

A Second Order Energy Stable Linear Scheme for a Thin Film Model Without Slope Selection

Weijia Li¹ · Wenbin Chen² · Cheng Wang³ ·
Yue Yan⁴  · Ruijian He⁵

Received: 2 November 2017 / Revised: 1 February 2018 / Accepted: 5 March 2018 /
Published online: 10 March 2018
© Springer Science+Business Media, LLC, part of Springer Nature 2018

Abstract In this paper we present a second order accurate, energy stable numerical scheme for the epitaxial thin film model without slope selection, with a mixed finite element approximation in space. In particular, an explicit treatment of the nonlinear term, $\frac{\nabla u}{1+|\nabla u|^2}$, greatly simplifies the computational effort; only one linear equation with constant coefficients needs to be solved at each time step. Meanwhile, a second order Douglas–Dupont regularization term, $A\tau\Delta^2(u^{n+1} - u^n)$, is added in the numerical scheme, so that an unconditional long time energy stability is assured. In turn, we perform an $\ell^\infty(0, T; L^2)$ convergence analysis for the proposed scheme, with an $O(\tau^2 + h^q)$ error estimate derived. In addition, an optimal convergence analysis is provided for the nonlinear term using Q_q finite elements, which shows that the spatial convergence order can be improved to $q + 1$ on regular rectangular mesh. A few numerical experiments are presented, which confirms the efficiency and accuracy of the proposed second order numerical scheme.

✉ Yue Yan
yan.yue@mail.shufe.edu.cn

Weijia Li
15110180013@fudan.edu.cn

Wenbin Chen
wbchen@fudan.edu.cn

Cheng Wang
cwang1@umassd.edu

Ruijian He
hejian010@gmail.com

¹ School of Mathematical Sciences, Fudan University, Shanghai 200433, China

² School of Mathematical Sciences and Shanghai Key Laboratory for Contemporary Applied Mathematics, Fudan University, Shanghai 200433, China

³ Mathematics Department, University of Massachusetts, North Dartmouth, MA 02747, USA

⁴ School of Mathematics, Shanghai University of Finance and Economics, Shanghai 200433, China

⁵ Department of Chemical and Petroleum Engineering, Schulich School of Engineering, University of Calgary, Calgary, AB T2N 1N4, Canada

Keywords Epitaxial thin film growth · Slope selection · Energy stability · Optimal convergence analysis · Mixed finite element

1 Introduction

In this article we consider an epitaxial thin film growth model. The equation is the gradient flow associated with the following energy functional

$$E(u) \equiv \int_{\Omega} \left(-\frac{1}{2} \ln(1 + |\nabla u|^2) + \frac{\varepsilon^2}{2} |\Delta u|^2 \right) dx, \quad (1)$$

where $\Omega \subset \mathbb{R}^2$ is a bounded domain, ε is a positive constant, and $u : \Omega \rightarrow \mathbb{R}$ is the height function. And also, we denote $\partial_n u$ the exterior normal derivative. In order to eliminate the boundary integral terms in the variational derivative of the energy, the following natural boundary conditions are considered:

$$\partial_n u = \partial_n \Delta u = 0, \quad \text{on } \partial\Omega \times (0, T). \quad (2)$$

Then we have the following variational derivative of the energy

$$\delta_u E = \nabla \cdot \left(\frac{\nabla u}{1 + |\nabla u|^2} \right) + \varepsilon^2 \Delta^2 u. \quad (3)$$

Herein we consider the L^2 gradient flow

$$\partial_t u = -\delta_u E = -\nabla \cdot \left(\frac{\nabla u}{1 + |\nabla u|^2} \right) - \varepsilon^2 \Delta^2 u, \quad (4)$$

with the following initial conditions

$$\begin{aligned} u(\mathbf{x}, 0) &= u_0(\mathbf{x}), \quad \text{in } \Omega, \\ (u_0, 1) &= 0, \end{aligned} \quad (5)$$

where (\cdot, \cdot) represents the L^2 inner product. We refer to (4) as the no-slope-selection equation, following most other references. Now we introduce $w = -\Delta u$ and Eq. (4) can be rewritten as

$$\partial_t u + \nabla \cdot \left(\frac{\nabla u}{1 + |\nabla u|^2} \right) - \varepsilon^2 \Delta w = 0, \quad \text{in } \Omega \times (0, T], \quad (6)$$

$$w + \Delta u = 0, \quad \text{in } \Omega \times (0, T]. \quad (7)$$

From (2) and (5), we have the following boundary and initial conditions

$$\partial_n u = \partial_n w = 0, \quad \text{on } \partial\Omega \times (0, T], \quad (8)$$

$$u(\cdot, \mathbf{x}) = u_0(\mathbf{x}), \quad \text{in } \Omega. \quad (9)$$

$$(u, 1) = (w, 1) = 0. \quad (10)$$

Meanwhile, we see that the energy functional can be written as

$$E(u, w) \equiv \int_{\Omega} \left(-\frac{1}{2} \ln(1 + |\nabla u|^2) + \frac{\varepsilon^2}{2} |w|^2 \right) dx. \quad (11)$$

Also, note that the first term in the energy functional represents the Ehrlich–Schwoebel (ES) effect [14], therefore, we denote it by $E^{ES}(u)$:

$$E^{ES}(u) \equiv \int_{\Omega} -\frac{1}{2} \ln(1 + |\nabla u|^2) \, dx.$$

In [15], the global in time well-posedness for two nonlinear models of epitaxial thin film growth, with or without slope selection, was established. And also, the gradient bound and the energy asymptotic law for the epitaxial growth equation with or without slope selection have been studied in [13, 16], as $\varepsilon \rightarrow 0$. In addition, the large-system asymptotic form of the minimum energy and the magnitude of gradients of energy-minimizing surfaces for epitaxial growth models are analysed in [14], with infinite or finite Ehrlich–Schwoebel (ES) barrier. Specially, for the case of a finite ES effect (corresponding to the model in this article), the well-posedness of the initial-boundary-value problem is proved and the bounds for the scaling laws of interface width, surface slope and energy are obtained.

In terms of the numerical simulations for the epitaxial thin film growth model, there have been many efforts to devise and analyse schemes for both the slope selection and no-slope selection equations; see the related references [6, 20, 23, 27], etc. In particular, the numerical schemes with high order accuracy and energy stability have attracted a great deal of attentions, due to the long time nature of the gradient flow coarsening process. In the paper of Li and Liu [15], a classical second order accurate semi-implicit numerical scheme, combined with Galerkin spectral approximation in space, is used for solving the equations, while a theoretical justification of the numerical energy stability is not available. Among the energy stable numerical approaches, the idea of convex splitting is worthy of discussion. For the epitaxial thin film growth models, the first such work was reported in [26], in which the authors studied unconditionally energy stable schemes, based on the convex–concave decomposition of the energy, motivated by Eyre’s pioneering work [7]. Also see the other related works on the energy stable schemes for MBE models [11, 17, 19, 21, 22, 31], etc. Meanwhile, there are two obvious shortcomings of the schemes reported in [26]: only first order accuracy (in time), and the high degree of nonlinearity of the numerical scheme, due to the implicit treatment of the nonlinear term. In particular, a direct convex splitting solver for the no-slope selection equation (4) is even more challenging, since the nonlinear term appears in the denominator part. Subsequently, an efficient linear, unconditionally stable, unconditionally solvable scheme was proposed in [3] for (4), to overcome this prominent difficulty, based on an alternate, and more advantageous, way of the convex–concave energy decomposition; such an alternate decomposition places the nonlinear part of the chemical potential in the concave part instead of the convex part. In turn, the implicit part of the scheme is completely linear, which greatly improves the numerical efficiency, in comparison with the one in [26]. In fact, such a linear stabilization approach has been reported in an earlier work [27], and a theoretical justification of this linear stability has been available in a more recent work [17]. Moreover, the linear operator involved in the scheme, which is positive elliptic with constant coefficients, can be efficiently inverted by FFT.

Of course, the linear scheme reported in [3] is only first-order accurate in time. Many efforts have been devoted in the past few years to develop second order accurate, energy stable schemes for epitaxial thin film growth models, with or without slope selection. For example, the second-order convex splitting scheme, in the modified Crank–Nicolson version, is proposed and discussed in [25]. A careful analysis indicates the unconditional energy stability and unique nonlinear solvability. An alternate second order accurate scheme, in the backward differentiation formula (BDF) version, has been proposed and analysed in a more recent work [8]. On the other hand, it is noted that, these second order schemes are highly

nonlinear, and the numerical implementations become highly challenging. For the slope-selection model, the nonlinear term keeps in the polynomial format so that either a nonlinear conjugate gradient or preconditioned steepest descent (PSD) solver can be efficiently applied. However, for the no-slope-selection model (4), the numerical difficulty associated with the high degree of nonlinearity is much more prominent, due to the complicated terms appearing in the fractional quotients. To overcome this difficulty, a linear iteration solver is proposed in [5] to implement the highly nonlinear second order numerical scheme associated with the no-slope-selection equation (4). In more details, a second order accurate $O(\tau^2)$ artificial diffusion term, in the form of Douglas–Dupont regularization, is introduced to the numerical scheme, and a linear iteration algorithm is proposed to implement the highly nonlinear scheme in [25], associated with the no-slope-selection model. As a result, the highly nonlinear numerical scheme can be very efficiently solved by a linear iteration algorithm, and a geometric convergence order is assured for this linear iteration under a constraint associated with the artificial diffusion coefficient.

Meanwhile, it is observed that, although only a linear equation is needed at each iteration stage in the linear algorithm proposed in [5], the overall computational cost is still a few times of a linear equation at each time step, implied by the geometric convergence. Subsequently, a question naturally arises: could one derive a second order accurate, energy stable numerical scheme for the no-slope-selection model (4), with only one linear equation (with constant coefficients) involved at each time step? In this article, we propose and analyse such a numerical scheme. In more details, a second order backward differentiation formula (BDF) is applied to approximate the temporal derivative, while the surface diffusion term is treated implicitly. On the other hand, the nonlinear chemical potential is approximated by an explicit extrapolation formula at time step t^{n+1} , with second order temporal accuracy. Moreover, a second order accurate $O(\tau^2)$ artificial term, $A\tau\Delta(u^n - u^{n-1})$, is added in the numerical scheme for the sake of stability analysis. In turn, the numerical scheme is linear, with constant coefficients, at each time step, due to the explicit extrapolation approach used in the nonlinear term. Furthermore, a careful energy estimate indicates that, the energy stability could be justified for the proposed numerical scheme at a theoretical level, provided that the given constant $A \geq \frac{25}{16}$. Therefore, all the desired properties have been established for the proposed numerical scheme.

A mixed finite element approximation is taken in space, based on a mixed weak formulation of the no-slope-selection model (6)–(7). In this approach, the numerical solutions for both the phase variable u and the chemical potential variable w belong to the same finite element space X_h , which is a piecewise polynomial subspace of H^1 . In combination with the second order temporal approximation, the resulting numerical scheme preserves the properties of unique solvability and unconditional energy stability. With a help of this uniform-in-time energy bound, we are able to establish the convergence analysis, with an error estimate of $O(\tau^2 + h^q)$ accuracy in the $\ell^\infty(0, T; L^2) \cap \ell^2(0, T; H_h^2)$ norm. Since the nonlinear term and all its derivatives have a direct ℓ^∞ bound, this convergence is unconditional; no scaling law is needed between τ and h to ensure its validity.

Moreover, it is observed that, the h^q convergence order in space is not optimal; such a loss of accuracy comes from the gradient structure of the no-slope-selection equation (4). Meanwhile, there have been many articles obtaining spatial full order convergence for fourth order elliptic equations on regular rectangular meshes, such as [18, 28]. Based on these preliminary estimates, we are able to apply similar techniques and obtain an optimal spatial convergence order for the proposed numerical scheme, using regular rectangular mesh.

The rest of the article is organized as follows. In Sect. 2, we present the fully discrete scheme. The unique solvability and unconditional long time energy stability are proved in

Sects. 2.1 and 2.2, respectively, and the $O(\tau^2 + h^q)$ convergence analysis is presented in Sect. 2.3. Subsequently, the optimal convergence analysis is provided in Sect. 3. Besides, some numerical results are presented in Sect. 4. Finally, the concluding remarks are given in Sect. 5.

2 The Numerical Scheme

Referring to [1], we denote the standard norms for the Sobolev spaces $W^{m,p}(\Omega)$ by $\|\cdot\|_{m,p}$. Let $H^m(\Omega)$ denote $W^{m,2}(\Omega)$. We replace $\|\cdot\|_{q,2}$ by $\|\cdot\|_q$, and $\|\cdot\|_{L^2}$ by $\|\cdot\|$. Also, $L^2_0(\Omega) \equiv \{\varphi \in L^2(\Omega) \mid (\varphi, 1) = 0\}$. Now we introduce the Sobolev space $X = \{v \in H^q(\Omega) \mid (v, 1) = 0\}$.

We denote by $L^2(0, T; X)$ the set of all the quadratic integrable functions from $[0, T]$ to X . Similarly we use notation $L^2(0, T; H^{-q}(\Omega))$, where H^{-q} is the dual space of H^q . Then the weak form of (6)–(7) is to find u and $w \in L^2(0, T; X)$, with $u_t \in L^2(0, T; H^{-q}(\Omega))$, satisfying

$$(\partial_t u, v) + \varepsilon^2(\nabla w, \nabla v) - \left(\frac{\nabla u}{1 + |\nabla u|^2}, \nabla v \right) = 0, \quad \forall v \in X. \tag{12}$$

$$(w, \psi) - (\nabla u, \nabla \psi) = 0, \quad \forall \psi \in X. \tag{13}$$

Taking $v = \partial_t u$ in (12), $\psi = \varepsilon^2 \partial_t w$ in (13) and adding up the two equations, we have

$$\|\partial_t u\|^2 + \frac{\varepsilon^2}{2} \frac{d}{dt} \|w\|^2 - \frac{1}{2} \frac{d}{dt} (\ln(1 + |\nabla u|^2), 1) = 0. \tag{14}$$

Integrating (14) from t_0 to t_1 for any $0 \leq t_0 < t_1$, we have

$$\int_{t_0}^{t_1} \|\partial_t u(s)\|^2 ds + E(u(t_1), w(t_1)) = E(u(t_0), w(t_0)), \tag{15}$$

which shows that the system (6)–(7) is energy stable.

Let $\tau = \frac{T}{N}$ be time step size. As for the mesh $\mathcal{T}_h = \{K\}$ on Ω with related finite function space \mathcal{P} , we let either \mathcal{T}_h be a quasi-uniform triangulation with $\mathcal{P} = \mathcal{P}_q(K)$ or \mathcal{T}_h be a regular rectangular mesh with $\mathcal{P} = \mathcal{Q}_q$, where h stands for a discretization parameter, $\mathcal{Q}_q \equiv \text{span}\{x^i y^j : 0 \leq i, j \leq q\}$ and $\mathcal{P}_q(K)$ is the set of polynomials of degree $\leq q$. Define the piecewise polynomial space $X_h \equiv \{v \in X \cap C^0(\Omega) \mid v|_K \in \mathcal{P}(K), \forall K \in \mathcal{T}_h\} \subset X$.

In order to set the initialization step of our scheme, we introduce the Ritz projection $R_h : L^2_0(\Omega) \rightarrow X_h$,

$$(\nabla(R_h \varphi - \varphi), \nabla \chi) = 0, \quad (R_h \varphi - \varphi, 1) = 0, \quad \forall \chi \in X_h. \tag{16}$$

Now we propose the fully discrete numerical scheme. Let $u_h^0 = R_h u_0, w_h^0 = R_h(-\Delta u_0)$. Denote the numerical solution of u and w at time t_n by u_h^n and w_h^n respectively. The initialization step is defined as below: given u_h^0, w_h^0 , find $u_h^1, w_h^1 \in X_h$, such that for any v_h and $\psi_h \in X_h$,

$$\left(\frac{u_h^1 - u_h^0}{\tau}, v_h \right) + \varepsilon^2 (\nabla w_h^1, \nabla v_h) - \left(\frac{\nabla u_h^0}{1 + |\nabla u_h^0|^2}, \nabla v_h \right) + A^{(0)} (\nabla (u_h^1 - u_h^0), \nabla v_h) = 0, \tag{17}$$

$$(w_h^1, \psi_h) - (\nabla u_h^1, \nabla \psi_h) = 0. \tag{18}$$

The unique solvability comes from the positive-definite property of the involved linear operators. In addition, an unconditional energy stability, $E(u_h^1, w_h^1) \leq E(u_h^0, w_h^0)$, follows a similar analysis as given by [3], provided that $A^{(0)} \geq 1$. And also, the first order temporal accuracy in the first time step does not affect the overall second order accuracy, which will be analysed in later sections.

For $n \geq 1$, given u_h^{n-1}, u_h^n and $w_h^n \in X_h$, find $u_h^{n+1}, w_h^{n+1} \in X_h$, such that for any v_h and $\psi_h \in X_h$,

$$\left(\frac{3u_h^{n+1} - 4u_h^n + u_h^{n-1}}{2\tau}, v_h \right) + \varepsilon^2 \left(\nabla w_h^{n+1}, \nabla v_h \right) - \left(\frac{\nabla \left(2u_h^n - u_h^{n-1} \right)}{1 + \left| \nabla \left(2u_h^n - u_h^{n-1} \right) \right|^2}, \nabla v_h \right) + A\tau \left(\nabla \left(w_h^{n+1} - w_h^n \right), \nabla v_h \right) = 0, \tag{19}$$

$$\left(w_h^{n+1}, \psi_h \right) - \left(\nabla u_h^{n+1}, \nabla \psi_h \right) = 0. \tag{20}$$

Note that we have used explicit extrapolation formula for the nonlinear term, and an artificial term $A\tau \left(\nabla \left(w_h^{n+1} - w_h^n \right), \nabla v_h \right)$ is added in the numerical scheme. In turn, the unconditional unique solvability is assured by the fact that, all the implicit terms are associated with linear elliptic operators with positive eigenvalues.

The idea of the modified BDF method has been similarly applied to the Cahn–Hilliard equation [29] and the slope-selection (SS) epitaxial thin film model [8]. Meanwhile, for the NSS model, such a BDF approach has a very different feature from these existing works, in terms of the fully explicit treatment of the nonlinear term, in comparison with the fully implicit ones in [8,29]. This explicit treatment does not cause any stability difficulty, due to a subtle feature of the nonlinearity in the NSS model, as demonstrated in the analysis below.

2.1 Unique Solvability

The unique solvability analysis is available.

Theorem 1 *The scheme (19)–(20) is unconditionally uniquely solvable.*

Proof We rewrite scheme (19) and (20) as below:

$$\begin{cases} \frac{3}{2\tau} \left(u_h^{n+1}, v_h \right) + \varepsilon^2 \left(\nabla w_h^{n+1}, \nabla v_h \right) + A\tau \left(\nabla w_h^{n+1}, \nabla v_h \right) = f^{n,n-1}(v_h), \\ \left(w_h^{n+1}, \psi_h \right) - \left(\nabla u_h^{n+1}, \nabla \psi_h \right) = 0, \end{cases} \tag{21}$$

with

$$f^{n,n-1}(v_h) \equiv \frac{1}{2\tau} \left(4u_h^n - u_h^{n-1}, v_h \right) + \left(\frac{\nabla \left(2u_h^n - u_h^{n-1} \right)}{1 + \left| \nabla \left(2u_h^n - u_h^{n-1} \right) \right|^2}, \nabla v_h \right) + A\tau \left(\nabla w_h^n, \nabla v_h \right).$$

It is clear that $f^{n,n-1}(v_h)$ is a continuous linear functional of v_h . Let $\mathbf{u} = [u_h^{n+1}, w_h^{n+1}]$ and $\mathbf{q} = [\psi_h, v_h]$, and we define the bilinear form

$$a(\mathbf{u}, \mathbf{q}) = \frac{3}{2\tau} \left(u_h^{n+1}, v_h \right) + \varepsilon^2 \left(\nabla w_h^{n+1}, \nabla v_h \right) + A\tau \left(\nabla w_h^{n+1}, \nabla v_h \right) + \left(w_h^{n+1}, \psi_h \right) - \left(\nabla u_h^{n+1}, \nabla \psi_h \right).$$

Thus, Eq. (21) is equivalent to finding $\mathbf{u} \in X_h \times X_h$ such that

$$a(\mathbf{u}, \mathbf{q}) = f(\mathbf{q}), \quad \forall \mathbf{q} \in X_h \times X_h,$$

where $f(\mathbf{q}) = f^{n,n-1}(v_h)$ is a continuous linear functional on $X_h \times X_h$.

It is easy to verify that $|a(\mathbf{u}, \mathbf{q})|$ is bounded above by $C_{\varepsilon,\tau} \|\mathbf{u}\|_1 \|\mathbf{q}\|_1$, where $C_{\varepsilon,\tau}$ is a positive constant that only depends on ε, τ . Thus we only need to prove the coercive requirement. Notice that $(w_h^{n+1}, u_h^{n+1}) = \|\nabla u_h^{n+1}\|^2$, then for any $\tau > 0$,

$$\begin{aligned} a(\mathbf{u}, \mathbf{u}) &= \frac{3}{2\tau} (u_h^{n+1}, w_h^{n+1}) + (\varepsilon^2 + A\tau) \|\nabla w_h^{n+1}\|^2 + (w_h^{n+1}, u_h^{n+1}) - \|\nabla u_h^{n+1}\|^2 \\ &= \frac{3}{2\tau} \|\nabla u_h^{n+1}\|^2 + (\varepsilon^2 + A\tau) \|\nabla w_h^{n+1}\|^2 \\ &\geq c_{\varepsilon,\tau} \|\mathbf{u}\|_1, \end{aligned}$$

where we have used Poincaré’s inequality in the last step, and $c_{\varepsilon,\tau}$ is a positive constant only dependent on ε, τ . Now by the Lax-Milgram theorem in [2], (21) admits a unique solution $\mathbf{u} = [u_h^{n+1}, w_h^{n+1}] \in X_h \times X_h$. □

2.2 Long Time Energy Stability

Here we define the discrete energy functional as

$$\begin{aligned} \tilde{E}(u_h^{n+1}, u_h^n, w_h^{n+1}) &\equiv \int_{\Omega} -\frac{1}{2} \ln \left(1 + |\nabla u_h^{n+1}|^2 \right) dx + \frac{\varepsilon^2}{2} \|w_h^{n+1}\|^2 \\ &\quad + \frac{1}{4\tau} \|u_h^{n+1} - u_h^n\|^2 + \|\nabla (u_h^{n+1} - u_h^n)\|^2, \quad n \geq 0. \end{aligned} \tag{22}$$

Now we have the following energy stability estimate.

Theorem 2 *Given $A \geq \frac{25}{16}$, the second-order numerical scheme (19)–(20) has the energy-decay property*

$$\tilde{E}(u_h^{n+1}, u_h^n, w_h^{n+1}) \leq \tilde{E}(u_h^n, u_h^{n-1}, w_h^n), \quad n \geq 1. \tag{23}$$

Proof Firstly we deduce from (18) and (20) that: for $n \geq 0$,

$$(w_h^{n+1} - w_h^n, \psi_h) - (\nabla (u_h^{n+1} - u_h^n), \nabla \psi_h) = 0, \quad \forall \psi_h \in X_h. \tag{24}$$

Taking $v_h = u_h^{n+1} - u_h^n$ in (19), $\psi_h = \varepsilon^2 w_h^{n+1}$ in (24) and adding them up,

$$\begin{aligned} 0 &= \left(\frac{3u_h^{n+1} - 4u_h^n + u_h^{n-1}}{2\tau}, u_h^{n+1} - u_h^n \right) + A\tau (\nabla (w_h^{n+1} - w_h^n), \nabla (u_h^{n+1} - u_h^n)) \\ &\quad + (w_h^{n+1} - w_h^n, \varepsilon^2 w_h^{n+1}) - \left(\frac{\nabla (2u_h^n - u_h^{n-1})}{1 + |\nabla (2u_h^n - u_h^{n-1})|^2}, \nabla (u_h^{n+1} - u_h^n) \right). \end{aligned} \tag{25}$$

Below we estimate the four terms on the right side of (25), successively. As for the first term, we apply the Cauchy–Schwarz inequality:

$$\left(\frac{3u_h^{n+1} - 4u_h^n + u_h^{n-1}}{2\tau}, u_h^{n+1} - u_h^n \right) \geq \frac{1}{\tau} \left(\frac{5}{4} \|u_h^{n+1} - u_h^n\|^2 - \frac{1}{4} \|u_h^n - u_h^{n-1}\|^2 \right). \tag{26}$$

Taking $\psi_h = w_h^{n+1} - w_h^n$ in (24) yields

$$A\tau \left(\nabla \left(w_h^{n+1} - w_h^n \right), \nabla \left(u_h^{n+1} - u_h^n \right) \right) = A\tau \left\| w_h^{n+1} - w_h^n \right\|^2. \tag{27}$$

The third term can be directly computed as

$$\begin{aligned} \left(w_h^{n+1} - w_h^n, \varepsilon^2 w_h^{n+1} \right) &= \frac{\varepsilon^2}{2} \left(\left\| w_h^{n+1} \right\|^2 - \left\| w_h^n \right\|^2 + \left\| w_h^{n+1} - w_h^n \right\|^2 \right) \\ &\geq \frac{\varepsilon^2}{2} \left(\left\| w_h^{n+1} \right\|^2 - \left\| w_h^n \right\|^2 \right). \end{aligned} \tag{28}$$

As for the fourth term (I) $\equiv - \left(\frac{\nabla(2u_h^n - u_h^{n-1})}{1 + |\nabla(2u_h^n - u_h^{n-1})|^2}, \nabla(u_h^{n+1} - u_h^n) \right)$, we notice that $\ln(1+x) \leq x$ for $x > -1$, so that

$$-\ln \left(1 + \left| \nabla u_h^{n+1} \right|^2 \right) + \ln \left(1 + \left| \nabla u_h^n \right|^2 \right) = \ln \left(\frac{1 + \left| \nabla u_h^n \right|^2}{1 + \left| \nabla u_h^{n+1} \right|^2} \right) \leq \frac{\left| \nabla u_h^n \right|^2 - \left| \nabla u_h^{n+1} \right|^2}{1 + \left| \nabla u_h^{n+1} \right|^2}.$$

Therefore, recalling that $E^{ES}(u_h) = \int_{\Omega} -\frac{1}{2} \ln(1 + |\nabla u_h|^2) \, dx$, we get

$$\begin{aligned} E^{ES} \left(u_h^{n+1} \right) - E^{ES} \left(u_h^n \right) &= \int_{\Omega} \left(-\frac{1}{2} \ln \left(1 + \left| \nabla u_h^{n+1} \right|^2 \right) + \frac{1}{2} \ln \left(1 + \left| \nabla u_h^n \right|^2 \right) \right) dx \\ &\leq \frac{1}{2} \int_{\Omega} \frac{\left| \nabla u_h^n \right|^2 - \left| \nabla u_h^{n+1} \right|^2}{1 + \left| \nabla u_h^{n+1} \right|^2} dx \\ &= -\frac{1}{2} \int_{\Omega} \frac{\nabla \left(u_h^{n+1} - u_h^n \right) \cdot \nabla \left(u_h^{n+1} + u_h^n \right)}{1 + \left| \nabla u_h^{n+1} \right|^2} dx. \end{aligned}$$

As a consequence, the following inequality is valid:

$$\begin{aligned} \text{(I)} &- \left[E^{ES} \left(u_h^{n+1} \right) - E^{ES} \left(u_h^n \right) \right] \\ &\geq - \left(\frac{\nabla \left(2u_h^n - u_h^{n-1} \right)}{1 + \left| \nabla \left(2u_h^n - u_h^{n-1} \right) \right|^2}, \nabla \left(u_h^{n+1} - u_h^n \right) \right) \\ &\quad + \frac{1}{2} \int_{\Omega} \frac{\nabla \left(u_h^{n+1} - u_h^n \right) \cdot \nabla \left(u_h^{n+1} + u_h^n \right)}{1 + \left| \nabla u_h^{n+1} \right|^2} dx \\ &\geq \int_{\Omega} \left(\frac{-\nabla \left(u_h^{n+1} - u_h^n \right) \cdot \nabla \left(2u_h^n - u_h^{n-1} \right)}{1 + \left| \nabla \left(2u_h^n - u_h^{n-1} \right) \right|^2} \right. \\ &\quad \left. + \frac{\frac{1}{2} \nabla \left(u_h^{n+1} - u_h^n \right) \cdot \nabla \left(u_h^{n+1} + u_h^n \right)}{1 + \left| \nabla u_h^{n+1} \right|^2} \right) dx \equiv \int_{\Omega} \text{(II)} \, dx. \end{aligned} \tag{29}$$

Next we separate the term (II) into two parts and estimate them respectively:

$$\begin{aligned}
 \text{(II)} &= \left(\frac{-\nabla(u_h^{n+1} - u_h^n) \cdot \nabla(2u_h^n - u_h^{n-1})}{1 + |\nabla(2u_h^n - u_h^{n-1})|^2} - \frac{-\nabla(u_h^{n+1} - u_h^n) \cdot \nabla(2u_h^n - u_h^{n-1})}{1 + |\nabla u_h^{n+1}|^2} \right) \\
 &+ \left(\frac{-\nabla(u_h^{n+1} - u_h^n) \cdot \nabla(2u_h^n - u_h^{n-1})}{1 + |\nabla u_h^{n+1}|^2} + \frac{\frac{1}{2}\nabla(u_h^{n+1} - u_h^n) \cdot \nabla(u_h^{n+1} + u_h^n)}{1 + |\nabla u_h^{n+1}|^2} \right) \\
 &\equiv g_1 + g_2.
 \end{aligned}$$

As for the first part,

$$\begin{aligned}
 g_1 &\geq - \frac{|\nabla(u_h^{n+1} - u_h^n)| |\nabla(u_h^{n+1} - 2u_h^n + u_h^{n-1})|}{\left(1 + |\nabla(2u_h^n - u_h^{n-1})|^2\right) \left(1 + |\nabla u_h^{n+1}|^2\right)} \\
 &\quad \times \frac{|\nabla(2u_h^n - u_h^{n-1})| |\nabla(u_h^{n+1} + 2u_h^n - u_h^{n-1})|}{\left(1 + |\nabla(2u_h^n - u_h^{n-1})|^2\right) \left(1 + |\nabla u_h^{n+1}|^2\right)} \\
 &\geq - \frac{|\nabla(u_h^{n+1} - u_h^n)| |\nabla(u_h^{n+1} - 2u_h^n + u_h^{n-1})|}{\left(1 + |\nabla(2u_h^n - u_h^{n-1})|^2\right) \left(1 + |\nabla u_h^{n+1}|^2\right)} \\
 &\quad \times \frac{|\nabla(2u_h^n - u_h^{n-1})| \left(|\nabla u_h^{n+1}| + |\nabla(2u_h^n - u_h^{n-1})|\right)}{\left(1 + |\nabla(2u_h^n - u_h^{n-1})|^2\right) \left(1 + |\nabla u_h^{n+1}|^2\right)}.
 \end{aligned}$$

For any real number $a, b \geq 0$, we have $a(a + b) < (1 + a^2)(1 + b^2)$. Applying this property to the inequality above, we obtain

$$\begin{aligned}
 g_1 &\geq - \frac{|\nabla(u_h^{n+1} - u_h^n)| |\nabla(u_h^{n+1} - 2u_h^n + u_h^{n-1})|}{\left(1 + |\nabla(2u_h^n - u_h^{n-1})|^2\right) \left(1 + |\nabla u_h^{n+1}|^2\right)} \\
 &\geq - \frac{|\nabla(u_h^{n+1} - u_h^n)|^2 - |\nabla(u_h^{n+1} - u_h^n)| |\nabla(u_h^n - u_h^{n-1})|}{\left(1 + |\nabla(2u_h^n - u_h^{n-1})|^2\right) \left(1 + |\nabla u_h^{n+1}|^2\right)} \tag{30} \\
 &\geq -\frac{3}{2} \frac{|\nabla(u_h^{n+1} - u_h^n)|^2}{\left(1 + |\nabla(2u_h^n - u_h^{n-1})|^2\right) \left(1 + |\nabla u_h^{n+1}|^2\right)} + \frac{1}{2} \frac{|\nabla(u_h^n - u_h^{n-1})|^2}{\left(1 + |\nabla(2u_h^n - u_h^{n-1})|^2\right) \left(1 + |\nabla u_h^{n+1}|^2\right)}.
 \end{aligned}$$

Using the Cauchy–Schwarz inequality on g_2 yields

$$\begin{aligned}
 g_2 &= \frac{1}{1 + |\nabla u_h^{n+1}|^2} \left[\frac{1}{2} |\nabla(u_h^{n+1} - u_h^n)|^2 - \nabla(u_h^{n+1} - u_h^n) \cdot \nabla(u_h^n - u_h^{n-1}) \right] \\
 &\geq \frac{1}{1 + |\nabla u_h^{n+1}|^2} \left(-\frac{1}{2} |\nabla(u_h^n - u_h^{n-1})|^2 \right) \tag{31} \\
 &\geq -\frac{1}{2} \frac{|\nabla(u_h^n - u_h^{n-1})|^2}{1 + |\nabla u_h^{n+1}|^2}.
 \end{aligned}$$

Substituting (30) and (31) into (II), we have

$$\text{(II)} \geq -\frac{3}{2} \frac{|\nabla(u_h^{n+1} - u_h^n)|^2}{1 + |\nabla u_h^{n+1}|^2} + \frac{1}{2} \frac{|\nabla(u_h^n - u_h^{n-1})|^2}{1 + |\nabla u_h^{n+1}|^2}. \tag{32}$$

Therefore, (29) could be rewritten as

$$(I) - \left[E^{ES} \left(u_h^{n+1} \right) - E^{ES} \left(u_h^n \right) \right] \geq -\frac{3}{2} \left\| \nabla \left(u_h^{n+1} - u_h^n \right) \right\|^2 - \left\| \nabla \left(u_h^n - u_h^{n-1} \right) \right\|^2. \tag{33}$$

Substituting the estimates (26)–(28) and (33) for the four terms into (25), we obtain

$$\begin{aligned} 0 &\geq \frac{1}{4\tau} \left(\left\| u_h^{n+1} - u_h^n \right\|^2 - \left\| u_h^n - u_h^{n-1} \right\|^2 \right) + \frac{1}{\tau} \left\| u_h^{n+1} - u_h^n \right\|^2 + A\tau \left\| w_h^{n+1} - w_h^n \right\|^2 \\ &\quad + \frac{\varepsilon^2}{2} \left(\left\| w_h^{n+1} \right\|^2 - \left\| w_h^n \right\|^2 \right) + E^{ES} \left(u_h^{n+1} \right) - E^{ES} \left(u_h^n \right) - \frac{3}{2} \left\| \nabla \left(u_h^{n+1} - u_h^n \right) \right\|^2 \\ &\quad - \left\| \nabla \left(u_h^n - u_h^{n-1} \right) \right\|^2. \end{aligned} \tag{34}$$

Meanwhile, a combination of the Cauchy–Schwarz inequality and (24) implies that

$$\begin{aligned} &\frac{1}{\tau} \left\| u_h^{n+1} - u_h^n \right\|^2 + A\tau \left\| w_h^{n+1} - w_h^n \right\|^2 \\ &\quad \geq 2A^{\frac{1}{2}} \left\| u_h^{n+1} - u_h^n \right\| \left\| w_h^{n+1} - w_h^n \right\| \\ &\quad \geq 2A^{\frac{1}{2}} \left\| \nabla \left(u_h^{n+1} - u_h^n \right) \right\|^2 \geq \frac{5}{2} \left\| \nabla \left(u_h^{n+1} - u_h^n \right) \right\|^2, \end{aligned}$$

where we have used the assumption that $A \geq \frac{25}{16}$. Therefore, we can rewrite the inequality (34) as

$$\begin{aligned} 0 &\geq \frac{1}{4\tau} \left(\left\| u_h^{n+1} - u_h^n \right\|^2 - \left\| u_h^n - u_h^{n-1} \right\|^2 \right) \\ &\quad + \left(\left\| \nabla \left(u_h^{n+1} - u_h^n \right) \right\|^2 - \left\| \nabla \left(u_h^n - u_h^{n-1} \right) \right\|^2 \right) \\ &\quad + \frac{\varepsilon^2}{2} \left(\left\| w_h^{n+1} \right\|^2 - \left\| w_h^n \right\|^2 \right) + E^{ES} \left(u_h^{n+1} \right) - E^{ES} \left(u_h^n \right) \\ &= \tilde{E} \left(u_h^{n+1}, u_h^n, w_h^{n+1} \right) - \tilde{E} \left(u_h^n, u_h^{n-1}, w_h^n \right). \end{aligned}$$

which is the conclusion we need. □

Remark 1 For the no-slope-selection model (4), there have been other second order numerical schemes, in which the energy stability is defined over an “alternate” energy functional; see the related works of [21,31], etc. In these numerical approaches, the energy functional is numerically defined, and the nonlinear energy density is based on an alternate numerical variable. In turn, such an energy stability does not justify an H^2 bound for the numerical solution, at a theoretical level. In comparison, the energy stability analysis derived in this section is based on the original phase variable u , so that the desired bound is available for our proposed numerical scheme. This subtle property will provide a great deal of convenience in the convergence analysis given by later sections.

Remark 2 The energy stability analysis is established in the framework of Galerkin approximation, which comes from the finite element spatial approximation. On the other hand, it is observed that, the nonlinear integral values could hardly be exactly computed in the practical computations, due to the highly complicated nature of the denominator form. Instead, a numerical approximation to these nonlinear integral values, which corresponds to the collocation approach, has to be taken into consideration.

In the case of a uniform spatial mesh, the Fourier collocation spectral scheme has been analysed in an existing work [3], and unconditional energy stability has been proved for the first order linear splitting scheme. An extension to the second order accurate linear iteration algorithm has been reported in [5]. For the second order linear scheme proposed in this article, the energy stability analysis for the collocation approximation is expected to be available, and the details will be considered in the future works.

2.3 The $\ell^\infty(\mathbf{0}, T; L^2) \cap \ell^2(\mathbf{0}, T; H_h^2)$ Convergence Analysis

We present the convergence analysis in this section. First of all, referring to [2], the following estimate holds for the Ritz projection $R_h : \forall \varphi \in H^{q+1}(\Omega) \cap X$,

$$\|\varphi - R_h\varphi\| + h\|\nabla(\varphi - R_h\varphi)\| \leq Ch^{q+1}\|\varphi\|_{q+1}. \tag{35}$$

The following discrete Gronwall inequality [24] is needed in the error analysis.

Lemma 1 Assume that $\tau > 0, B > 0, \{a_n\}, \{b_n\}, \{\gamma_n\}$ are non-negative sequences such that

$$a_m + \tau \sum_{n=1}^m b_n \leq \tau \sum_{n=1}^{m-1} \gamma_n a_n + B, \quad m \geq 1.$$

Then,

$$a_m + \tau \sum_{n=1}^m b_n \leq B \exp\left(\tau \sum_{n=1}^{m-1} \gamma_n\right), \quad m \geq 1.$$

Besides, we will also use the inverse estimate in [2, p.111, Lemma 4.5.3].

Lemma 2 Given a quasi-uniform triangulation \mathcal{T}_h ($h \leq 1$) on domain $\Omega \subset \mathbb{R}^n$ and the related finite dimensional function subspace $X_h \subset W^{l,p} \cap W^{m,q}$ with $1 \leq p \leq \infty, 1 \leq q \leq \infty$ and $0 \leq m \leq l$, there exists a constant C such that for any $\chi \in X_h, K \in \mathcal{T}_h$, we have

$$\|\chi\|_{W^{l,p}(K)} \leq \tilde{C}h^{m-l+n/p-n/q}\|\chi\|_{W^{m,q}(K)},$$

where \tilde{C} is independent of h and χ .

The following lemma is also needed in the analysis.

Lemma 3 Given functions $\varphi_1 \in X, \varphi_2 \in X$ and $v \in X$, let us define function $\Phi : [0, 1] \rightarrow \mathbb{R}$

$$\Phi_{\varphi_1, \varphi_2, v}(s) = \left(\frac{\nabla(\varphi_1 + s(\varphi_2 - \varphi_1))}{1 + |\nabla(\varphi_1 + s(\varphi_2 - \varphi_1))|^2}, \nabla v \right), \tag{36}$$

then we have

$$|\Phi_{\varphi_1, \varphi_2, v}(1) - \Phi_{\varphi_1, \varphi_2, v}(0)| \leq \|\nabla(\varphi_2 - \varphi_1)\| \|\nabla v\|.$$

Proof The derivative of Φ with respect to s is

$$\begin{aligned} &\Phi'_{\varphi_1, \varphi_2, v}(s) \\ &= \left(\frac{\nabla(\varphi_2 - \varphi_1)}{1 + |\nabla(\varphi_1 + s(\varphi_2 - \varphi_1))|^2}, \nabla v \right) - \left(\frac{2|\nabla(\varphi_1 + s(\varphi_2 - \varphi_1))|^2 \nabla(\varphi_2 - \varphi_1)}{(1 + |\nabla(\varphi_1 + s(\varphi_2 - \varphi_1))|^2)^2}, \nabla v \right) \\ &= \left(\frac{1 - |\nabla(\varphi_1 + s(\varphi_2 - \varphi_1))|^2}{(1 + |\nabla(\varphi_1 + s(\varphi_2 - \varphi_1))|^2)^2} \nabla(\varphi_2 - \varphi_1), \nabla v \right) \end{aligned}$$

Thus we have the estimate

$$|\Phi'_{\varphi_1, \varphi_2, v}(s)| \leq \|\nabla(\varphi_2 - \varphi_1)\| \|\nabla v\|.$$

Applying the above estimate, we have

$$|\Phi_{\varphi_1, \varphi_2, v}(1) - \Phi_{\varphi_1, \varphi_2, v}(0)| \leq \int_0^1 |\Phi'_{\varphi_1, \varphi_2, v}(s)| \, ds \leq \|\nabla(\varphi_2 - \varphi_1)\| \|\nabla v\|.$$

which is the conclusion we need. □

We denote by (u, w) the exact solution pair to the original equation (4), and all the upper bounds for the exact solution are denoted as C_0 . We say that the solution pair is in the regularity class \mathcal{C} if and only if

$$\begin{aligned} u &\in L^\infty(0, T; H^{q+1}) \cap H^1(0, T; H^{q+1}) \cap H^2(0, T; H^1) \\ &\quad \cap W^{1,\infty}(0, T; H^2) \cap W^{2,\infty}(0, T; L^2) \cap H^3(0, T; L^2), \\ w &\in L^\infty(0, T; H^{q+1}) \cap H^1(0, T; H^1). \end{aligned} \tag{37}$$

The following theorem is the main result of this section.

Theorem 3 *Suppose that the exact solution pair (u, w) is in the regularity class \mathcal{C} , for a fixed final time $T > 0$. Denote $u(t_n)$ by u^n and let u_h^n be the solution to the fully discrete numerical scheme (17)–(20) at time $t_m = m\tau$, for $1 \leq m \leq N$ with $N\tau = T$. Assume that*

$$0 < \tau < \frac{\varepsilon^2}{2\alpha^2}, \quad 0 < h \leq 1,$$

where $\alpha > 3$ is a constant, then we have the following error estimate

$$\|u^m - u_h^m\| + \left(\tau \sum_{n=1}^m \|w^n - w_h^n\|^2 \right)^{\frac{1}{2}} \leq C_{\varepsilon, T} (h^q + \tau^2). \tag{38}$$

Proof The error functions are defined as

$$\begin{aligned} e_u^n &\equiv \rho_u^n + \sigma_u^n \equiv (u^n - R_h u^n) + (R_h u^n - u_h^n) = u^n - u_h^n, \\ e_w^n &\equiv \rho_w^n + \sigma_w^n \equiv (w^n - R_h w^n) + (R_h w^n - w_h^n) = w^n - w_h^n. \end{aligned} \tag{39}$$

Subtracting the numerical scheme formulation (19)–(20) from the weak form (12)–(13), we obtain the following error equations:

$$\begin{aligned} (\delta_\tau^{n+1} e_u, v_h) + \varepsilon^2 (\nabla e_w^{n+1}, \nabla v_h) + A\tau (\nabla (e_w^{n+1} - e_w^n), \nabla v_h) \\ = (\mathcal{R}_1^{n+1}, v_h) + A\tau (\mathcal{R}_2^{n+1}, \nabla v_h) + (\mathcal{N}^{n+1}, \nabla v_h), \quad \forall v_h \in X_h, \\ (e_w^{n+1}, \psi_h) - (\nabla e_u^{n+1}, \nabla \psi_h) = 0, \quad \forall \psi_h \in X_h, \end{aligned}$$

for any $n \geq 1$, where

$$\begin{aligned} \delta_\tau^{n+1} v &= \frac{3v^{n+1} - 4v^n + v^{n-1}}{2\tau}, \quad \mathcal{R}_1^{n+1} = \delta_\tau^{n+1} u - \partial_t u^{n+1}, \\ \mathcal{R}_2^{n+1} &= \nabla (w^{n+1} - w^n), \quad \mathcal{N}^{n+1} = \frac{\nabla u^{n+1}}{1 + |\nabla u^{n+1}|^2} - \frac{\nabla (2u_h^n - u_h^{n-1})}{1 + |\nabla (2u_h^n - u_h^{n-1})|^2}. \end{aligned}$$

Notice that the definition (16) of R_h indicates that: $(\nabla \rho_u^{n+1}, \nabla \chi) = 0, \forall \chi \in X_h$. Thus the error equations can be rewritten as: For $n \geq 1$, and for any $v_h \in X_h, \psi_h \in X_h$,

$$\begin{aligned}
 & (\delta_\tau^{n+1} \sigma_u, v_h) + \varepsilon^2 (\nabla \sigma_w^{n+1}, \nabla v_h) + A\tau (\nabla (\sigma_w^{n+1} - \sigma_w^n), \nabla v_h) \\
 & = -(\delta_\tau^{n+1} \rho_u, v_h) + (\mathcal{R}_1^{n+1}, v_h) + A\tau (\mathcal{R}_2^{n+1}, \nabla v_h) + (\mathcal{N}^{n+1}, \nabla v_h),
 \end{aligned} \tag{40}$$

and

$$(\sigma_w^{n+1}, \psi_h) - (\nabla \sigma_u^{n+1}, \nabla \psi_h) = -(\rho_w^{n+1}, \psi_h). \tag{41}$$

With a slight modification, we obtain the scheme for the initialization step: For $n = 0$, for any $v_h \in X_h, \psi_h \in X_h$,

$$\begin{aligned}
 & (\delta_\tau^1 \sigma_u^1, v_h) + \varepsilon^2 (\nabla \sigma_w^1, \nabla v_h) + A^{(0)} (\nabla (\sigma_u^1 - \sigma_u^0), \nabla v_h) \\
 & = -(\delta_\tau^1 \rho_u, v_h) + (\mathcal{R}_1^1, v_h) + A^{(0)} (\mathcal{R}_2^1, \nabla v_h) + (\mathcal{N}^1, \nabla v_h),
 \end{aligned} \tag{42}$$

$$(\sigma_w^1, \psi_h) - (\nabla \sigma_u^1, \nabla \psi_h) = -(\rho_w^1, \psi_h), \tag{43}$$

where

$$\delta_\tau^1 v = \frac{v^1 - v^0}{\tau}, \mathcal{R}_1^1 = \delta_\tau^1 u - \partial_t u^1, \mathcal{R}_2^1 = \nabla (u^1 - u^0), \mathcal{N}^1 = \frac{\nabla u^1}{1 + |\nabla u^1|^2} - \frac{\nabla u_h^0}{1 + |\nabla u_h^0|^2}.$$

Now we focus on the case when $n \geq 1$. Taking $v_h = \sigma_u^{n+1}$ in (40), $\psi_h = \varepsilon^2 \sigma_w^{n+1}$ in (41) and adding up the two equations lead to

$$\begin{aligned}
 & (\delta_\tau^{n+1} \sigma_u, \sigma_u^{n+1}) + \varepsilon^2 \|\sigma_w^{n+1}\|^2 + A\tau (\sigma_w^{n+1} - \sigma_w^n, \sigma_w^{n+1}) \\
 & = -A\tau (\sigma_w^{n+1} - \sigma_w^n, \rho_w^{n+1}) - \varepsilon^2 (\rho_w^{n+1}, \sigma_w^{n+1}) - (\delta_\tau^{n+1} \rho_u, \sigma_u^{n+1}) + (\mathcal{R}_1^{n+1}, \sigma_u^{n+1}) \\
 & \quad + A\tau (\mathcal{R}_2^{n+1}, \nabla \sigma_u^{n+1}) + (\mathcal{N}^{n+1}, \nabla \sigma_u^{n+1}),
 \end{aligned} \tag{44}$$

where we have used the transformation

$$\begin{aligned}
 A\tau (\nabla (\sigma_w^{n+1} - \sigma_w^n), \nabla \sigma_u^{n+1}) & = A\tau (\nabla (\sigma_w^{n+1} - \sigma_w^n), \nabla e_u^{n+1}) \\
 & = A\tau (\sigma_w^{n+1} - \sigma_w^n, e_w^{n+1}) \\
 & = A\tau (\sigma_w^{n+1} - \sigma_w^n, \sigma_w^{n+1}) + A\tau (\sigma_w^{n+1} - \sigma_w^n, \rho_w^{n+1}).
 \end{aligned}$$

First we focus on the terms on the left-hand side. In order to estimate the first term, we recall the G-norm introduced in [4]. Let $\mathbf{p}^{k+1} \equiv [\sigma_u^k, \sigma_u^{k+1}]^T$, and

$$\|\mathbf{p}^{k+1}\|_{\mathbf{G}}^2 \equiv (\mathbf{p}^{k+1}, \mathbf{G}\mathbf{p}^{k+1}), \quad \mathbf{G} = \begin{pmatrix} \frac{1}{2} & -1 \\ -1 & \frac{5}{2} \end{pmatrix}. \tag{45}$$

Applying this notation to the first term, we have

$$(\delta_\tau^{n+1} \sigma_u, \sigma_u^{n+1}) = \frac{1}{2\tau} (\|\mathbf{p}^{n+1}\|_{\mathbf{G}}^2 - \|\mathbf{p}^n\|_{\mathbf{G}}^2) + \frac{1}{4\tau} \|\sigma_u^{n+1} - 2\sigma_u^n + \sigma_u^{n-1}\|^2. \tag{46}$$

The third term can be represented as

$$A\tau (\sigma_w^{n+1} - \sigma_w^n, \sigma_w^{n+1}) = \frac{A\tau}{2} (\|\sigma_w^{n+1}\|^2 - \|\sigma_w^n\|^2 + \|\sigma_w^{n+1} - \sigma_w^n\|^2). \tag{47}$$

Now we estimate the terms on the right-hand side. As for the first two terms, applying the property of Ritz projection in (35) and the Cauchy–Schwarz inequality yields

$$\begin{aligned}
 -A\tau (\sigma_w^{n+1} - \sigma_w^n, \rho_w^{n+1}) &\leq \frac{A\tau}{2} \|\sigma_w^{n+1} - \sigma_w^n\|^2 + C\tau h^{2(q+1)} \|w^{n+1}\|_{1+q}^2, \\
 -\varepsilon^2 (\rho_w^{n+1}, \sigma_w^{n+1}) &\leq \frac{\varepsilon^2}{2} \|\sigma_w^{n+1}\|^2 + C\varepsilon^2 h^{2(q+1)} \|w^{n+1}\|_{1+q}^2.
 \end{aligned}
 \tag{48}$$

To analyse the third and fourth terms, we use Taylor expansion and the Cauchy–Schwarz inequality:

$$\begin{aligned}
 -(\delta_\tau^{n+1} \rho_u, \sigma_u^{n+1}) &= \left(\frac{3\rho_u^{n+1} - 4\rho_u^n + \rho_u^{n-1}}{2\tau}, \sigma_u^{n+1} \right) \\
 &\leq \frac{C_1}{2} \|\sigma_u^{n+1}\|^2 + C \frac{h^{2(q+1)}}{\tau} \int_{t_{n-1}}^{t_{n+1}} \|\partial_t u\|_{q+1}^2 dt, \\
 (\mathcal{R}_1^{n+1}, \sigma_u^{n+1}) &= \left(\frac{3u^{n+1} - 4u^n + u^{n-1}}{2\tau} - \partial_t u^{n+1}, \sigma_u^{n+1} \right) \\
 &\leq \frac{C_1}{2} \|\sigma_u^{n+1}\|^2 + C\tau^3 \int_{t_{n-1}}^{t_{n+1}} \|\partial_{ttt} u\|^2 dt.
 \end{aligned}
 \tag{49}$$

Applying Taylor expansion and the Cauchy–Schwarz inequality, for the fifth term we get

$$\begin{aligned}
 A\tau (\mathcal{R}_2^{n+1}, \nabla \sigma_u^{n+1}) &= A\tau (\nabla (w^{n+1} - w^n), \nabla \sigma_u^{n+1}) \\
 &\leq \frac{C_2}{6} \|\nabla \sigma_u^{n+1}\|^2 + C\tau^3 \int_{t_n}^{t_{n+1}} \|\nabla \partial_t w\|^2 dt.
 \end{aligned}
 \tag{50}$$

Notice that

$$\|\nabla \sigma_u^{n+1}\|^2 = (\nabla \sigma_u^{n+1}, \nabla e_u^{n+1}) = (\sigma_u^{n+1}, e_w^{n+1}) \leq \|\sigma_u^{n+1}\| (\|\sigma_w^{n+1}\| + \|\rho_w^{n+1}\|).$$

Therefore, using again the Cauchy–Schwarz inequality, we have:

$$\begin{aligned}
 &\frac{C_2}{6} \|\nabla \sigma_u^{n+1}\|^2 \\
 &\leq \frac{C_2}{6} \|\sigma_u^{n+1}\| (\|\sigma_w^{n+1}\| + \|\rho_w^{n+1}\|) \\
 &\leq \frac{C_2}{6} \left(\frac{1}{4C_3} \|\sigma_u^{n+1}\|^2 + 2C_3 \|\sigma_w^{n+1}\|^2 + 2C_3 C h^{2(q+1)} \|w^{n+1}\|_{1+q}^2 \right),
 \end{aligned}
 \tag{51}$$

which leads to

$$\begin{aligned}
 A\tau (\mathcal{R}_2^{n+1}, \nabla \sigma_u^{n+1}) &\leq \frac{C_2}{24C_3} \|\sigma_u^{n+1}\|^2 + \frac{C_2 C_3}{3} \|\sigma_w^{n+1}\|^2 \\
 &\quad + C h^{2(q+1)} \|w^{n+1}\|_{1+q}^2 + C\tau^3 \int_{t_n}^{t_{n+1}} \|\nabla \partial_t w\|^2 dt.
 \end{aligned}
 \tag{52}$$

As for the nonlinear term, we recall the function Φ defined in (36) and Lemma 3, and arrive at

$$\begin{aligned} & (\mathcal{N}^{n+1}, \nabla \sigma_u^{n+1}) \\ &= \Phi(2u_h^n - u_h^{n-1}, u^{n+1}, \sigma_u^{n+1})(1) - \Phi(2u_h^n - u_h^{n-1}, u^{n+1}, \sigma_u^{n+1})(0) \\ &\leq \left\| \nabla \left(u^{n+1} - 2u_h^n + u_h^{n-1} \right) \right\| \left\| \nabla \sigma_u^{n+1} \right\| \\ &\leq \left(\left\| \nabla \left(u^{n+1} - 2u_h^n + u_h^{n-1} \right) \right\| + 2 \left\| \nabla \sigma_u^n \right\| + 2 \left\| \nabla \rho_u^n \right\| + \left\| \nabla \sigma_u^{n-1} \right\| + \left\| \nabla \rho_u^{n-1} \right\| \right) \left\| \nabla \sigma_u^{n+1} \right\|. \end{aligned}$$

Now we use the Cauchy–Schwarz inequality and estimate $\|\nabla \sigma_u^{n+1}\|^2$ as in (51):

$$\begin{aligned} (\mathcal{N}^{n+1}, \nabla \sigma_u^{n+1}) &\leq \frac{5C_2}{24C_3} \|\sigma_u^{n+1}\|^2 + \frac{27}{2C_2C_5} \|\sigma_u^n\|^2 + \frac{27}{8C_2C_4} \|\sigma_u^{n-1}\|^2 \\ &\quad + \frac{5C_2C_3}{6} \|\sigma_w^{n+1}\|^2 + \frac{54C_5}{C_2} \|\sigma_w^n\|^2 + \frac{27C_4}{2C_2} \|\sigma_w^{n-1}\|^2 \\ &\quad + Ch^{2(q+1)} \left(\|w^{n+1}\|_{1+q}^2 + \|w^n\|_{1+q}^2 + \|w^{n-1}\|_{1+q}^2 \right) \\ &\quad + C\tau^3 \int_{t_n}^{t_{n+1}} \|\nabla \partial_{tt} u\|^2 dt + Ch^{2q} \left(\|u^n\|_{1+q}^2 + \|u^{n-1}\|_{1+q}^2 \right). \end{aligned} \tag{53}$$

Substituting estimates (46)–(53) into the error Eq. (44), we obtain

$$\begin{aligned} & \frac{1}{2\tau} \left(\|\mathbf{p}^{n+1}\|_{\mathbf{G}}^2 - \|\mathbf{p}^n\|_{\mathbf{G}}^2 \right) + \frac{\varepsilon^2}{2} \|\sigma_w^{n+1}\|^2 + \frac{A\tau}{2} \left(\|\sigma_w^{n+1}\|^2 - \|\sigma_w^n\|^2 \right) \\ &\leq \left(C_1 + \frac{C_2}{4C_3} \right) \|\sigma_u^{n+1}\|^2 + \frac{27}{2C_2C_5} \|\sigma_u^n\|^2 + \frac{27}{8C_2C_4} \|\sigma_u^{n-1}\|^2 \\ &\quad + \frac{7C_2C_3}{6} \|\sigma_w^{n+1}\|^2 + \frac{54C_5}{C_2} \|\sigma_w^n\|^2 + \frac{27C_4}{2C_2} \|\sigma_w^{n-1}\|^2 \\ &\quad + C \left(\tau h^{2(q+1)} + h^{2(q+1)} \right) \|w\|_{L^\infty(0,T;H^{q+1})}^2 + C \frac{h^{2(q+1)}}{\tau} \int_{t_{n-1}}^{t_{n+1}} \|\partial_{tt} u\|_{q+1}^2 dt \\ &\quad + C\tau^3 \int_{t_{n-1}}^{t_{n+1}} \left(\|\partial_{ttt} u\|^2 + \|\nabla \partial_{tt} u\|^2 + \|\nabla \partial_t w\|^2 \right) dt + Ch^{2q} \|u\|_{L^\infty(0,T;H^{1+q})}^2. \end{aligned}$$

Summing up from $n = 1$ to $n = m$ and multiplying by 2τ on both sides, we get

$$\begin{aligned} & \|\mathbf{p}^{m+1}\|_{\mathbf{G}}^2 - \|\mathbf{p}^1\|_{\mathbf{G}}^2 + 2\tau\varepsilon^2 \sum_{n=1}^m \|\sigma_w^{n+1}\|^2 + A\tau^2 \left(\|\sigma_w^{m+1}\|^2 - \|\sigma_w^1\|^2 \right) \\ &\leq \left(2C_1 + \frac{C_2}{2C_3} \right) \tau \|\sigma_u^{m+1}\|^2 + \left(2C_1 + \frac{C_2}{2C_3} + \frac{27}{C_2C_5} + \frac{27}{4C_2C_4} \right) \tau \sum_{n=1}^{m-1} \|\sigma_u^{n+1}\|^2 \\ &\quad + \left(\frac{7C_2C_3}{3} + \frac{108C_5}{C_2} + \frac{27C_4}{C_2} \right) \tau \sum_{n=1}^m \|\sigma_w^{n+1}\|^2 + \left(\frac{27}{C_2C_5} + \frac{27}{4C_2C_4} \right) \tau \|\sigma_u^1\|^2 \\ &\quad + \left(\frac{108C_5}{C_2} + \frac{27C_4}{C_2} \right) \tau \|\sigma_w^1\|^2 + C_{\varepsilon,T} \left(\tau h^{2(q+1)} + h^{2(q+1)} + h^{2q} + \tau^4 \right). \end{aligned}$$

It is easy to verify that $\|\mathbf{P}^{m+1}\|_{\mathbf{G}}^2 \geq \frac{1}{2}\|\sigma_u^{m+1}\|^2$, $\|\mathbf{P}^1\|_{\mathbf{G}}^2 = \frac{5}{2}\|\sigma_u^1\|^2$. Taking $C_1 = \frac{C_2^2}{\varepsilon^2}$, $C_3 = \frac{\varepsilon^2}{4C_2}$, $C_4 = \frac{C_2\varepsilon^2}{108}$, $C_5 = \frac{C_2\varepsilon^2}{432}$, we have

$$\begin{aligned} & \frac{\varepsilon^2 - 8C_2^2\tau}{2\varepsilon^2} \|\sigma_u^{m+1}\|^2 + \frac{11\tau\varepsilon^2}{12} \sum_{n=1}^m \|\sigma_w^{n+1}\|^2 \\ & \leq C_{\varepsilon,T} (h^{2q} + \tau^4) + \left(4C_1 + \frac{27}{C_2C_5} + \frac{27}{4C_2C_4}\right) \tau \sum_{n=1}^{m-1} \|\sigma_u^{n+1}\|^2 \tag{54} \\ & \quad + \left(\frac{5}{2} + \frac{27\tau}{C_2C_5} + \frac{27\tau}{4C_2C_4}\right) \|\sigma_u^1\|^2 + \left(A\tau^2 + \frac{\tau\varepsilon^2}{2}\right) \|\sigma_w^1\|^2. \end{aligned}$$

In order to estimate $\|\sigma_w^1\|$ and $\|\sigma_u^1\|$, we take $v_h = \sigma_u^1$, $\psi_h = \varepsilon^2\sigma_w^1$ in (42)–(43) and add up:

$$\begin{aligned} & \frac{1}{\tau} \|\sigma_u^1\|^2 + \varepsilon^2 \|\sigma_w^1\|^2 + A^{(0)} \|\nabla\sigma_u^1\|^2 \\ & = -(\delta_\tau^1 \rho_u, \sigma_u^1) + (\mathcal{R}_1^1, \sigma_u^1) + A^{(0)} (\mathcal{R}_2^1, \nabla\sigma_u^1) + (\mathcal{N}^1, \nabla\sigma_u^1) - \varepsilon^2 (\rho_w^1, \sigma_w^1). \tag{55} \end{aligned}$$

Similar estimates of the right-hand side terms could be obtained as when $n \geq 1$. For the first two terms, we use Taylor expansion and the Cauchy–Schwarz inequality:

$$\begin{aligned} -(\delta_\tau^1 \rho_u, \sigma_u^1) & \leq \frac{\tilde{C}_1}{3\tau} \|\sigma_u^1\|^2 + Ch^{2(q+1)} \int_0^\tau \|\partial_{tt}u\|_{q+1}^2 dt. \\ (\mathcal{R}_1^1, \sigma_u^1) & = \left(-\frac{1}{\tau} \int_0^\tau t \partial_{tt}u dt, \sigma_u^1\right) \\ & \leq \frac{\tilde{C}_1}{3\tau} \|\sigma_u^1\|^2 + C\tau^2 \int_0^\tau \|\partial_{tt}u\|^2 dt \\ & \leq \frac{\tilde{C}_1}{3\tau} \|\sigma_u^1\|^2 + C\tau^3 \|u\|_{W^{2,\infty}(0,T;L^2)}^2. \end{aligned}$$

For the third term we make use of integration by parts, the regularity of the exact solution and the Cauchy–Schwarz inequality:

$$\begin{aligned} A^{(0)} (\mathcal{R}_2^1, \nabla\sigma_u^1) & = A^{(0)} (\nabla(u^1 - u^0), \nabla\sigma_u^1) = -A^{(0)} (\Delta(u^1 - u^0), \sigma_u^1) \\ & \leq 2A^{(0)} \tau \|u\|_{W^{1,\infty}(0,T;H^2)} \|\sigma_u^1\| \leq \frac{\tilde{C}_1}{3\tau} \|\sigma_u^1\|^2 + C\tau^3. \end{aligned}$$

As for the nonlinear term, applying Lemma 3 and the technique used in the above inequality yields:

$$\begin{aligned} (\mathcal{N}^1, \nabla\sigma_u^1) & = \left(\frac{\nabla u^1}{1+|\nabla u^1|^2} - \frac{\nabla u^0}{1+|\nabla u^0|^2}, \nabla\sigma_u^1\right) + \left(\frac{\nabla u^0}{1+|\nabla u^0|^2} - \frac{\nabla u_h^0}{1+|\nabla u_h^0|^2}, \nabla\sigma_u^1\right) \\ & \leq -\left(\nabla \cdot \left(\frac{\nabla u^1}{1+|\nabla u^1|^2} - \frac{\nabla u^0}{1+|\nabla u^0|^2}\right), \sigma_u^1\right) + \|\nabla\rho_u^0\| \|\nabla\sigma_u^1\| \\ & \leq C\tau \|\sigma_u^1\| + Ch^q \|u_0\|_{q+1} \|\nabla\sigma_u^1\| \\ & \leq \frac{\tilde{C}_2}{\tau} \|\sigma_u^1\|^2 + \frac{C\tau^3}{\tilde{C}_2} + \frac{Ch^{2q}}{\tilde{C}_3} \|u_0\|_{q+1}^2 + \tilde{C}_3 \|\nabla\sigma_u^1\|^2. \tag{56} \end{aligned}$$

For the last term we use the Cauchy–Schwarz inequality again:

$$-\varepsilon^2 (\rho_w^1, \sigma_w^1) \leq \tilde{C}_4 \varepsilon^2 \|\sigma_w^1\|^2 + Ch^{2(q+1)} \|w^1\|_{q+1}^2.$$

Substituting the above estimates into (55), taking $\tilde{C}_1 = \tilde{C}_2 = \tilde{C}_4 = \frac{1}{4}$, $\tilde{C}_3 = \frac{A^{(0)}}{2}$ and multiplying by τ on both sides of (55) yields

$$\frac{1}{2} \|\sigma_u^1\|^2 + \frac{3\tau\varepsilon^2}{4} \|\sigma_w^1\|^2 + \frac{A^{(0)}\tau}{2} \|\nabla\sigma_u^1\|^2 \leq C_\varepsilon (\tau h^{2(q+1)} + \tau^4 + \tau h^{2q}). \tag{57}$$

Substituting (57) into (54), we get

$$\begin{aligned} & \frac{\varepsilon^2 - 2\alpha^2\tau}{2\varepsilon^2} \|\sigma_u^{m+1}\|^2 + \frac{\tau\varepsilon^2}{4} \sum_{n=1}^m \|\sigma_w^{n+1}\|^2 \\ & \leq C_{\varepsilon,T} (h^{2q} + \tau^4) + \left(4C_1 + \frac{27}{C_2C_5} + \frac{27}{4C_2C_4}\right) \tau \sum_{n=1}^{m-1} \|\sigma_u^{n+1}\|^2. \end{aligned} \tag{58}$$

Applying the Gronwall inequality (Lemma 1), we obtain

$$\frac{\varepsilon^2 - 2\alpha^2\tau}{2\varepsilon^2} \|\sigma_u^{m+1}\| + \left(\frac{\tau\varepsilon^2}{4} \sum_{n=1}^m \|\sigma_w^{n+1}\|^2\right)^{\frac{1}{2}} \leq C_{\varepsilon,T} (h^q + \tau^2).$$

A combination of the above estimate for $\|\sigma_u^{m+1}\|$ and $\|\sigma_w^{n+1}\|$ with (35) yields the conclusion we need. □

3 Optimal Convergence Analysis

The error estimate for the proposed numerical scheme in the previous section has indicated an h^q spatial convergence order, and a $(q + 1)$ th convergence order has not been theoretically available because of the difficulty in analysing the nonlinear term, while the numerical results shown in Table 1 indicate that the scheme has a $(q + 1)$ th convergence order. In this section, for rectangular domains aligned with x–y axis, by using \mathcal{Q}_q finite elements on rectangular meshes, this gap between the numerical results and the theoretical analysis can be overcome.

Recall the notations introduced in Sect. 2: Given a regular rectangular mesh $\mathcal{T}_h = \{K\}$ on a rectangular domain $\Omega \subset \mathbb{R}^2$ aligned with x–y axis, of which $\{a_i\}$ and $\{l_j\}$ denote the element vertices and edges. Set shape function space $\mathcal{P} = \mathcal{Q}_q = span\{x^i y^j : 0 \leq i, j \leq q\}$, and we define the finite function space $X_h \equiv \{v \in X \cap C^0(\Omega) \mid v|_K \in \mathcal{Q}_q(K), \forall K \in \mathcal{T}_h\} \subset X$ with $X = \{v \in H^q(\Omega) \mid (v, 1) = 0\}$.

Now we introduce the interpolation operator $i_h^q : C^0(\bar{\Omega}) \rightarrow X_h$ defined in [9, p.108]:

$$i_h^q w(a_i) = w(a_i), \tag{59}$$

$$\int_{l_j} (i_h^q w - w) v \, ds = 0, \quad \forall v|_{l_j} \in \mathcal{P}_{q-2}, \tag{60}$$

$$\int_K (i_h^q w - w) v \, dx = 0, \quad \forall v|_K \in \mathcal{Q}_{q-2}, \tag{61}$$

which, according to [9, p.108], satisfies [9, p.101, Lemma A.4]:

$$|i_h^q w - w|_{s,l,\Omega} \leq Ch^{q+1-s} |w|_{q+1,l,\Omega}, \quad 0 \leq s \leq q + 1, \quad 1 < l < \infty. \tag{62}$$

Using the results in [18] the following convergence property for i_h^q can be obtained:

Lemma 4 Assume that $a(x)$ is Lipschitz continuous in Ω , for any $v \in X_h$ and $w \in H^{q+2}(\Omega)$, $q \geq 1$, we have

$$|(a(x, y)\partial_y(i_h^q w - w), \partial_y v)| + |(a(x, y)\partial_x(i_h^q w - w), \partial_x v)| \leq Ch^{q+1}\|w\|_{q+2}|v|_1. \tag{63}$$

And,

$$|(a(x, y)\partial_x(i_h^q w - w), \partial_y v)| + |(a(x, y)\partial_y(i_h^q w - w), \partial_x v)| \leq Ch^{q+1}\|w\|_{q+2}|v|_1. \tag{64}$$

where $w = 0$ on $\partial\Omega$.

Proof Since (64) comes directly from [18, p. 341, Lemma 3(I)], here we only consider the proof of (63), which follows the proof of [18, Lemma 3(I)].

Firstly we separate $(a(x, y)\partial_x(i_h^q w - w), \partial_x v)$ into two parts:

$$\begin{aligned} & (a(x, y)\partial_x(i_h^q w - w), \partial_x v) \\ &= \sum_K \int_K a(x, y)\partial_x(i_h^q w - w) \partial_x v \, d\mathbf{x} \\ &= \sum_K \int_K a(x_K, y_K)\partial_x(i_h^q w - w) \partial_x v \, d\mathbf{x} \\ & \quad + \sum_K \int_K (a(x, y) - a(x_K, y_K)) \partial_x(i_h^q w - w) \partial_x v \, d\mathbf{x} \\ &\leq \sum_K \int_K a(x_K, y_K)\partial_x(i_h^q w - w) \partial_x v \, d\mathbf{x} + Ch \cdot h^q \|w\|_{q+1} |v|_1 \\ &\equiv I_a + Ch^{q+1} \|w\|_{q+1} |v|_1, \end{aligned} \tag{65}$$

in which (x_K, y_K) denotes the center of element K , and the continuity of $a(x, y)$ has been used. To analyse I_a , we need to make use of Lemma 1(I) and Lemma 2(I) of [18]:

$$\int_K \partial_x(i_h^q w - w) \partial_x v \, d\mathbf{x} = O(h^{q+1}) |w|_{k+2, K} |v|_{1, K}, \quad \forall v \in X_h, q \geq 2, \tag{66}$$

$$\int_K \partial_x(i_h^q w - w) \partial_x v \, d\mathbf{x} = O(h^2) |w|_{3, K} |v|_{1, K}, \quad \forall v \in X_h, q = 1. \tag{67}$$

Since all the boundary integral terms appeared in the proof of (66) and (67), i.e. [18, equation (21),(29),(37),(48),(52),(55),(58)], are equal to 0, thus to estimate I_a we simply need to sum up the result for each $K \in \mathcal{T}_h$:

$$\begin{aligned} I_a &= \sum_K \int_K a(x_K, y_K)\partial_x(i_h^q w - w) \partial_x v \, d\mathbf{x} \\ &\leq Ch^{q+1} \sum_K |w|_{k+2, K} |v|_{1, K} \leq Ch^{q+1} \|w\|_{q+2} |v|_1. \end{aligned}$$

Substituting the above result into (65) yields the conclusion. □

Referring to Theorem 2.2.3 and Theorem 2.3.3 in [28, p. 29 and p. 47], for w without boundary restrictions, an estimate of the same integral as in (64) is available:

Lemma 5 Assume that $a(x) \in W^{1,\infty}(\Omega)$, for any $v \in X_h$ and $w \in H^{q+2}(\Omega)$, $q \geq 1$, we have

$$|(a(x, y)\partial_x (i_h^q w - w), \partial_y v)| \leq Ch^{q+\frac{1}{2}} \|w\|_{q+2} |v|_1. \tag{68}$$

Below we denote by $u^n \equiv u(t_n)$ and u_h^n the value of the exact solution and the numerical solution at t_n , respectively. Define $\eta_u^n \equiv u^n - i_h^q u^n$. We hope to obtain the same error order as (64) without the Dirichlet boundary restriction. Here we consider a special case when $a(x, y) = 0$ on $\partial\Omega$, which is the case we shall encounter in the following treatment of the nonlinear term. Referring to the proof of (64) in [18] and the proof of (68) in [28], and with a careful treatment for the boundary term, we are able to derive the desired result.

Lemma 6 Given function $a(x, y) \in W^{1,\infty}(\Omega)$ with $a(x, y) = 0$ on $\partial\Omega$, for any $v \in X_h$ and $w \in H^{q+2}(\Omega)$, $q \geq 1$, we have

$$|(a(x, y)\partial_x (i_h^q w - w), \partial_y v)| \leq Ch^{q+1} \|w\|_{q+2} |v|_1. \tag{69}$$

Proof As in [28], on each rectangular element we define $\bar{a} = \int_K \frac{a}{|K|} dx$, where $|K|$ is the area of K , then the original integral can thus be separated into two parts. Applying $|\bar{a} - a| \leq Ch|a|_{1,\infty}$ and (62) yields

$$\begin{aligned} (a(x, y)\partial_x \eta_w^n, \partial_y v) &= - \sum_K ((\bar{a} - a)\partial_x \eta_w^n, \partial_y v) + \sum_K (\bar{a}\partial_x \eta_w^n, \partial_y v) \\ &\leq Ch^{q+1} |a|_{1,\infty} \|w\|_{q+1} \|\nabla v\| + \sum_K (\bar{a}\partial_x \eta_w^n, \partial_y v). \end{aligned} \tag{70}$$

The second part on the right-hand side can be analysed in the same way as in the proof of (64) in [18]. We separate the Taylor expansion of $\partial_y v$ on the middle point of K into three parts, and analyse them respectively. Only boundary terms require Dirichlet boundary condition of w to achieve the $O(h^{q+1})$ error order:

$$\begin{aligned} \sum_K (\bar{a}\partial_x \eta_w^n, \partial_y v) &\leq Ch^{q+1} \|w\|_{q+2} \|\nabla v\| \\ &+ \left(\sum_{\partial K \cap \partial\Omega_4 \neq \emptyset} \int_{l_4} - \sum_{\partial K \cap \partial\Omega_3 \neq \emptyset} \int_{l_3} \right) (\bar{a}w_{x^{q+1}} E^q(x)v_{x^q}) dx. \end{aligned} \tag{71}$$

Here l_3, l_4 denote the lower and upper boundary of K , respectively. Similarly $\partial\Omega_3, \partial\Omega_4$ represent respectively the lower and upper boundary of Ω . The auxiliary function $E^q(x) = \frac{[(x-x_K)^2 - h_K^2]^q}{2^q}$ is $O(h^{2q})$, in which x_K is the x-coordinate of the middle point of element K , h_K is half of the width of element K in the x-direction. And also, $w_{x^{q+1}} \equiv \frac{\partial^{q+1} w}{\partial x^{q+1}}$. In [28], for w without boundary conditions, applying trace theorem

$$\|\gamma w\|_{q+1,2,\partial K} \|\gamma v\|_{q,2,\partial K} \leq \|w\|_{q+\frac{3}{2},2,K} \|v\|_{q+\frac{1}{2},2,K},$$

inverse estimate

$$\|v\|_{q+\frac{1}{2},2,K} \leq h^{\frac{1}{2}-q} \|v\|_{1,2,K}$$

and Poincaré’s inequality yields

$$\left(\sum_{\partial K \cap \partial\Omega_4 \neq \emptyset} \int_{l_4} - \sum_{\partial K \cap \partial\Omega_3 \neq \emptyset} \int_{l_3} \right) (\bar{a}w_{x^{q+1}} E^q(x)v_{x^q}) dx \leq Ch^{q+\frac{1}{2}} \|w\|_{q+2} \|\nabla v\|.$$

However, with a careful application of the condition that $a = 0$ on $\partial\Omega$ and $a \in W^{1,\infty}$, we have that: on boundary rectangular elements, $|\bar{a}| = O(h)$. Specifically, without loss of generality, we assume edge $l \in \partial K \cap \partial\Omega$ is parallel to x-axis and denote the y-coordinate on l by y_l , then

$$|\bar{a}| = \left| \int_K \frac{a}{|K|} \, dx \right| = \left| \int_K \frac{\int_{y_l}^y \partial_y a \, ds}{|K|} \, dx \right| \leq |a|_{1,\infty} \int_K \frac{\int_{y_l}^y ds}{|K|} \, dx \leq Ch|a|_{1,\infty}.$$

Thus

$$\left(\sum_{\partial K \cap \partial\Omega_4 \neq \emptyset} \int_{l_4} - \sum_{\partial K \cap \partial\Omega_3 \neq \emptyset} \int_{l_3} \right) (\bar{a} w_{x^{q+1}} E^q(x) v_{x^q}) \, dx \leq Ch^{q+\frac{3}{2}} |a|_{1,\infty} \|w\|_{q+2} \|\nabla v\|.$$

Substituting this estimate into (71), we thus obtain $O(h^{q+1})$ estimate. □

Since our exact solutions to Eqs. (6) and (7) satisfy Neumann boundary condition $\partial_n w = \partial_n u = 0$ on $\partial\Omega$ and Ω is aligned with x–y axis, for $a(x, y) = \frac{\partial_x u^n \partial_y u^n}{(1+|\nabla u^n|^2)^2}$ we have $a|_{\partial\Omega} = 0$. Thus Lemma 6 leads to the the following conclusion needed in the proof of Theorem 4.

Corollary 1 *Let Ω be a rectangular domain aligned with x–y axis. Given function $a(x, y) = \frac{\partial_x u^n \partial_y u^n}{(1+|\nabla u^n|^2)^2}$, with $u \in L^\infty(0, T; W^{2,\infty})$ and $\partial_n u = 0$ on $\partial\Omega$, for any $v \in X_h$ and $w \in H^{q+2}(\Omega)$, $q \geq 1$, we have*

$$|(a(x, y) \partial_x (i_h^q w - w), \partial_y v)| \leq Ch^{q+1} \|w\|_{q+2} |v|. \tag{72}$$

Here we introduce the standard Lagrange interpolation operator defined in [2, p. 77, (3.3.2)]: For the finite element $(\Omega, \mathcal{Q}_q, \mathcal{N})$, the basis $\{\phi_i\}$ dual to \mathcal{N} , and $\{N_i\} \in \mathcal{N}$, denote $I_h : C^0(\bar{\Omega}) \rightarrow \mathcal{Q}_q(K)$ the standard interpolation $I_h v \equiv \sum_{i=1}^k N_i(v) \phi_i$. Then (62) and (63) in turn yield:

Lemma 7 *For any $w \in H^{q+2}(\Omega) \cap W^{q+1,\infty}(\Omega)$, we have*

$$\|\nabla (i_h^q w - w)\|_{L^\infty} \leq Ch^{q-1} |w|_{q+1,\infty}, \quad \|\nabla (i_h^q w - R_h w)\| \leq Ch^{q+1} |w|_{q+2},$$

where C is a constant independent of w .

Proof Assume that the maximum element diameter h of mesh T_h is sufficiently small. It is well-known that interpolant I_h has the error estimate $|I_h w - w|_{i,p} \leq Ch^{m-i} |w|_{m,p}$ for any $1 \leq p \leq \infty$ and $m - l - n/p > 0$; see [2, p. 105, Theorem 4.4.4]. Also, from the inverse estimate Lemma 2 we obtain that $\|v\|_{1,\infty} \leq Ch^{-2} \|v\|$. Thus we get

$$\begin{aligned} \|\nabla (i_h^q w - w)\|_{L^\infty} &\leq \|\nabla (i_h^q w - I_h w)\|_{L^\infty} + \|\nabla (I_h w - w)\|_{L^\infty} \\ &\leq Ch^{-1} |i_h^q w - I_h w|_1 + Ch^q |w|_{q+1,\infty} \\ &\leq Ch^{-1} (|i_h^q w - w|_1 + |w - I_h w|_1) + Ch^q |w|_{q+1,\infty} \\ &\leq Ch^{q-1} |w|_{q+1} + Ch^q |w|_{q+1,\infty} \\ &\leq Ch^{q-1} (|w|_{q+1} + |w|_{q+1,\infty}). \end{aligned}$$

Since Ω is bounded, we can bound $|w|_{q+1}$ by $|w|_{q+1,\infty}$, thus the first inequality has been proved. As for the second inequality, we make use of the definition of R_h :

$$\begin{aligned} \|\nabla (i_h^q w - R_h w)\|^2 &= (\nabla (i_h^q w - w), \nabla (i_h^q w - R_h w)) \\ &\leq Ch^{q+1} |w|_{q+2} |i_h^q w - R_h w|_1. \end{aligned}$$

Thus we obtain $\|\nabla(i_h^q w - R_h w)\| \leq Ch^{q+1} |w|_{q+2}$. □

Now we state the main conclusion of this section.

Theorem 4 *Given a rectangular domain Ω that is aligned with x - y axis, and a regular rectangular mesh T_h on Ω , define finite function space X_h as in the beginning of this section and interpolation operator i_h^q as in [9]. Assume that the exact solution pair (u, w) and time step size τ satisfy the assumptions in Theorem 3. In addition, assume that $u \in L^\infty(0, T; H^{q+2}) \cap L^\infty(0, T; W^{2,\infty})$, then we have the following error estimate for scheme (19) and (20):*

$$\|u^n - u_h^n\| + \left(\tau \sum_{m=1}^n \|w^m - w_h^m\|^2 \right)^{\frac{1}{2}} \leq C_{\varepsilon,T} (h^{q+1} + \tau^2). \tag{73}$$

Proof For simplicity, below we denote $f(|\nabla u|^2) \equiv \frac{1}{1+|\nabla u|^2}$. Recalling the analysis in Sect. 2.3, we only need to improve the analysis of the nonlinear terms $(\mathcal{N}^1, \nabla\sigma_u^1)$ and $(\mathcal{N}^{n+1}, \nabla\sigma_u^{n+1})$. Firstly, separate $(\mathcal{N}^1, \nabla\sigma_u^1)$ into two parts as in (56):

$$\begin{aligned} (\mathcal{N}^1, \nabla\sigma_u^1) &\equiv \left(\frac{\nabla u^1}{1+|\nabla u^1|^2} - \frac{\nabla u^0}{1+|\nabla u^0|^2}, \nabla\sigma_u^1 \right) + \left(\frac{\nabla u^0}{1+|\nabla u^0|^2} - \frac{\nabla u_h^0}{1+|\nabla u_h^0|^2}, \nabla\sigma_u^1 \right) \\ &= - \left(\nabla \cdot \left(\frac{\nabla u^1}{1+|\nabla u^1|^2} - \frac{\nabla u^0}{1+|\nabla u^0|^2} \right), \sigma_u^1 \right) + (\tilde{\mathcal{N}}^0, \nabla\sigma_u^1) \\ &\leq C\tau^3 + \frac{C}{\tau} \|\sigma_u^1\|^2 + (\tilde{\mathcal{N}}^0, \nabla\sigma_u^1), \end{aligned}$$

where $\tilde{\mathcal{N}}^0 \equiv \frac{\nabla u^0}{1+|\nabla u^0|^2} - \frac{\nabla u_h^0}{1+|\nabla u_h^0|^2}$. Here we consider the general form of $(\tilde{\mathcal{N}}^0, \nabla\sigma_u^1)$:

$$\begin{aligned} (\tilde{\mathcal{N}}^n, \nabla\sigma_u^{n+1}) &\equiv \left(\frac{\nabla u^n}{1+|\nabla u^n|^2} - \frac{\nabla u_h^n}{1+|\nabla u_h^n|^2}, \nabla\sigma_u^{n+1} \right) \\ &= (f(|\nabla u^n|^2) \nabla e_u^n, \nabla\sigma_u^{n+1}) + ([f(|\nabla u^n|^2) - f(|\nabla u_h^n|^2)] \nabla u_h^n, \nabla\sigma_u^{n+1}) \\ &= (f(|\nabla u^n|^2) \nabla e_u^n, \nabla\sigma_u^{n+1}) - ([f(|\nabla u^n|^2) - f(|\nabla u_h^n|^2)] \nabla e_u^n, \nabla\sigma_u^{n+1}) \\ &\quad + ([f(|\nabla u^n|^2) - f(|\nabla u_h^n|^2)] \nabla u^n, \nabla\sigma_u^{n+1}) \\ &\equiv \text{(III)} + \text{(IV)} + \text{(V)}. \end{aligned}$$

Separating (3) into three parts, applying (63) and Lemma 7, we have

$$\begin{aligned} \text{(III)} &= (f(|\nabla u^n|^2) \nabla \eta_u^n, \nabla\sigma_u^{n+1}) + (f(|\nabla u^n|^2) \nabla (i_h^q - R_h) u^n, \nabla\sigma_u^{n+1}) \\ &\quad + (f(|\nabla u^n|^2) \nabla \sigma_u^n, \nabla\sigma_u^{n+1}) \\ &\leq Ch^{q+1} |u^n|_{q+2} \|\nabla\sigma_u^{n+1}\| + Ch^{q+1} |u^n|_{q+2} \|\nabla\sigma_u^{n+1}\| + \|\nabla\sigma_u^{n+1}\| \|\nabla\sigma_u^n\| \\ &\leq Ch^{2(q+1)} + C \|\nabla\sigma_u^{n+1}\|^2 + C \|\nabla\sigma_u^n\|^2. \end{aligned}$$

Similarly, we separate the ∇e_u^n term in (IV) into three parts

$$\begin{aligned}
 \text{(IV)} &= \left(\frac{\nabla(u^n + u_h^n) \cdot \nabla e_u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla \eta_u^n, \nabla \sigma_u^{n+1} \right) \\
 &\quad - ([f(|\nabla u^n|^2) - f(|\nabla u_h^n|^2)] \nabla(i_h^q - R_h)u^n, \nabla \sigma_u^{n+1}) \\
 &\quad - ([f(|\nabla u^n|^2) - f(|\nabla u_h^n|^2)] \nabla \sigma_u^n, \nabla \sigma_u^{n+1}) \\
 &\leq \left(\frac{\nabla(u^n + u_h^n) \cdot \nabla e_u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla \eta_u^n, \nabla \sigma_u^{n+1} \right) + 2 \|\nabla \sigma_u^{n+1}\| \|\nabla \sigma_u^n\| \\
 &\quad + 2 \|\nabla(i_h^q - R_h)u^n\| \|\nabla \sigma_u^{n+1}\|.
 \end{aligned}$$

Again, we separate the first term on the right-hand side into three parts. A combination of the inequality $\frac{|a|+|b|}{(1+a^2)(1+b^2)} < 1$ and (62) yields:

$$\begin{aligned}
 \text{(IV)} &\leq \left(\frac{\nabla(u^n + u_h^n) \cdot \nabla \eta_u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla \eta_u^n, \nabla \sigma_u^{n+1} \right) \\
 &\quad + \left(\frac{\nabla(u^n + u_h^n) \cdot \nabla(i_h^q - R_h)u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla \eta_u^n, \nabla \sigma_u^{n+1} \right) \\
 &\quad + \left(\frac{\nabla(u^n + u_h^n) \cdot \nabla \sigma_u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla \eta_u^n, \nabla \sigma_u^{n+1} \right) \\
 &\quad + 2 \|\nabla \sigma_u^{n+1}\| \|\nabla \sigma_u^n\| + 2 \|\nabla(i_h^q - R_h)u^n\| \|\nabla \sigma_u^{n+1}\| \\
 &\leq \|\nabla \eta_u^n\|_{L^4}^2 \|\nabla \sigma_u^{n+1}\| + (2 + \|\nabla \eta_u^n\|_{L^\infty}) \|\nabla(i_h^q - R_h)u^n\| \|\nabla \sigma_u^{n+1}\| \\
 &\quad + (\|\nabla \eta_u^n\|_{L^\infty} + 2) \|\nabla \sigma_u^{n+1}\| \|\nabla \sigma_u^n\| \\
 &\leq C(h^{4q} + h^{2(q+1)}) + C(2 + h^{q-1}) \|\nabla \sigma_u^{n+1}\|^2 + C \|\nabla \sigma_u^n\|^2 \\
 &\leq C(h^{4q} + h^{2(q+1)}) + C \|\nabla \sigma_u^{n+1}\|^2 + C \|\nabla \sigma_u^n\|^2,
 \end{aligned} \tag{74}$$

in which we have made use of (62) and 2-D Sobolev embedding (from H^{q+2} into $W^{q+1,4}$) in the third step:

$$\|\nabla \eta_u^n\|_{L^4}^2 \leq Ch^{2(q+1-1)}|u|_{q+1,4} \leq Ch^{2q}|u|_{q+2,2}.$$

In addition, we have used the restriction $h < 1$, raised in Sect. 2.3, in the derivation of (74). Since $4q \geq 2(q + 1)$ for $q \geq 1$, we have obtained $O(h^{q+1})$ spatial convergence for (IV).

As for term (V), we separate it into three parts:

$$\begin{aligned}
 \text{(V)} &= - \left(\frac{\nabla(u^n + u_h^n) \cdot \nabla \eta_u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla u^n, \nabla \sigma_u^{n+1} \right) \\
 &\quad - \left(\frac{\nabla(u^n + u_h^n) \cdot \nabla(i_h^q - R_h)u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla u^n, \nabla \sigma_u^{n+1} \right) \\
 &\quad - \left(\frac{\nabla(u^n + u_h^n) \cdot \nabla \sigma_u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla u^n, \nabla \sigma_u^{n+1} \right) \\
 &\equiv (V_1) + (V_2) + (V_3).
 \end{aligned}$$

The last two terms can be analysed similarly:

$$\begin{aligned} (V_2) &\leq \|\nabla u^n\|_{L^\infty} \| \nabla (i_h^q - R_h) u^n \| \|\nabla \sigma_u^{n+1}\| \leq C \|\nabla \sigma_u^{n+1}\|^2 + Ch^{2(q+1)}, \\ (V_3) &\leq \|\nabla u^n\|_{L^\infty} \|\nabla \sigma_u^{n+1}\| \|\nabla \sigma_u^n\| \leq C \|\nabla \sigma_u^{n+1}\|^2 + C \|\nabla \sigma_u^n\|^2. \end{aligned}$$

In order to apply (63) and Corollary 1 to (V₁), we separate the term twice so as to obtain a continuous coefficient function as $a(x, y)$ in Corollary 1.

$$\begin{aligned} (V_1) &= - \left(\frac{2\nabla u^n \cdot \nabla \eta_u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla u^n, \nabla \sigma_u^{n+1} \right) \\ &\quad + \left(\frac{\nabla e_u^n \cdot \nabla \eta_u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla u^n, \nabla \sigma_u^{n+1} \right) \\ &= - \left(\frac{2\nabla u^n \cdot \nabla \eta_u^n}{(1 + |\nabla u^n|^2)^2} \nabla u^n, \nabla \sigma_u^{n+1} \right) \\ &\quad - \left(\frac{2(\nabla u^n \cdot \nabla \eta_u^n)(\nabla(u^n + u_h^n) \cdot \nabla e_u^n)}{(1 + |\nabla u^n|^2)^2(1 + |\nabla u_h^n|^2)} \nabla u^n, \nabla \sigma_u^{n+1} \right) \\ &\quad + \left(\frac{\nabla e_u^n \cdot \nabla \eta_u^n}{(1 + |\nabla u^n|^2)(1 + |\nabla u_h^n|^2)} \nabla u^n, \nabla \sigma_u^{n+1} \right). \end{aligned}$$

Note that the last two terms above can be analysed in the same way as (74), for brevity we only present the results here. Also, for the convenience of later analysis, we write the first term component wise:

$$\begin{aligned} (V_1) &\leq (a_{12}(x, y) \partial_x \eta_u^n, \partial_y \sigma_u^{n+1}) + (a_{12}(x, y) \partial_y \eta_u^n, \partial_x \sigma_u^{n+1}) + (a_1(x, y) \partial_x \eta_u^n, \partial_x \sigma_u^{n+1}) \\ &\quad + (a_2(x, y) \partial_y \eta_u^n, \partial_y \sigma_u^{n+1}) + C(h^{4q} + h^{2(q+1)}) + C \|\nabla \sigma_u^{n+1}\|^2 + C \|\nabla \sigma_u^n\|^2. \end{aligned}$$

Here $a_1(x, y) = \frac{\partial_x u^n \partial_x u^n}{(1+|\nabla u^n|^2)^2}$, $a_2(x, y) = \frac{\partial_y u^n \partial_y u^n}{(1+|\nabla u^n|^2)^2}$, and $a_{12}(x, y) = \frac{\partial_x u^n \partial_y u^n}{(1+|\nabla u^n|^2)^2}$. Notice that $a_{12}(x, y)$ is exactly the coefficient function in Corollary 1. Therefore, applying Corollary 1 to the first two terms and (63) to the third and fourth term leads to

$$(V_1) \leq C(h^{2(q+1)} + h^{4q}) + C \|\nabla \sigma_u^{n+1}\|^2 + C \|\nabla \sigma_u^n\|^2.$$

Therefore, we have proved that the nonlinear term in the initial step (17) has $(q + 1)$ th order spatial convergence.

As for the nonlinear term $(\mathcal{N}^{n+1}, \nabla \sigma_u^{n+1})$ with $n \geq 1$, we rewrite the term as follows:

$$\begin{aligned} &(\mathcal{N}^{n+1}, \nabla \sigma_u^{n+1}) \\ &= \left(\frac{\nabla u^{n+1}}{1 + |\nabla u^{n+1}|^2} - \frac{\nabla(2u_h^n - u_h^{n-1})}{1 + |\nabla(2u_h^n - u_h^{n-1})|^2} \right), \nabla \sigma_u^{n+1} \\ &= (f(|\nabla u^{n+1}|^2) \nabla u^{n+1} - f(|\nabla(2u^n - u^{n-1})|^2) \nabla(2u^n - u^{n-1}), \nabla \sigma_u^{n+1}) \\ &\quad + (f(|\nabla(2u^n - u^{n-1})|^2) \nabla(2u^n - u^{n-1}) \\ &\quad - f(|\nabla(2u_h^n - u_h^{n-1})|^2) \nabla(2u_h^n - u_h^{n-1}), \nabla \sigma_u^{n+1}) \\ &\equiv (VI) + (VII). \end{aligned}$$

Notice that (VI) can be estimated using Lemma 3:

$$\begin{aligned} \text{(VI)} &\leq \|\nabla(u^{n+1} - 2u^n + u^{n-1})\| \|\nabla\sigma_u^{n+1}\| \\ &\leq C\tau^3 \int_{t_{n-1}}^{t_{n+1}} \|\nabla\partial_{tt}u\|^2 dt + C \|\nabla\sigma_u^{n+1}\|^2. \end{aligned}$$

And (VII) can be analysed as $(\tilde{\mathcal{N}}^n, \nabla\sigma_u^{n+1})$ and be bounded as

$$\text{(VII)} \leq C \left(\|\nabla\sigma_u^{n+1}\|^2 + \|\nabla\sigma_u^n\|^2 + \|\nabla\sigma_u^{n-1}\|^2 \right) + C \left(h^{2(q+1)} + h^{4q} \right).$$

Therefore, the nonlinear terms in the two schemes both have $(q + 1)$ th order spatial convergence when $q \geq 1$. □

Remark 3 In this section we have obtained the optimal convergence analysis for scheme (17)–(20) on rectangular mesh, which is based on the previous works by [18, 28]. One should notice that (1) plays a crucial role in obtaining the optimal convergence order. On the other hand, the convergence analysis for the case when T_h is a triangulation is still ongoing. Recently, Yan et al. has obtained the optimal error estimates for a class of linear fourth-order elliptic problems in [30], using a super-closeness relation between the numerical solution and the Ritz projection of the exact solution. And we hope to extend the results to the MBE equations.

4 Numerical Results

4.1 Convergence Test

In this subsection we present some numerical tests to check the theoretical convergence of the proposed scheme (19)–(20). Firstly, we set $\Omega = [0, 1]^2$, $T = 1$ and $\varepsilon^2 = 0.05$. The exact solution is given by

$$u_e(x, y, t) = \cos(\pi x) \cos(\pi y)e^{-t}. \tag{75}$$

Next, in order to satisfy the PDE (4) and boundary conditions (2), we add an artificial, time-dependent forcing term on the right hand side:

$$\begin{aligned} \partial_t u_e + \nabla \cdot \left(\frac{\nabla u_e}{1 + |\nabla u_e|^2} \right) - \varepsilon^2 \Delta w_e &= g, \quad (x, y, t) \in \Omega \times (0, T], \\ \text{with } g &= (-1 + 4\pi^4 \varepsilon^2)u - \frac{4\pi^2 u}{2 + \pi^2 e^{-2t} [1 - \cos(2\pi x) \cos(2\pi y)]} \\ &+ \frac{4\pi^4 e^{-2t} u}{[1 - \cos(2\pi x) \cos(2\pi y)]^2} [\cos(2\pi x) + \cos(2\pi y) - 2 \cos(2\pi x) \cos(2\pi y)]. \end{aligned} \tag{76}$$

In the building of finite elements, we use both \mathcal{P}_1 and \mathcal{P}_2 elements on uniform meshes with grid size $h = 1/32, 1/64, 1/128, 1/256$. The L^2 -norm of errors $\|u_e - u_h\|$ are recorded at $T = 1$. As for the time step, we firstly set $\tau = h^2$, so that the spatial error counts the most part in $\|u_e - u_h\|$ and we can see that both kinds of elements achieve full order spatial convergence. Table 1 shows the L^2 error and convergence order for our proposed scheme under this condition. As we can see from the table, second and third order accuracies are observed for the spatial approximations with \mathcal{P}_1 and \mathcal{P}_2 elements, respectively. Next we set $\tau = h/2$, and the results are showed in Table 2, from which we can observe a clear second order convergence for the temporal approximation.

Table 1 Scheme convergence with triangular mesh

Spatial convergence				Temporal convergence			
$\tau = h^2$	h	L2 error	Order	$\tau = h/2$	h	L2 error	Order
\mathcal{P}_1 element	1/32	1.04623E-03	1.95662	\mathcal{P}_1 element	1/32	6.40797E-04	2.04527
	1/64	2.63590E-04	1.98883		1/64	1.58241E-04	2.01774
	1/128	6.60295E-05	1.99711		1/128	3.93490E-05	2.00772
	1/256	1.65159E-05	1.99925		1/256	9.81302E-06	2.00356
\mathcal{P}_2 element	1/32	6.36500E-06	3.89642	\mathcal{P}_2 element	1/32	1.58738E-03	2.0209
	1/64	4.95739E-07	3.68251		1/64	3.94786E-04	2.0075
	1/128	4.83008E-08	3.3595		1/128	9.84749E-05	2.0032
	1/256	5.53185E-09	3.12621				

Table 2 Outputs of Q1 MFEM

Mesh	Error			
	$\ u^N - u_h^N\ $	$\ w^N - w_h^N\ $	$\ \nabla(u^N - u_h^N)\ $	$\ \nabla(w^N - w_h^N)\ $
$h = \tau = \frac{1}{32}$	3.049E-2	6.131E-1	1.383E-1	3.058E0
$h = \tau = \frac{1}{64}$	8.483E-3	1.720E-1	3.968E-2	9.182E-1
$h = \tau = \frac{1}{128}$	2.194E-3	4.463E-2	1.140E-2	2.615E-1
$2hx = hy = \tau = \frac{1}{32}$	3.018E-2	6.123E-1	1.409E-1	3.129E0
$2hx = hy = \tau = \frac{1}{64}$	8.405E-3	1.716E-1	4.205E-2	9.645E-1
$2hx = hy = \tau = \frac{1}{128}$	2.174E-3	4.455E-2	1.340E-2	2.976E-1

Mesh	Order			
	$\ u^N - u_h^N\ $	$\ w^N - w_h^N\ $	$\ \nabla(u^N - u_h^N)\ $	$\ \nabla(w^N - w_h^N)\ $
$h = \tau = \frac{1}{64}$	1.8457	1.83411	1.80151	1.73553
$h = \tau = \frac{1}{128}$	1.95131	1.94592	1.79974	1.8121
$2hx = hy = \tau = \frac{1}{64}$	1.84433	1.83475	1.74433	1.69806
$2hx = hy = \tau = \frac{1}{128}$	1.95096	1.94602	1.64946	1.69654

On the other hand, we test the Q_1 finite element methods introduced in Sect. 3. Parameters and exact solutions are the same as above, except for $\varepsilon^2 = 0.01$. Mesh sizes are selected to test both the uniform and quasi-uniform cases. Error norm and convergence order are shown in Table 2. The (q+1)-order spatial convergence is clearly observed.

Note that the H^1 semi-norms $\|\nabla(u(T) - u_h^N)\|$ and $\|\nabla(w(T) - w_h^N)\|$ lie between $O(h^q)$ and $O(h^{q+1})$. To find which term has caused this super-convergence, we run the same Q_1 mixed finite element $((w_{n+1}, \psi_h) - (\nabla u_{n+1}, \psi_h) = 0)$ on two linear parabolic systems:

$$\begin{aligned} \text{Test 1 : } & \left(\frac{3u_{n+1} - 4u_n + u_{n-1}}{2\tau}, v_h \right) + \varepsilon^2(\nabla w_{n+1}, v_h) = (g_{n+1}, v_h); \\ \text{Test 2 : } & \left(\frac{3u_{n+1} - 4u_n + u_{n-1}}{2\tau}, v_h \right) + \varepsilon^2(\nabla w_{n+1}, v_h) \\ & + A\tau(\nabla(w_{n+1} - w_n), v_h) = (g_{n+1}, v_h). \end{aligned}$$

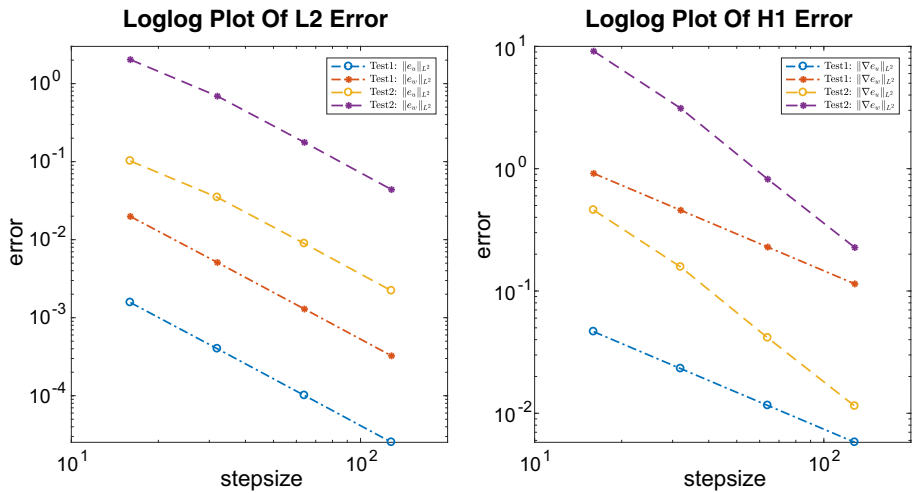


Fig. 1 Convergence results on uniform rectangular mesh with $h = \tau = [\frac{1}{16}, \frac{1}{32}, \frac{1}{64}, \frac{1}{128}]$

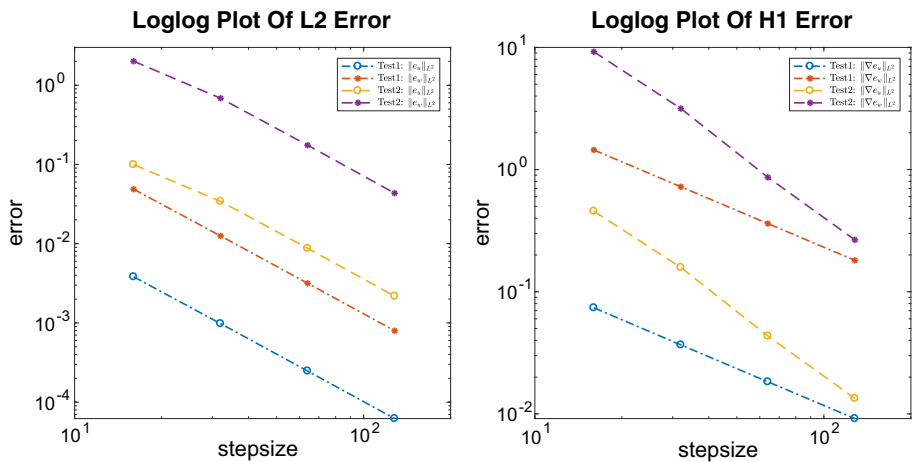


Fig. 2 Convergence results on quasi-uniform rectangular mesh with $2h_x = h_y = \tau = [\frac{1}{16}, \frac{1}{32}, \frac{1}{64}, \frac{1}{128}]$

Firstly we run the experiments with $\varepsilon^2 = 0.01$. The results of $\|u(T) - u_h^N\|$, $\|w(T) - w_h^N\|$, $\|\nabla(u(T) - u_h^N)\|$ and $\|\nabla(w(T) - w_h^N)\|$ on uniform and quasi-uniform meshes are shown in Figs. 1, 2, and Table 3. For the case with $\varepsilon^2 = 0.05$, the results are shown in Table 4. It is observed that the “super-convergence” phenomenon of $\|\nabla(u(T) - u_h^N)\|$ and $\|\nabla(w(T) - w_h^N)\|$ in Test 2 vanishes in Table 4, when $A\tau = \frac{25}{16}\tau < 0.05 = \varepsilon^2$. Therefore, the stabilization term doesn’t lead to a better convergence order. Instead, the difference in convergence order between Test 1 and Test 2 in Table 3 may have been caused by the stabilization error, which is no longer dominant when $A\tau < \varepsilon^2$.

Table 3 Convergence order of Test 1 and Test 2 on rectangular mesh with $\varepsilon^2 = 0.01$

Convergence order	$\ u^N - u_h^N\ $		$\ w^N - w_h^N\ $	
	Test 1	Test 2	Test 1	Test 2
$h = \tau = 1/32$	1.96864	1.55050	1.95773	1.55619
$h = \tau = 1/64$	1.98768	1.96652	1.98285	1.96753
$h = \tau = 1/128$	1.99458	2.00865	1.99225	2.00883
$2h_x = h_y = \tau = 1/32$	1.96553	1.54019	1.96047	1.55459
$2h_x = h_y = \tau = 1/64$	1.98661	1.96576	1.98320	1.96828
$2h_x = h_y = \tau = 1/128$	1.99423	2.00859	1.99229	2.00904
Convergence order	$\ \nabla(u^N - u_h^N)\ $		$\ \nabla(w^N - w_h^N)\ $	
	Test 1	Test 2	Test 1	Test 2
$h = \tau = 1/32$	1.00384	1.54599	1.00020	1.55166
$h = \tau = 1/64$	1.00106	1.92414	1.00008	1.92547
$h = \tau = 1/128$	1.00027	1.85432	1.00003	1.85556
$2h_x = h_y = \tau = 1/32$	1.00966	1.52928	1.00142	1.54355
$2h_x = h_y = \tau = 1/64$	1.00267	1.86558	1.00042	1.86988
$2h_x = h_y = \tau = 1/128$	1.00069	1.69612	1.00011	1.70087

Table 4 Convergence order of Test 1 and Test 2 on rectangular mesh $\varepsilon^2 = 0.05$

Convergence order	$\ u^N - u_h^N\ $		$\ w^N - w_h^N\ $	
	Test 1	Test 2	Test 1	Test 2
$h = \tau = 1/32$	2.13883	1.99674	2.13901	1.99953
$h = \tau = 1/64$	2.04797	1.99928	2.04763	2.00002
$h = \tau = 1/128$	2.01476	1.99986	2.01460	2.00007
$2h_x = h_y = \tau = 1/32$	2.13952	1.99381	2.13947	2.00218
$2h_x = h_y = \tau = 1/64$	2.04904	1.99850	2.04797	2.00060
$2h_x = h_y = \tau = 1/128$	2.01530	1.99964	2.01481	2.00017
Convergence order	$\ \nabla(u^N - u_h^N)\ $		$\ \nabla(w^N - w_h^N)\ $	
	Test 1	Test 2	Test 1	Test 2
$h = \tau = 1/32$	1.83884	1.00185	1.84953	0.99969
$h = \tau = 1/64$	1.41662	1.00047	1.42862	0.99992
$h = \tau = 1/128$	1.14354	1.00012	1.14972	0.99998
$2h_x = h_y = \tau = 1/32$	1.58859	1.00568	1.62324	1.00047
$2h_x = h_y = \tau = 1/64$	1.20992	1.00144	1.23142	1.00011
$2h_x = h_y = \tau = 1/128$	1.05859	1.00036	1.06603	1.00003

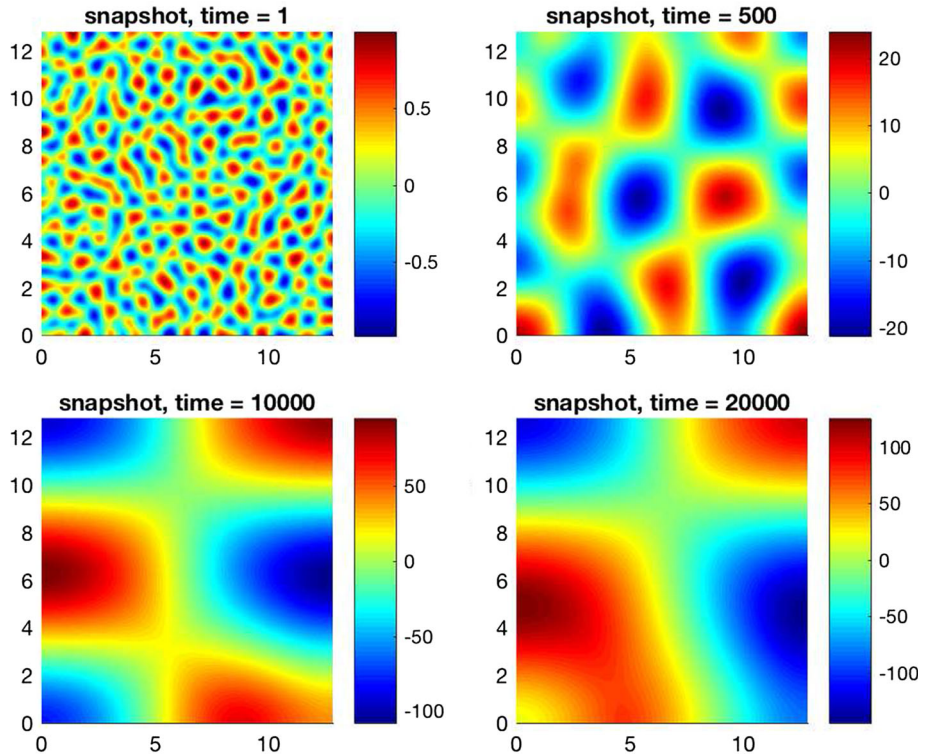


Fig. 3 Snapshots of u at indicated time

4.2 Energy Decay Simulation

In this subsection we aim at simulating the energy decay process of u in the PDE (4). Recall the energy functional (1) proposed at the beginning:

$$E(u) \equiv \int_{\Omega} \left(-\frac{1}{2} \ln(1 + |\nabla u|^2) + \frac{\varepsilon^2}{2} |\Delta u|^2 \right) dx. \tag{77}$$

We set the surface diffusion coefficient parameter as $\varepsilon^2 = 0.005$, and the computational domain is taken to be $\Omega = (0, 12.8)^2$ with time interval $[0, 20,000]$. For $\mathbf{x} \in \Omega$, let $u(\mathbf{x}, 0)$ have random value between $(-0.05, 0.05)$. For the spatial discretization, we use a resolution of $N = 256$, and as for the time step size, we set $\tau = 0.004$ ($t < 200$), $\tau = 0.04$ ($1000 > t \geq 200$), $\tau = 0.08$ ($2000 > t \geq 1000$), $\tau = 0.16$ ($t \geq 2000$). When time step size changes, we reset the last solution as the initial solution and invoke the initialization scheme again with the new time step size.

Snapshots of the solution at $t = 1, 500, 5000, 12,000$ are shown in Fig. 3. In Fig. 4 we show the energy evolution of the scheme (19)–(20) in time interval $[0, 12,250]$.

As is shown in Fig. 4, energy decays at a much faster speed at the beginning. In order to obtain the energy decay rate, below we provide the semi-log plot of the temporal evolution of E in $[1, 400]$ using both P_1 and P_2 elements (Figs. 5, 6).

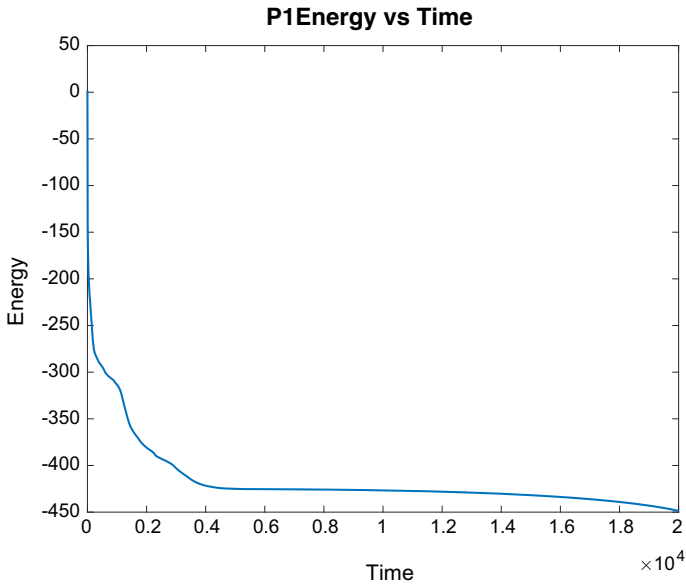


Fig. 4 Plot of the temporal evolution of E for $\varepsilon^2 = 0.005$ using P_1 elements

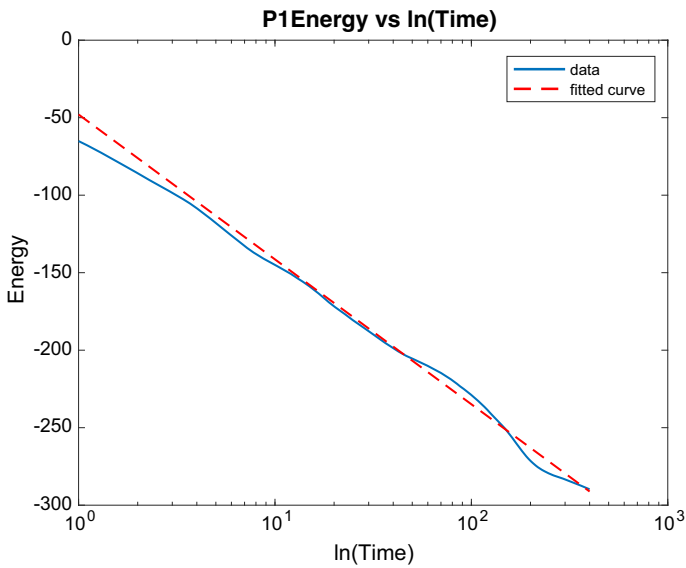


Fig. 5 Semi-log plot of the temporal evolution of E for $\varepsilon^2 = 0.005$ using P_1 elements. The blue line represents the energy obtained by the numerical simulation, while the dashed red line is a least square approximation to the energy data. The fitted line has the form $a \ln(t) + b$, with $a = -40.59$, $b = -47.94$ (Color figure online)

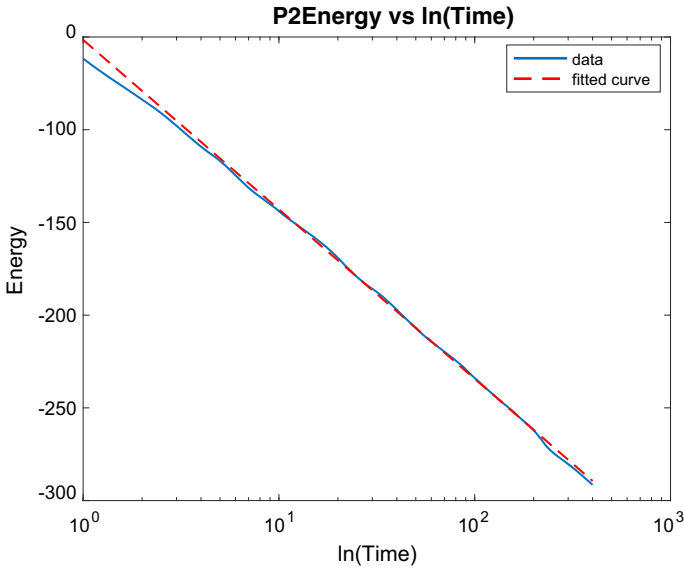


Fig. 6 Semi-log plot of the temporal evolution of E for $\varepsilon^2 = 0.005$ using P_2 elements. The blue line represents the energy obtained by the numerical simulation, while the dashed red line is a least square approximation to the energy data. The fitted line has the form $a \ln(t) + b$, with $a = -39.7$, $b = -51.5$ (Color figure online)

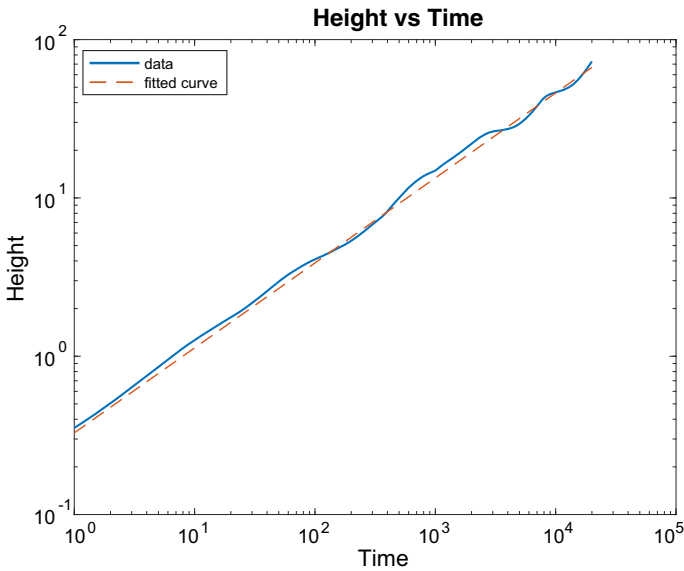


Fig. 7 The log–log plot of the temporal evolution of the average height of u with $\varepsilon^2 = 0.005$ using P_1 elements. The blue line represents the data obtained by the numerical simulation, while the dashed red line is a least square approximation to the height data. The fitted line has the form at^b , with $a = 0.336$, $b = 0.5341$ (Color figure online)

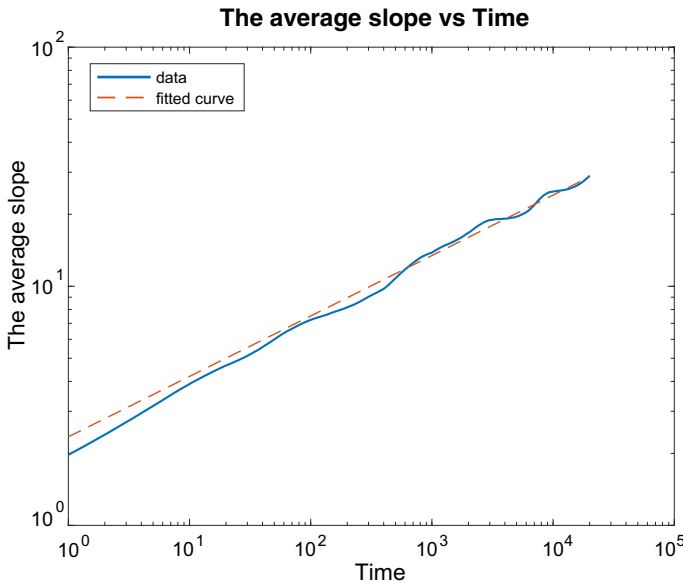


Fig. 8 The log–log plot of the temporal evolution of the average slope of u with $\varepsilon^2 = 0.005$ using P_1 elements. The blue line represents the data obtained by the numerical simulation, while the dashed red line is a least square approximation to the slope data. The fitted line has the form at^b , with $a = 2.348$, $b = 0.2526$ (Color figure online)

4.3 Other Physical Quantities

In this subsection we look at two more physically interesting quantities, i.e., the characteristic height $h(t)$ and the average slope $m(t)$, of which the expressions are as below:

$$h(t) = \sqrt{\frac{1}{|\Omega|} \int_{\Omega} |u(\mathbf{x}, t) - \bar{u}(t)|^2 \, d\mathbf{x}}, \quad \text{with } \bar{u}(t) := \frac{1}{|\Omega|} \int_{\Omega} u(\mathbf{x}, t) \, d\mathbf{x}.$$

$$m(t) = \sqrt{\frac{1}{|\Omega|} \int_{\Omega} |\nabla u(\mathbf{x}, t)|^2 \, d\mathbf{x}}.$$

For the no-slope-selection model (4), one could obtain $h \sim O(t^{1/2})$, $m(t) \sim O(t^{1/4})$, and $E \sim O(-\ln(t))$ as $t \rightarrow \infty$. (See [5, 10, 15, 16] and references therein.) This implies that the characteristic (average) length $\ell(t) := h(t)/m(t) \sim O(t^{1/4})$ as $t \rightarrow \infty$, so that the average length and average slope scale the same with increasing time. Of course, the average mound height $h(t)$ grows faster than the average length $\ell(t)$, which is expected because there is no preferred slope of the height function.

At a theoretical level, the detailed analyses in [12, 13, 16] have indicated (at best) lower bounds for the energy dissipation and, conversely, upper bounds for the average height. On the other hand, the rates quoted as the upper or lower bounds are typically observed for the averaged values of the quantities of interest. To adequately capture the full range of coarsening behaviors, numerical simulations for the coarsening process require short- and long-time accuracy and stability, in addition to high spatial accuracy for small values of ε .

Under the same mesh settings as in Sect. 4.2, here we provide the log–log plot for $h(t)$ and $m(t)$ in time interval $[1, 19,839]$. In fact, these two quantities can be easily measured experimentally. Rigorously, the lower bound for the energy decay rate is of the order of $-\ln(t)$, the upper bounds for the average height and average slope/average length are of the order of $t^{1/2}$, $t^{1/4}$, respectively. Figures 7 and 8 present the log–log plots for the average height versus time, and average slope versus time, respectively. The detailed scaling “exponents” are obtained using least squares fits of the computed data up to $t = 400$. A clear observation of the $t^{1/2}$ and $t^{1/4}$ scaling laws can be made, with different coefficients dependent upon ε , or, equivalently, the domain size, L .

5 Concluding Remarks

In this paper we have presented a second order accurate, linear energy stable numerical scheme for a thin film model without slope selection, with a mixed finite element approximation in space. The unconditional unique solvability and unconditional long time energy stability have been justified at a theoretical level. And also, we obtain an $O(h^q + \tau^2)$ -order convergence analysis in the $\ell^\infty(0, T; H^2)$ norm, when τ is sufficiently small and $h < 1$. In addition, using regular rectangular mesh, we can improve the spatial convergence to $(q + 1)$ th order under current Neumann boundary condition assumptions. Furthermore, the numerical experiments showed that the proposed second-order scheme is able to produce accurate long time numerical results with a reasonable computational cost.

Acknowledgements This work is supported in part by the Grants NSFC 11671098, 11331004, 91630309, a 111 Project B08018 (W. Chen), NSF DMS-1418689 (C. Wang), and Grant 2017110715 by Shanghai University of Finance and Economics (Y. Yan). C. Wang also thanks Shanghai Center for Mathematical Sciences and Shanghai Key Laboratory for Contemporary Applied Mathematics, Fudan University, for support during his visit. We thank the anonymous reviewers for their careful reading of our manuscript and their many insightful comments and suggestions, which lead to substantial improvements of this paper.

References

1. Adams, R.A., Fournier, J.J.F.: Sobolev Spaces, 2nd edn. Academic Press, Singapore (2003)
2. Brenner, S., Scott, R.: The Mathematical Theory of Finite Element Methods, 3rd edn. Springer, New York (2007)
3. Chen, W., Conde, S., Wang, C., Wang, X., Wise, S.: A linear energy stable scheme for a thin film model without slope selection. *J. Sci. Comput.* **52**, 546–562 (2012)
4. Chen, W., Gunzburger, M., Sun, D., Wang, X.: Efficient and long-time accurate second-order methods for the Stokes–Darcy system. *SIAM J. Numer. Anal.* **51**, 2563–2584 (2013)
5. Chen, W., Wang, C., Wang, X., Wise, S.: A linear iteration algorithm for energy stable second order scheme for a thin film model without slope selection. *J. Sci. Comput.* **59**, 574–601 (2014)
6. Chen, W., Wang, Y.: A mixed finite element method for thin film epitaxy. *Numer. Math.* **122**, 771–793 (2012)
7. Eyre, D.: Unconditionally gradient stable time marching the Cahn–Hilliard equation. *MRS. Symp. Proc.* **529**, 39 (1998)
8. Feng, W., Wang, C., Wise, S.M., Zhang, Z.: A second-order energy stable Backward Differentiation Formula method for the epitaxial thin film equation with slope selection. *Numer. Methods Partial Differ. Equ.* (2017, in review)
9. Girault, V., Raviart, P.A.: Finite Element Methods for Navier–Stokes Equations: Theory and Algorithms. Springer, Berlin (2012)
10. Golubović, L.: Interfacial coarsening in epitaxial growth models without slope selection. *Phys. Rev. Lett.* **78**, 90–93 (1997)

11. Ju, L., Li, X., Qiao, Z., Zhang, H.: Energy stability and convergence of exponential time differencing schemes for the epitaxial growth model without slope selection. *Math. Comput.* (2017). <https://doi.org/10.1090/mcom/3262>
12. Kohn, R.: Energy-driven pattern formation. In: Sanz-Sole, M., Soria, J., Varona, J.L., Verdera, J. (eds.) *Proceedings of the International Congress of Mathematicians*, vol. 1, pp. 359–384. European Mathematical Society Publishing House, Madrid (2007)
13. Kohn, R., Yan, X.: Upper bound on the coarsening rate for an epitaxial growth model. *Commun. Pure Appl. Math.* **56**, 1549–1564 (2003)
14. Li, B.: High-order surface relaxation versus the Ehrlich–Schwoebel effect. *Nonlinearity* **19**, 25812603 (2006)
15. Li, B., Liu, J.: Thin film epitaxy with or without slope selection. *Eur. J. Appl. Math.* **14**, 713–743 (2003)
16. Li, B., Liu, J.: Epitaxial growth without slope selection: energetics, coarsening, and dynamic scaling. *J. Nonlinear Sci.* **14**, 429–451 (2004)
17. Li, D., Qiao, Z., Tang, T.: Characterizing the stabilization size for semi-implicit Fourier-spectral method to phase field equations. *SIAM J. Numer. Anal.* **54**, 1653–1681 (2016)
18. Li, J.: Full-order convergence of a mixed finite element method for fourth-order elliptic equations. *J. Math. Anal. Appl.* **230**, 329–349 (1999)
19. Li, X., Qiao, Z., Zhang, H.: Convergence of a fast explicit operator splitting method for the epitaxial growth model with slope selection. *SIAM J. Numer. Anal.* **55**, 265–285 (2017)
20. Qiao, Z., Sun, Z., Zhang, Z.: The stability and convergence of two linearized finite difference schemes for the nonlinear epitaxial growth model. *Numer. Methods Partial Differ. Equ.* **28**, 1893–1915 (2012)
21. Qiao, Z., Sun, Z., Zhang, Z.: Stability and convergence of second-order schemes for the nonlinear epitaxial growth model without slope selection. *Math. Comput.* **84**, 653–674 (2015)
22. Qiao, Z., Wang, C., Wise, S., Zhang, Z.: Error analysis of a finite difference scheme for the epitaxial thin film growth model with slope selection with an improved convergence constant. *Int. J. Numer. Anal. Model.* **14**, 283–305 (2017)
23. Qiao, Z., Zhang, Z., Tang, T.: An adaptive time-stepping strategy for the molecular beam epitaxy models. *SIAM J. Sci. Comput.* **33**, 1395–1414 (2012)
24. Shen, J.: Long time stability and convergence for fully discrete nonlinear galerkin methods. *Appl. Anal.* **38**, 201–229 (1990)
25. Shen, J., Wang, C., Wang, X., Wise, S.: Second-order convex splitting schemes for gradient flows with Ehrlich–Schwoebel type energy: application to thin film epitaxy. *SIAM J. Numer. Anal.* **50**, 105–125 (2012)
26. Wang, C., Wang, X., Wise, S.: Unconditionally stable schemes for equations of thin film epitaxy. *Discrete Contin. Dyn. Syst.* **28**, 405–423 (2010)
27. Xu, C., Tang, T.: Stability analysis of large time-stepping methods for epitaxial growth models. *SIAM J. Numer. Anal.* **44**(4), 1759–1779 (2006)
28. Yan, N.: *Superconvergence Analysis and a Posteriori Error Estimation in Finite Element Methods*. Science Press, New York (2008)
29. Yan, Y., Chen, W., Wang, C., Wise, S.: A second-order energy stable BDF numerical scheme for the Cahn–Hilliard equation. *Commun. Comput. Phys.* **23**, 572–602 (2018)
30. Yan, Y., Li, W., Chen, W., Wang, Y.: Optimal convergence analysis of a mixed finite element method for fourth-order elliptic problems. *Commun. Comput. Phys.* (2017, accepted)
31. Yang, X., Zhao, J., Wang, Q.: Numerical approximations for the molecular beam epitaxial growth model based on the invariant energy quadratization method. *J. Comput. Phys.* **333**, 104–127 (2017)