u_1(\partial_1 u_2 - AK\partial_2 u_2) + Ku_2\partial_2 u_2) \\ & - \nu(\partial_{11} u_2 + (1 + A^2)K^2\partial_{22} u_2 - 2AK\partial_{12} u_2 + (AK^2\partial_2 A - A\partial_1 K - K\partial_1 A)\partial_2 u_2) \\ & + K\partial_2 p, \end{aligned} \tag{5.6}$$

where  $0 < Z < 1$  can be any fixed constant,  $A$  and  $K$  are defined as in (1.7), and differential operators  $\partial_t$ ,  $\partial_i$  and  $\partial_{ij}$  are defined in the usual sense. In the following test, we always take  $Z = 0.1$ ,  $\nu = 0.05$  and  $T = 1/8$ .

### 5.1. Accuracy tests with SIPG

The following are the error tables for our numerical scheme with SIPG:

TABLE 1.  $L^2$  error table for  $X_h^1 - M_h^0 - S_h^1$  formulation with SIPG.

(a) Free Surface Error $\eta - \eta_h$			(b) Velocity Error $u - u_h$			(c) Pressure Error $p - p_h$		
$N$	Error	Order	$N$	Error	Order	$N$	Error	Order
4	9.4028E-4	–	4	4.8233E-2	–	4	3.1028E-3	–
8	2.5493E-4	1.8830	8	1.3863E-2	1.7988	8	1.6255E-3	0.9327
16	6.6612E-5	1.9362	16	3.6182E-3	1.9378	16	8.2709E-4	0.9748
32	1.7409E-5	1.9359	32	9.4871E-4	1.9312	32	4.1570E-4	0.9925
64	4.9397E-6	1.8173	64	2.4593E-4	1.9477	64	2.0818E-4	0.9977

TABLE 2.  $L^2$  error table for  $X_h^2 - M_h^1 - S_h^2$  formulation with SIPG.

(a) Free Surface Error $\eta - \eta_h$			(b) Velocity Error $u - u_h$			(c) Pressure Error $p - p_h$		
$N$	Error	Order	$N$	Error	Order	$N$	Error	Order
4	1.4435E-4	–	4	5.8590E-3	–	4	9.6076E-4	–
8	1.6283E-5	3.1481	8	7.7408E-4	2.9201	8	2.3182E-4	2.0512
16	2.1495E-6	2.9213	16	9.8674E-5	2.9717	16	5.4527E-5	2.0880
32	2.9961E-7	2.8428	32	1.2586E-5	2.9709	32	1.2942E-5	2.0749

### 5.2. Accuracy tests with NIPG

The following are the error tables for our numerical scheme with NIPG:

TABLE 3.  $L^2$  error table for  $X_h^1 - M_h^0 - S_h^1$  formulation with NIPG.

(a) Free Surface Error $\eta - \eta_h$			(b) Velocity Error $u - u_h$			(c) Pressure Error $p - p_h$		
$N$	Error	Order	$N$	Error	Order	$N$	Error	Order
4	9.3809E-4	–	4	4.5428E-2	–	4	3.2035E-3	–
8	2.5179E-4	1.8975	8	1.8218E-2	1.3182	8	1.6481E-3	0.9588
16	6.5861E-5	1.9347	16	6.7565E-3	1.4310	16	8.4487E-4	0.9640
32	1.7258E-5	1.9322	32	2.1758E-3	1.6347	32	4.2453E-4	0.9929
64	4.9108E-6	1.8132	64	7.5074E-4	1.5352	64	2.1239E-4	0.9992

TABLE 4.  $L^2$  error table for  $X_h^2 - M_h^1 - S_h^2$  formulation with NIPG.

(a) Free Surface Error $\eta - \eta_h$			(b) Velocity Error $u - u_h$			(c) Pressure Error $p - p_h$		
$N$	Error	Order	$N$	Error	Order	$N$	Error	Order
4	1.4810E-4	–	4	6.0580E-3	–	4	9.2305E-4	–
8	1.7895E-5	3.0489	8	1.2693E-3	2.2548	8	2.8213E-4	1.7100
16	3.3218E-6	2.4295	16	1.9305E-4	2.7170	16	7.5701E-5	1.8980
32	8.5392E-7	1.9598	32	2.7534E-5	2.8097	32	2.0531E-5	1.8825

### 5.3. Discussion on the numerical tests

In above accuracy tests, we can obtain the optimal order of convergence when applying scheme (2.14) with SIPG, and the sub-optimal order with NIPG. However, in both cases, our numerical results are much better than the analytical result obtained in Theorem 4.5.

## 6. CONCLUSIONS AND REMARKS

In this paper, we construct a stable numerical scheme to solve the system (1.9) with discontinuous Galerkin method, and discuss the error analysis in the fluid with certain assumptions.

Although we focus on the 2-D fluid throughout this paper, it is easy to see this scheme can be naturally extended to the 3-D case with periodic settings for two horizontal directions. The main restriction to our scheme is that we require the free surface to be a single-valued function of the horizontal variables, which is not always true in practice, especially when the topological structure of the free surface varies during the evolution. Our scheme is non-conservative. Hence, it might not give the qualitatively correct simulation if the exact solution possesses singularities. However, for smooth exact solutions, our scheme can give a quite good approximation. Here we use the piecewise polynomial  $P^k$  due to the connection between the bulk discretization and surface discretization. Note that  $Q^k$ , which denotes the piecewise polynomials of degree at most  $k$  for each variables, is acceptable in the scheme, but not suitable in analysis since it cannot satisfy the desired inf-sup condition for pressure term in the error analysis.

*Acknowledgements.* The author thanks the editor and referees for their constructive comments and suggestions.

## REFERENCES

- [1] M.J. Ahn, H.Y. Lee and M.R. Ohm, Error estimates for fully discrete approximation to a free boundary problem in polymer technology. *Appl. Math. Comput.* **138** (2003) 227–238.
- [2] J. Thomas Beale, The initial value problem for the Navier–Stokes equations with a free surface. *Commun. Pure Appl. Math.* **34** (1981) 359–392.
- [3] B. Cockburn and C.-W. Shu, TVB Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws. II. General framework. *Math. Comput.* **52** (1989) 411–435.
- [4] B. Cockburn and C.-W. Shu, Runge-Kutta discontinuous Galerkin methods for convection-dominated problems. *J. Sci. Comput.* **16** (2001) 173–261.
- [5] B. Cockburn, S. Hou and C.-W. Shu, The Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws. IV. The multidimensional case. *Math. Comput.* **54** (1990) 545–581.
- [6] R. Danchin, Estimates in Besov spaces for transport and transport-diffusion equations with almost Lipschitz coefficients. *Rev. Mat. Iberoamericana* **21** (2005) 863–888.
- [7] V. Girault, B. Riviere and M.F. Wheeler, A splitting method using discontinuous Galerkin for the transient incompressible Navier-Stokes equations. *ESAIM: M2AN* **39** (2005) 1115–1147.
- [8] V. Girault, B. Riviere and M.F. Wheeler, A discontinuous Galerkin method with nonoverlapping domain decomposition for the Stokes and Navier-Stokes problems. *Math. Comput.* **74** (2005) 249–53–84.
- [9] J. Grooss and J.S. Hesthaven, A level set discontinuous Galerkin method for free surface flows. *Comput. Methods Appl. Mech. Engrg.* **195** (2006) 3406–3429.
- [10] Y. Guo and I. Tice, Local well-posedness of the viscous surface wave problem without surface tension. *Anal. PDE* **6** (2013) 287–369.
- [11] Y. Guo and I. Tice, Almost exponential decay of periodic viscous surface waves without surface tension. *Arch. Ration. Mech. Anal.* **207** (2013) 459–531.
- [12] E. Hairer, C. Lubich and M. Roche, The numerical solution of differential-algebraic systems by Runge-Kutta methods. In vol. 1409 of *Lect. Notes Math.* Springer-Verlag, Berlin (1989).
- [13] F.H. Harlow and J. Eddie Welch, Numerical calculation of time-dependent viscous incompressible flow of fluid with free surface. *Phys. Fluids* **8** (1965) 2182–2189.
- [14] C.W. Hirt and B.D. Nichols, Volume of fluid (VOF) method for the dynamics of free boundaries. *J. Comput. Phys.* **39** (1981) 201–225.
- [15] H.Y. Lee, Error analysis of finite element approximation of a Stefan problem with nonlinear free boundary condition. *J. Appl. Math. Comput.* **22** (2006) 223–235.
- [16] R.H. Nochetto and C. Verdi, An efficient linear scheme to approximate parabolic free boundary problems: error estimates and implementation. *Math. Comput.* **51** (1988) 27–53.
- [17] M. Sussman and M.Y. Hussaini, A discontinuous spectral element method for the level set equation. *J. Sci. Comput.* **19** (2003) 479–500.
- [18] L.-H. Wang, On Korn’s inequality. *J. Comput. Math.* **21** (2003) 321–324.
- [19] M.F. Wheeler, An elliptic collocation-finite element method with interior penalties. *SIAM J. Numer. Anal.* **15** (1978) 152–161.
- [20] L. Wu, Well-posedness and decay of the viscous surface wave. *SIAM J. Math. Anal.* **46** (2014) 2084–2135.