# Optimal error estimates for the scalar auxiliary variable finite-element schemes for gradient flows

Hongtao Chen[1,2] · Jingjing Mao[1,2] · Jie Shen[1,3]

## Abstract

We carry out in this paper a rigorous error analysis for a finite element discretization of the scalar auxiliary variable (SAV) schemes. The finite-element method we study is a Galerkin method with standard Lagrange elements based on a mixed variational formulation. We derive optimal error estimates for both the first- and second-order SAV schemes with the finite-element method in space.

## 1 Introduction

The scalar auxiliary variable (SAV) approach is recently proposed for solving a large class of gradient flows [17]. The SAV approach is inspired by the IEQ approach (cf. for instance [20,21]), it inherits its essential advantages and fixes many of its shortcomings. In particular, it enjoys the following remarkable advantages: (i) unconditionally energy stable, (ii) at each time step only decoupled systems with constant coefficients need to be solved. It has been already applied to several challenging gradient flows [16], and shown great potential in numerical simulation of complex systems governed by gradient flows.

✉ Jie Shen
shen7@purdue.edu

1   School of Mathematical Sciences, Xiamen University, Xiamen 361005, China

2   Fujian Provincial Key Laboratory of Mathematical Modeling and High-Performance Scientific Computing, Xiamen University, Xiamen 361005, China

3   Department of Mathematics, Purdue University, West Lafayette, IN 47906, USA

While the procedure of constructing unconditionally energy stable SAV schemes is quite standard, one should be aware that energy stability alone does not ensure that the numerical solution will converge to the exact solution, particularly when auxiliary variables, such as in the SAV approach, are introduced. Moreover, the fact that SAV schemes are energy stable with a modified energy further complicates the convergence analysis.

Recently, convergence and error analysis for a first-order semi-discrete SAV scheme is carried out in [15]. The analysis essentially relies on the $H^2$ bound, which implies the $L^\infty$ bound, of the numerical solution. The aim of this paper to derive error estimates for first- and second-order fully discretized SAV schemes with a mixed finite element discretization for the space variables.

Compared with [15], we have to overcome several additional difficulties: (i) The finite element space only belongs to $H^1$, so it is impossible to require $H^2$ bound for the finite element solution, and we have to derive its $L^\infty$ bound by another approach (see the proof of Theorem 1). (ii) A direct procedure as in [15] only leads to a suboptimal error estimate with respect to the finite-element space, since in the analysis the term $\|\Delta_h(\mu - R_h\mu)\|_{H^{-1}}$ will appear, and it can only be bounded by $\|\nabla(\mu - R_h\mu)\|$ where $\Delta_h$ and $R_h$ denote the discrete Laplace and Ritz projection operators, respectively. However, based on the $L^\infty$ bound of the finite element solution (obtained associated with the suboptimal result) a stronger stability for the finite element solutions can be derived in Lemma 11. Then together with the discrete inverse inequality (see Lemma 7) and the commuting property between $\Delta_h$ and the projection operator, we can obtain optimal error estimates with respect to the given finite element space (see Theorem 2). (iii) The analysis in [15] is only valid for first-order SAV schemes, analysis for a second order SAV method is much more complicated. Our analysis can be also applied to second order scheme and produce the optimal error estimates (see Theorem 4).

Although some convergence and error analysis are available for fully implicit (such as backward Euler) [6,7,10] or nonlinearly implicit (such as convex splitting) [2,8] schemes without restrictive assumptions on the free energy, most of the convergence and error analysis for linearly implicit (such as semi-implicit or stabilized semi-implicit) [8,12] are based on a uniform Lipschitz assumption, i.e.,

$$|F'(x) - F'(y)| \le L|x - y|, \quad \forall x, y \in \mathbb{R}$$

where $F(u)$ is the nonlinear free energy density defined below. However, this assumption greatly limits its range of applicability as even the usual double-well potential does not satisfy the assumption. Our error analysis, as in [15], does not require the Lipschitz assumption.

The rest of the paper is organized as follows. In Sect. 2, we present a fully discrete SAV scheme and prove it is unconditionally energy stable with the modified energy. In Sect. 3, we recall the wellposedness and regularity results of the $H^{-1}$ gradient flows, and present some useful properties about the discrete Laplace operator $\Delta_h$. It is followed by the $L^\infty$ bound and error estimates for finite element solutions of the first order scheme will be derived in Sect. 4. In Sect. 5 the method will be extended to second order stable scheme, and the corresponding error estimates will be also analyzed. Some numerical experiments are presented in Sect. 6.

## 2 Fully discrete FEM SAV scheme and its stability

We start by introduce some notations to be used throughout the paper, see, e.g. [3]. We assume $\Omega \in \mathbb{R}^n (n = 2, 3)$ is a bounded Lipschitz domain. Then we denote the spaces $L^p(\Omega)$ associated with the $L^p$ norm $\|u\|_{L^p} := \left(\int_\Omega |u(x)|^p dx\right)^{1/p}$, and $L^p([0, T])$ associated with the $L^p$ norm $\|u\|_{L^p} := \left(\int_0^T |u(t)|^p dt\right)^{1/p}$. Both of them use the abbreviation $L^p$ as they can be easily distinguished by the context. Here we also introduce the space $L^\infty(\Omega)$ and $W_\infty^s(\Omega)$ with

$$\|v\|_{L^\infty} = \sup_{x \in \Omega} |v(x)|, \quad \|v\|_{W_\infty^s} = \max_{|\alpha| \leq s} \|D^\alpha v\|_{L^\infty}$$

and the space $L^\infty([0, T])$ and $W_\infty^s([0, T])$ with similar norms. The Sobolev spaces $H^s$ with the order s will also be used, and the norms are denoted by $\|\cdot\|_s$. The space $L^p(0, T; V)$ represents the $L^p$ space on the interval $(0, T)$ with values in the function space $V$. We denote by $(\cdot, \cdot)$ the inner product in $L^2$ and the $L^2$ norm without subscript. We also denote by C a general constant independent of mesh size $h$ and time step $\Delta t$.

Let $F(u)$ be a nonlinear free energy density, we focus on a typical energy functional $E[u(x)]$ given by

$$E[u] = \int_\Omega \left(\frac{\lambda}{2}u^2 + \frac{1}{2}|\nabla u|^2\right) dx + E_1[u],$$

where we assume that $\lambda \geq 0$ and $E_1[u] = \int_\Omega F(u)dx > -c_0$. Consider the gradient flow

$$\frac{\partial u}{\partial t} = G(-\Delta u + \lambda u + g(u)) \tag{1}$$

where $G = -I$ for $L^2$ gradient flow, $G = \Delta$ for the $H^{-1}$ gradient flow and $g(u) = F'(u)$. As an example, when $E_1[u] = \int_\Omega \alpha(1 - u^2)^2 dx$, the two gradient flows are the celebrated Allen–Cahn and Cahn–Hilliard equations [1,4]. The Eq. (1) is supplemented with the initial condition $u(x, 0) = u^0(x)$ and the homogeneous boundary condition

$$u|_{\partial\Omega} = 0, \quad \text{or} \quad \frac{\partial u}{\partial n} = 0, \quad \text{if } G = -I;$$

$$\frac{\partial u}{\partial n} = \frac{\partial \mu}{\partial n} = 0, \quad \text{if } G = \Delta, \tag{2}$$

or periodic boundary conditions where $\mu = \frac{\delta E}{\delta u}$.

Let $C_0 > c_0$ so that $E_1[u] + C_0 > 0$. Without loss of generality, we substitute $E_1$ with $E_1 + C_0$ without changing the gradient flow. In this setting, $E_1(u)$ always have a positive lower bound $C_0 - c_0$ for any $u$, which we still denote as $C_0$. In the SAV approach, we introduce a scalar variable $r(t) = \sqrt{E_1[u]}$, and rewrite (1) as

$$u_t = G\mu, \tag{3}$$

$$\mu = -\Delta u + \lambda u + \frac{r(t)}{\sqrt{E_1[u]}} g(u),  \tag{4}$$

$$r_t = \frac{1}{2\sqrt{E_1[u]}} \int_\Omega g(u) u_t dx.  \tag{5}$$

Below we only consider the $H^{-1}$ gradient flow, the analysis can be easily extended to $L^2$ gradient flow.

Let $\Gamma_h$ be a quasi-uniform partition of the domain $\Omega$ such that $\bar{\Omega} = \cup_{K \in \Gamma_h} K$, $K$ are triangles in the case $d = 2$ and tetrahedrons in the case $d = 3$. Assume that we are given a family $V_h \subset H^1(\Omega)$ which consists of piecewise polynomials based on $\Gamma_h$ with total order of all variables is less or equal to $k$, such that [3,5]

$$\inf_{\chi \in V_h} \{\|v - \chi\| + h\|\nabla(v - \chi)\|\} \le Ch^s \|v\|_s, \quad 1 \le s \le k+1.$$

Then a fully discrete mixed finite element discretization of (3)–(5) is: Find $(u_h^{n+1}, \mu_h^{n+1}, r_h^{n+1}) \in [V_h]^2 \times R$ such that

$$(u_h^{n+1} - u_h^n, v) = -\Delta t(\nabla \mu_h^{n+1}, \nabla v), \quad \forall v \in V_h,  \tag{6}$$

$$(\mu_h^{n+1}, \tau) = (\nabla u_h^{n+1}, \nabla \tau) + \lambda(u_h^{n+1}, \tau) + \frac{r_h^{n+1}}{\sqrt{E_1(u_h^n)}}(g(u_h^n), \tau), \quad \forall \tau \in V_h,  \tag{7}$$

$$r_h^{n+1} - r_h^n = \frac{1}{2\sqrt{E_1(u_h^n)}}(g(u_h^n), u_h^{n+1} - u_h^n),  \tag{8}$$

with $u_h^0$ to be some appropriate approximation of $u(0)$ and $r_h^0 = \sqrt{E_1(u_h^0)}$.

We first derive the unconditional energy stability for the above fully discrete SAV scheme, which plays an important role in subsequent error analysis.

**Lemma 1** *For the $H^{-1}$ gradient flow and all $N \le T/\Delta t$, we have*

$$\max_{0 \le n \le N} \frac{\lambda}{2}\|u_h^n\|^2 + \frac{1}{2}\|\nabla u_h^n\|^2 + (r_h^n)^2 + \sum_{n=0}^{N-1} \left( \frac{\lambda}{2}\|u_h^{n+1} - u_h^n\|^2 + \frac{1}{2}\|\nabla(u_h^{n+1} - u_h^n)\|^2 \right.$$

$$\left. + (r_h^{n+1} - r_h^n)^2 + \Delta t\|\nabla \mu_h^{n+1}\|^2 \right) \le \frac{\lambda}{2}\|u_h^0\|^2 + \frac{1}{2}\|\nabla u_h^0\|^2 + (r_h^0)^2.  \tag{9}$$

**Proof** Taking the inner product of (6) and (7) with $\mu_h^{n+1}$ and $u_h^{n+1} - u_h^n$ respectively, and multiplying (8) by $2r_h^{n+1}$, we derive that the fully discrete SAV scheme satisfies the following discrete energy law:

$$\frac{\lambda}{2}\|u_h^{n+1}\|^2 + \frac{1}{2}\|\nabla u_h^{n+1}\|^2 + (r_h^{n+1})^2 - \frac{\lambda}{2}\|u_h^n\|^2 - \frac{1}{2}\|\nabla u_h^n\|^2 - (r_h^n)^2$$

$$+\frac{\lambda}{2}\|u_h^{n+1} - u_h^n\|^2 + \frac{1}{2}\|\nabla(u_h^{n+1} - u_h^n)\|^2 + (r_h^{n+1} - r_h^n)^2 = -\Delta t\|\nabla \mu_h^{n+1}\|^2 \le 0$$

from which (9) follows immediately. □

Hence, by (9) the SAV is unconditionally energy stable with the modified energy

$$\frac{\lambda}{2}\|u_h^n\|^2 + \frac{1}{2}\|\nabla u_h^n\|^2 + (r_h^n)^2.$$

An important fact is that the SAV scheme is easy to implement. Indeed, if we define $\triangle_h u_h \in V_h$ with $\int_\Omega \triangle_h u_h dx = 0$ satisfying

$$(\triangle_h u_h, v) = -(\nabla u_h, \nabla v), \quad \forall v \in V_h, \tag{10}$$

and the $L^2$ projection by $P_h : L^2 \to V_h$ and $b_h^n = \frac{P_h g(u_h^n)}{\sqrt{E_1(u_h^n)}}$, then eliminating $r_h^{n+1}$ and $\mu_h^{n+1}$ from (6)–(8) to obtain

$$(u_h^{n+1} - u_h^n, v) = \Delta t \triangle_h(-\triangle_h u_h^{n+1} + \lambda u_h^{n+1} + b_h^n(r_h^n + \frac{1}{2}(b_h^n, u_h^{n+1} - u_h^n)), v).$$

Then the above equation can be written as

$$(I - \lambda \Delta t \triangle_h + \Delta t \triangle_h^2)u_h^{n+1} - \frac{\Delta t}{2}\triangle_h b_h^n(b_h^n, u_h^{n+1}) = u_h^n + \Delta t r_h^n \triangle_h b_h^n - \frac{\Delta t}{2}\triangle_h b_h^n(b_h^n, u_h^n), \tag{11}$$

which can be solved using the Sherman–Morrison–Woodbury formula:

$$(A + UV^T)^{-1} = A^{-1} - A^{-1}U(I + V^T A^{-1}U)^{-1}V^T A^{-1}. \tag{12}$$

That is, if we denote the righthand side of (11) by $c_h^n$. Multiplying both sides of (11) with $(I - \lambda \Delta t \triangle_h + \Delta t \triangle_h^2)^{-1}$ and taking the inner product with $b_h^n$, we obtain

$$(b_h^n, u_h^{n+1}) + \frac{\Delta t}{2}\gamma_h^n(b_h^n, u_h^{n+1}) = (b_h^n, (I - \lambda \Delta t \triangle_h + \Delta t \triangle_h^2)^{-1}c_h^n)$$

where we have from $\triangle_h$ is negative definite

$$\gamma_h^n = -(b_h^n, (I - \lambda \Delta t \triangle_h + \Delta t \triangle_h^2)^{-1}\triangle_h b_h^n) = (b_h^n, (-\triangle_h^{-1} + \lambda \Delta t - \Delta t \triangle_h)^{-1}b_h^n) > 0.$$

Therefore, we have

$$(b_h^n, u_h^{n+1}) = \frac{(b_h^n, (I - \lambda \Delta t \triangle_h + \Delta t \triangle_h^2)^{-1}c_h^n)}{1 + \frac{\Delta t}{2}\gamma_h^n}. \tag{13}$$

To summarize, we implement as follows:

(i) Compute $b_h^n$ and $c_h^n$;
(ii) Compute $(b_h^n, u_h^{n+1})$ from (13);
(iii) Compute $u_h^{n+1}$ from (11).

Note that in (ii) and (iii) we only need to solve twice a linear equation with constant coefficients of the form

$$(I - \lambda \Delta t \Delta_h + \Delta t \Delta_h^2)x = b. \tag{14}$$

Therefore, the above procedure is extremely efficient.

## 3 Well-posedness and properties

We assume that $F \in C^3(R)$. We first recall the existence, uniqueness and regularity results for $H^{-1}$ gradient flows, and there are similar results for $L^2$ gradient flows [15, Proposition 2.1].

**Proposition 1** [18] *For $H^{-1}$ gradient flow, if $u^0 \in L^2$, there exist constants $p_0, C > 0$ such that*

$$F''(s) = g'(s) \geq -C, \quad sg(s) \geq b|s|^{p_0} - c.$$

*Then there exists a unique solution $u$ for* (1) *such that*

$$u \in L^2(0, T; H^2) \cap L^{p_0}(0, T; H^{p_0}) \cap C([0, T]; L^2).$$

*Moreover, if $u^0 \in H^2$, and*

$$|g'(x)| < C(|x|^p + 1), \text{for any } p > 0 \text{ if } d = 1, 2; \quad 0 < p < 4 \text{ if } d = 3; \tag{15}$$
$$|g''(x)| < C(|x|^q + 1), \text{for any } q > 0 \text{ if } d = 1, 2; \quad 0 < q < 3 \text{ if } d = 3; \tag{16}$$

*there exists a unique solution $u$ for* (1) *in the space $L^2(0, T; H^4) \cap C([0, T]; H^2)$.*

We introduce the following negative norms [cf. (5.6) in [19]]:

$$\|v\|_{-s} = \sup\{\frac{(v, \phi)}{\|\phi\|_s}; \phi \in H^s\}, \quad \text{for } s \geq 0 \text{ integer.}$$

We prove below some properties of the discrete Laplace operator $\Delta_h$, which will be frequently used in the subsequent analysis [11,13].

**Lemma 2** *For the operator $\Delta_h$, we have the following bound*

$$\|\Delta_h v\|_{-1} \leq C\|\nabla v\|. \tag{17}$$

**Proof** Since

$$\|P_h v\|_1 \leq \|v\|_1 + \|v - P_h v\|_1 \leq C\|v\|_1, \tag{18}$$

by the definition (10) of the operator $\Delta_h$, (17) immediately follows from

$$
\begin{aligned}
\|\Delta_h v\|_{-1} &\leq \sup_{0 \neq \varphi \in H^1} \frac{(\Delta_h v, \varphi)}{\|\varphi\|_{H^1}} = \sup_{0 \neq \varphi \in H^1} \frac{(\Delta_h v, P_h \varphi)}{\|\varphi\|_{H^1}} \\
&= \sup_{0 \neq \varphi \in H^1} -\frac{(\nabla v, \nabla P_h \varphi)}{\|\varphi\|_{H^1}} \leq C\|\nabla v\|.
\end{aligned}
$$

□

The main ingredient in [15] is to derive a bound on the $H^2$ norm of the solution which automatically provides a $L^\infty$ bound. However, this approach can not follow since $V_h$ is not a subspace of $H^2$. Instead, thanks to the discrete Laplace operator $\Delta_h$, we can bound the numerical solution $u_h$ of (6)–(8) in $L^\infty$ norm by $\|\Delta_h u_h\|$.

**Lemma 3** [9] *For any $v \in V_h$, it holds that*

$$
\|v\|_\infty \leq C\|v\|^{\frac{4-d}{4}} (\|v\|^2 + \|\Delta_h v\|^2)^{\frac{d}{8}}. \tag{19}
$$

Let $R_h : H^1 \to V_h$ denote the Ritz projection operator onto the finite element spaces, i.e.,

$$
(\nabla(\phi - R_h \phi), \nabla \varphi) = 0, \quad \forall \phi \in H^1, \phi \in V_h \tag{20}
$$

with $\int_\Omega R_h \phi \, dx = \int_\Omega \phi \, dx$. Then we have from [3]

$$
\|\phi - R_h \phi\| + h\|\nabla(\phi - R_h \phi)\| \leq Ch^s \|\phi\|_{H^s}, \tag{21}
$$

and

$$
\|\phi - R_h \phi\|_{L^\infty} \leq Ch^s \ell_h \|\phi\|_{W^s_\infty}, \tag{22}
$$

where $\ell_h = \max(1, \log(1/h))$ and $1 \leq s \leq k+1$. For the general negative norms, the Ritz projection has the following error estimates.

**Lemma 4** [19] *Let $R_h u$ be the Ritz projection of $u$. It holds that*

$$
\|u - R_h u\|_{-s} \leq Ch^{s+q} \|u\|_q, \quad for \ 0 \leq s \leq k-1, \quad 1 \leq q \leq k+1. \tag{23}
$$

**Lemma 5** [19] $\Delta_h$ *is related to the projection operators by*

$$
\Delta_h R_h = P_h \Delta. \tag{24}
$$

It is well known that the $L^2$ projection $P_h$ is stable in $L^2$ norm. In fact it can be proved to be stable in $H^{-1}$ norm as follows.

**Lemma 6** *For any $u \in H^{-1}(\Omega)$, we have*

$$
\|P_h u\|_{-1} \leq C\|u\|_{-1}. \tag{25}
$$

**Proof** By (18) and the definition of $P_h$ we have

$$\|P_h u\|_{-1} \leq \sup_{v \in H^1} \frac{(P_h u, v)}{\|v\|_1} = \sup_{v \in H^1} \frac{(P_h u, v - P_h v) + (P_h u, P_h v)}{\|v\|_1}$$

$$= \sup_{v \in H^1} \frac{(u, P_h v)}{\|v\|_1} \leq C\|u\|_{-1},$$

which is the desired result. □

**Lemma 7** *For any $u_h \in V_h$, it holds that*

$$\|\Delta_h u_h\| \leq Ch^{-1}\|\nabla u_h\| \leq Ch^{-2}\|u_h\|. \tag{26}$$

**Proof** By the inverse inequality in the finite element space [3,14], we have

$$\|\Delta_h u_h\| \leq \sup_{v \in V_h} \frac{(\Delta_h u_h, v)}{\|v\|} = \sup_{v \in V_h} \frac{-(\nabla u_h, \nabla v)}{\|v\|}$$

$$\leq \sup_{v \in V_h} \frac{Ch^{-1}\|\nabla u_h\|\|v\|}{\|v\|} \leq Ch^{-1}\|\nabla u_h\| \leq Ch^{-2}\|u_h\|,$$

which implies (26). □

## 4 Error estimates for the first-order fully discrete scheme

In the first part of this section, we present error estimates of the finite element solution based on a mathematical induction on the $L^\infty$ norm of $u_h^n$, which is needed to control the nonlinear terms in the equations. Then in the second part, we will derive a stronger stability result for the finite element solution and improve the error estimate such that it is optimal in each norm.

### 4.1 A first error estimate

**Lemma 8** *In addition to the assumptions in Proposition 1, we assume the following regularity holds:*

$$u \in L^\infty([0, T]; W_\infty^1), \quad u_t \in L^\infty([0, T], H^1), \quad u_{tt} \in L^2((0, T); L^2). \tag{27}$$

*Then the truncation errors due to the time discretization satisfy that*

$$\max_{1 \leq n \leq N} \left( \|E_u^n\|_{-1} + \|E_\mu^n\|_1 + |E_r^n| \right) \leq C\Delta t.$$

**Proof** By direct calculation,

$$r_{tt} = -\frac{1}{4\sqrt{E_1[u]^3}} \left( \int_\Omega g(u)u_t dx \right)^2 + \frac{1}{2\sqrt{E_1[u]}} \int_\Omega (g'(u)u_t^2 + g(u)u_{tt})dx.$$

We know from Proposition 1 that $u \in L^2(0, T; H^4) \cap C([0, T]; H^2)$ (thus $u \in L^\infty([0, T], L^\infty(\Omega))$), from which together with (27), (15) and (16) we deduce that

$$\int_0^T |r_{tt}|^2 dt \leq C \int_0^T (\|u_t\|_{L^4}^4 + \|g'(u)\nabla u\|_{L^\infty}^2 \|u_{tt}\|_{-1}^2)dt$$

$$\leq C\|\nabla u\|_{L^\infty(L^\infty)}^2 \int_0^T (\|u_t\|_1^4 + \|u_{tt}\|_{-1}^2)dt \leq C. \qquad (28)$$

If we define $v^n = v(t^n)$ for any continuous function $v$ and $D_t\varphi^{n+1} := (\varphi^{n+1} - \varphi^n)/\Delta t$ for any function value $\{\varphi^n\}_{n=0}^N$, this gives that the exact solution $(u, \mu, r)$ satisfies the equations

$$(u^{n+1} - u^n, v) = -\Delta t(\nabla \mu^{n+1}, \nabla v) + \Delta t(E_u^n, v), \qquad (29)$$

$$(\mu^{n+1}, \tau) = (\nabla u^{n+1}, \nabla \tau) + \lambda(u^{n+1}, \tau) + \frac{r^{n+1}}{\sqrt{E_1(u^n)}}(g(u^n), \tau) + (E_\mu^n, \tau),$$
$$\qquad (30)$$

$$r^{n+1} - r^n = \frac{1}{2\sqrt{E_1[u^n]}} \int_\Omega g(u^n)(u^{n+1} - u^n)dx + \Delta t E_r^n, \qquad (31)$$

where

$$|E_u^n| = |D_t u^{n+1} - \partial_t u^{n+1}| \leq \int_{t^n}^{t^{n+1}} |u_{tt}(\cdot, s)|ds, \qquad (32)$$

$$|E_r^n| = |D_t r^{n+1} - \partial_t r^n - \frac{1}{2} \int_\Omega \frac{g(u^n)}{\sqrt{E_1[u^n]}}(D_t u^{n+1} - \partial_t u^n)dx|$$

$$\leq C \left( \int_{t^n}^{t^{n+1}} |r_{tt}(s)|ds + \int_{t^n}^{t^{n+1}} \int_\Omega |u_{tt}(x, s)|dsdx \right), \qquad (33)$$

and

$$E_\mu^n = r^{n+1} \left( \frac{g(u^{n+1})}{\sqrt{E_1(u^{n+1})}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}} \right).$$

Hence, from (15) we derive

$$\|E_\mu^n\|_s = |r^{n+1}| \|\| \frac{g(u^{n+1})}{\sqrt{E_1(u^{n+1})}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}} \|_s \le \sup_{t \in [0,T]} |r(t)|$$

$$\left( \|g(u^n)\|_s \frac{|E_1(u^n) - E_1(u^{n+1})|}{\sqrt{E_1(u^{n+1}) E_1(u^n)(E_1(u^n) + E_1(u^{n+1}))}} + \frac{\|g(u^{n+1}) - g(u^n)\|_s}{\sqrt{E_1(u^{n+1})}} \right)$$

$$\le C \left( \|r\|_{L^\infty}, \|u_t\|_{L^\infty(H^s)}, \|u\|_{L^\infty(L^\infty)} \right) \Delta t. \tag{34}$$

Together with (27) we get the desired results.                                          $\square$

If we define the discrete norms as follows:

$$\|u\|_{l^2(H^s)} := \left( \sum_{n=1}^N \Delta t \|u^n\|_s^2 \right)^{1/2}, \quad \|u\|_{l^\infty(H^s)} := \sup_{1 \le n \le N} \|u^n\|_s,$$

we can present the following error estimates for the finite element approximations.

**Theorem 1** *Under the same assumptions as in Lemma 8, let $(u, \mu, r)$ and $(u_h, \mu_h, r_h)$ be the solution of (3)– (5) and (6)–(8), respectively. Then we have the following error estimates*

$$\|u - u_h\|_{l^\infty(L^2)} + \|\nabla(u - u_h)\|_{l^\infty(L^2)} + \|\nabla(\mu - \mu_h)\|_{l^2(L^2)} + |r - r_h|_{l^\infty}$$
$$\le Ch^k \left( \|u\|_{l^\infty(H^{k+1})} + \|u\|_{H^1(H^{\max(1,k-1)})} + \|u\|_{L^2(H^{k+3})} \right) + C \Delta t. \tag{35}$$

**Proof** Let $\theta_u^{n+1} = u_h^{n+1} - R_h u^{n+1}, \theta_\mu^{n+1} = \mu_h^{n+1} - R_h \mu^{n+1}, e^{n+1} = r_h^{n+1} - r^{n+1}$ and $\rho_u^{n+1} = u^{n+1} - R_h u^{n+1}, \rho_\mu^{n+1} = \mu^{n+1} - R_h \mu^{n+1}$. By using (20) the difference between (6)–(8) and (29)–(31) gives that

$$(D_t \theta_u^{n+1}, v) + (\nabla \theta_\mu^{n+1}, \nabla v) = (D_t \rho_u^{n+1}, v) - (E_u^n, v), \tag{36}$$

$$(\theta_\mu^{n+1}, \tau) - (\nabla \theta_u^{n+1}, \nabla \tau) - \lambda(\theta_u^{n+1}, \tau) = (\rho_\mu^{n+1}, \tau) - \lambda(\rho_u^{n+1}, \tau)$$
$$+ \frac{e^{n+1}}{\sqrt{E_1(u^n)}} (g(u^n), \tau) + r_h^{n+1} \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}, \tau \right) - (E_\mu^n, \tau), \tag{37}$$

$$D_t e^{n+1} - \frac{1}{2} \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}}, D_t \theta_u^{n+1} \right) = \frac{1}{2} \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}, D_t u^{n+1} \right)$$

$$- \frac{1}{2} \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}}, D_t \rho_u^{n+1} \right) - E_r^n \tag{38}$$

for all $v, \tau \in V_h$. We can choose appropriate initial approximations such that $\|\theta_u^0\| + h\|\nabla \theta_u^0\| + \|\theta_\mu^0\| + h\|\nabla \theta_\mu^0\| + |e^0| \le Ch^{k+1}$.

The Eq. (37) implies that

$$
\|\Delta_h \theta_u^{n+1}\|^2 \le \left( \|\theta_\mu^{n+1}\| + \lambda \|\theta_u^{n+1}\| + \|\rho_\mu^{n+1}\| + \lambda \|\rho_u^{n+1}\| + |e^{n+1}| \left\| \frac{g(u^n)}{\sqrt{E_1(u^n)}} \right\| \right.
$$

$$
\left. + |r_h^{n+1}| \left\| \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}} \right\| + \|E_\mu^n\| \right)^2 := I_1. \tag{39}
$$

Then substitute $v = \theta_\mu^{n+1}, \theta_u^{n+1}$ into (36), $\tau = D_t \theta_u^{n+1}, \theta_\mu^{n+1}$ into (37) respectively. Then multiplying (38) with $2e^{n+1}$ and adding the resultant equalities with (39), we derive that

$$
\frac{1}{2}(\|\theta_u^{n+1}\|_1^2 - \|\theta_u^n\|_1^2) + \Delta t \|\Delta_h \theta_u^{n+1}\|^2 + \Delta t \|\theta_\mu^{n+1}\|_1^2 + (e^{n+1})^2 - (e^n)^2 \le \Delta t I_1 + \sum_{i=2}^{6} I_i, \tag{40}
$$

with

$$
\begin{aligned}
I_2 ={}& \Delta t (D_t \rho_u^{n+1}, \theta_\mu^{n+1}) - \Delta t (E_u^n, \theta_\mu^{n+1}) + \Delta t (D_t \rho_u^{n+1}, \theta_u^{n+1}) \\
& - \Delta t (E_u^n, \theta_u^{n+1}) + \lambda \Delta t (\theta_u^{n+1}, \theta_\mu^{n+1}) \\
& + \Delta t (\rho_\mu^{n+1}, \theta_\mu^{n+1}) - \lambda \Delta t (\rho_u^{n+1}, \theta_\mu^{n+1}) - \Delta t (E_\mu^n, \theta_\mu^{n+1}), \\
I_3 ={}& -\Delta t (\rho_\mu^{n+1}, D_t \theta_u^{n+1}) + \lambda \Delta t (\rho_u^{n+1}, D_t \theta_u^{n+1}) + \Delta t (E_\mu^n, D_t \theta_u^{n+1}), \\
I_4 ={}& -r^{n+1} \Delta t \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}, D_t \theta_u^{n+1} \right) \\
& + e^{n+1} \Delta t \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}, D_t u^{n+1} \right) \\
I_5 ={}& \Delta t r_h^{n+1} \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}, \theta_\mu^{n+1} \right), \\
I_6 ={}& -e^{n+1} \Delta t (\frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}}, D_t \rho_u^{n+1}) - 2\Delta t e^{n+1} E_r^n + \frac{e^{n+1} \Delta t}{\sqrt{E_1(u^n)}} (g(u^n), \theta_\mu^{n+1}).
\end{aligned}
$$

At this moment, we invoke a mathematical induction on

$$
\|u_h^n\|_{L^\infty} \le \|u\|_{L^\infty(L^\infty)} + 1. \tag{41}
$$

Since $n = 0$, $u_h^0$ can be defined by any appropriate approximation of $u^0$ satisfying (41), e.g., $u_h^0 = R_h u^0$, we derive from (22) that

$$
\|u^0 - u_h^0\|_{L^\infty} \le Ch^2 \ell_h \|u^0\|_{W_\infty^2}.
$$

Thus there exists a positive constant $h_1$ such that (41) holds for $n = 0$ when $h < h_1$. In the following, we present estimates of the finite element solution by assuming that (41) holds for $0 \leq n \leq m$, for some nonnegative integer $m$. We shall see that if (41) holds for $0 \leq n \leq m$, then it also holds for $n = m + 1$.

Below we bound each term on the righthand side of (40). By (23) and Lemma 8,

$$
\begin{aligned}
I_2 &\leq C\Delta t\big(\|D_t\rho_u^{n+1}\|_{-1}^2 + \|\rho_u^{n+1}\|_{-1}^2 + \|\rho_\mu^{n+1}\|_{-1}^2\big) + \frac{\Delta t}{8}\|\theta_\mu^{n+1}\|_1^2 + C\Delta t\big(\|E_u^n\|_{-1}^2 \\
&\quad + \|E_\mu^n\|_{-1}^2\big) + C\Delta t\|\theta_u^{n+1}\|_1^2 \\
&\leq Ch^{2k}\int_{t^n}^{t^{n+1}}\|u_t\|_{\max(1,k-1)}^2 dt + Ch^{2k}\Delta t\big(\|\mu^{n+1}\|_{\max(1,k-1)}^2 + \|u^{n+1}\|_{\max(1,k-1)}^2\big) \\
&\quad + \frac{\Delta t}{8}\|\theta_\mu^{n+1}\|_1^2 + C\Delta t^2\int_{t^n}^{t^{n+1}}\|u_{tt}\|_{-1}^2 dt + C\Delta t^3 + C\Delta t\|\theta_u^{n+1}\|_1^2.
\end{aligned}
$$

It follows from (36) that

$$
D_t\theta_u^{n+1} = \Delta_h\theta_\mu^{n+1} + P_h\left(D_t\rho_u^{n+1} - E_u^n\right). \tag{42}
$$

Together with (17), (25), (32) and (23) we can derive

$$
\begin{aligned}
\|D_t\theta_u^{n+1}\|_{-1} &\leq \|\nabla\theta_\mu^{n+1}\| + \|D_t\rho_u^{n+1}\|_{-1} + \|E_u^n\|_{-1} \\
&\leq \|\nabla\theta_\mu^{n+1}\| + Ch^k(\Delta t)^{-1}\int_{t^n}^{t^{n+1}}\|u_t\|_{\max(1,k-1)}dt + \int_{t^n}^{t^{n+1}}\|u_{tt}\|_{-1}dt.
\end{aligned} \tag{43}
$$

Then we derive from (43) and (34) that

$$
I_3 \leq C\Delta t h^{2k}\|\mu^{n+1}\|_{k+1}^2 + C\Delta t h^{2k}\|u^{n+1}\|_{k+1}^2 + C\Delta t^3 + Z^n, \tag{44}
$$

where

$$
\begin{aligned}
Z^n &:= \frac{\Delta t}{8}\|\nabla\theta_\mu^{n+1}\|^2 + Ch^{2k}\int_{t^n}^{t^{n+1}}\|u_t\|_{\max(1,k-1)}^2 dt \\
&\quad + C\Delta t^2\int_{t^n}^{t^{n+1}}\|u_{tt}\|_{-1}^2 dt.
\end{aligned}
$$

In view of (27),

$$
I_4 \leq Z^n + C\Delta t(e^{n+1})^2 + C\Delta t\left\|\frac{\nabla g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{\nabla g(u^n)}{\sqrt{E_1(u^n)}}\right\|^2.
$$

The last term on the right-hand side of the last inequality needs to be treated as follows

$$\frac{\nabla g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{\nabla g(u^n)}{\sqrt{E_1(u^n)}} = \nabla g(u^n)\frac{E_1(u^n) - E_1(u_h^n)}{\sqrt{E_1(u_h^n)E_1(u^n)}(E_1(u^n) + E_1(u_h^n))}$$
$$+\frac{\nabla g(u_h^n) - \nabla g(u^n)}{\sqrt{E_1(u_h^n)}} = A_1 + A_2. \tag{45}$$

The two terms on the righthand side are bounded by

$$\|A_1\| \le C\|\nabla g(u^n)\|\|u^n - u_h^n\| \le C(\|\theta_u^n\| + \|\rho_u^n\|), \tag{46}$$

and

$$\|A_2\| \le C\|\nabla g(u_h^n) - \nabla g(u^n)\|$$
$$\le C(\|(g'(u_h^n) - g'(u^n))\nabla u^n\| + \|g'(u_h^n)\nabla\rho_u^n\| + \|g'(u_h^n)\nabla\theta_u^n\|)$$
$$\le C(\|\nabla\rho_u^n\| + \|\rho_u^n\| + \|\theta_u^n\| + \|\nabla\theta_u^n\|), \tag{47}$$

where we have used the fact that $E_1[u_h^n] > C_0$ and $g'(u_h^n)$, $g''(u_h^n)$ have an uniform upper bound by (15), (16) and (41). Therefore, we combine (45), (46) and (47) to get

$$\|\frac{\nabla g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{\nabla g(u^n)}{\sqrt{E_1(u^n)}}\| \le C(\|\nabla(\rho_u^n)\| + \|\rho_u^n\| + \|\theta_u^n\| + \|\nabla\theta_u^n\|)$$
$$\le Ch^k\|u^n\|_{k+1} + C(\|\theta_u^n\| + \|\nabla\theta_u^n\|). \tag{48}$$

Similarly, it holds that

$$\frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}} \le C\|u^n - u_h^n\| \le C(\|\theta_u^n\| + \|\rho_u^n\|), \tag{49}$$

which gives that

$$\Delta t I_1 \le \frac{1}{8}\Delta t\|\theta_\mu^{n+1}\|^2 + C\Delta t(\|\theta_u^n\|^2 + \|\theta_u^{n+1}\|^2) + C\Delta t h^{2k}(\|\mu^{n+1}\|_k^2$$
$$+\|u^n\|_k^2 + \|u^{n+1}\|_k^2) + C\Delta t|e^{n+1}|^2 + C\Delta t^3$$

and by (9)

$$I_5 \le C\Delta t\|\theta_u^n\|^2 + C\Delta t h^{2k}\|u^n\|_k^2 + \frac{1}{8}\Delta t\|\theta_\mu^{n+1}\|^2.$$

Again by using (9), (41) and (33), we can derive that $\|g'(u_h^n)\nabla u_h^n\| \le C$ so that

$$I_6 \le C\Delta t(e^{n+1})^2 + C\Delta t^3 + Ch^{2k}\int_{t^n}^{t^{n+1}}\|u_t\|_{\max(1,k-1)}^2 dt + \frac{\Delta t}{8}\|\theta_\mu^{n+1}\|^2.$$

We combine the estimates for each term $I_i$, $1 \le i \le 6$, and (48) to get

$$\frac{1}{2}(\|\theta_u^{n+1}\|_1^2 - \|\theta_u^n\|_1^2) + \Delta t \|\Delta_h \theta_u^{n+1}\|^2 + \frac{1}{8} \Delta t \|\theta_\mu^{n+1}\|_1^2 + (e^{n+1})^2 - (e^n)^2$$

$$\le C \Delta t \left( \|\theta_u^{n+1}\|_1^2 + \|\theta_u^n\|_1^2 + (e^{n+1})^2 \right) + C \Delta t^2 \int_{t^n}^{t^{n+1}} \|u_{tt}(s)\|_{-1}^2 ds + C \Delta t^3 + C h^{2k}$$

$$\left( \int_{t^n}^{t^{n+1}} \|u_t\|_{\max(k-1,1)}^2 dt + \Delta t \|\mu^{n+1}\|_{k+1}^2 + \Delta t \|u^{n+1}\|_{k+1}^2 + \Delta t \|u^n\|_{k+1}^2 \right). \quad (50)$$

By applying Gronwall's inequality, we have

$$\max_{0 \le n \le m} \left( \|\theta_u^{n+1}\|_1^2 + (e^{n+1})^2 \right) + \Delta t \sum_{n=0}^{m} \left( \|\theta_\mu^{n+1}\|_1^2 + \|\Delta_h \theta_u^{n+1}\|^2 \right)$$

$$\le C h^{2k} \left( \|u\|_{H^1(H^{\max(1,k-1)})} + \|\mu\|_{l^2(H^{k+1})} + \|u\|_{l^2(H^{k+1})} \right)^2 + C \Delta t^2$$

$$\le C_1 (h^{2k} + \Delta t^2). \quad (51)$$

for some positive constant $C_1$.

Hence, there exists a positive constant $h_2$ such that when $h < h_2$ and if $\Delta t \le h$, then from the inverse inequality (c.f. [19, Lemma 6.4]) we have

$$\|\theta_u^{m+1}\|_{L^\infty}^2 \le C \log(1/h) \|\theta_u^{m+1}\|_1^2 \le C h^{-1}(h^{2k} + \Delta t^2) \le C(h^{2k-1} + \Delta t).$$

On the other hand, if $\Delta t \ge h$, then we have

$$\|\Delta_h \theta_u^{m+1}\|^2 \le \frac{1}{\Delta t} \sum_{n=0}^{m} \Delta t \|\Delta_h \theta_u^{n+1}\|^2 \le C \left( \frac{h^{2k}}{\Delta t} + \Delta t \right) \le C(h^{2k-1} + \Delta t).$$

Overall, together with (19) we have $\|\theta_u^{m+1}\|_{L^\infty}^2 \le C(h^{2k-1} + \Delta t)$, and so from (22)

$$\|u_h^{m+1} - u^{m+1}\|_{L^\infty} \le \|\theta_u^{m+1}\|_{L^\infty} + \|u^{m+1} - R_h u^{m+1}\|_{L^\infty} \le C(h^{k-1/2} + \Delta t^{1/2}),$$

where $C$ is independent of $m$. Thus there exists positive constants $h_3$, $\Delta t_3$ such that when $h < h_3$ and $\Delta t < \Delta t_3$ we have

$$\|u_h^{m+1} - u^{m+1}\|_{L^\infty} \le 1,$$

and this completes the mathematical induction on (41) in the case that $h < h_3$ and $\Delta t < \Delta t_3$. Thus (51) holds for $m = N - 1$ with the same constant $C_1$, provided $h < h_3$ and $\Delta t < \Delta t_3$.

If $h \ge h_3$ or $\Delta t \ge \Delta t_3$, from (9) we see that

$$\max_{0 \le n \le N-1} \left( \|\theta_u^{n+1}\|_1^2 + (e^{n+1})^2 \right) + \Delta t \sum_{n=0}^{N-1} \|\nabla \theta_\mu^{n+1}\|^2$$

$$\leq C_2 \leq C_2(\Delta t_3^{-2} + h_3^{-2k})(\Delta t^2 + h^{2k}) \tag{52}$$

for some positive constant $C_2$. From (51) and (52) we obtain for any $\Delta t$ and $h$,

$$\max_{0 \leq n \leq N-1} \left( \|\theta_u^{n+1}\|_1^2 + (e^{n+1})^2 \right) + \Delta t \sum_{n=0}^{N-1} \|\nabla \theta_\mu^{n+1}\|^2$$
$$\leq (C_1 + C_2(\Delta t_3^{-2} + h_3^{-2k}))(\Delta t^2 + h^{2k})$$

from which we obtain the desired result (35) by using (21) and the triangle inequality.
$\square$

### 4.2 An improved $L^2$-error estimate

Notice that the error estimate for $L^2$ norm in Theorem 1 is not optimal. Below, we shall derive an improved $L^2$-error estimate. For this purpose, we improve Lemma 2 first.

**Lemma 9** *For the operator $\Delta_h$, we have the following bound*

$$\|\Delta_h v\|_{-2} \leq C\|v\|, \quad \forall v \in V_h. \tag{53}$$

**Proof** By the definition (10) of the operator $\Delta_h$, (24), (17) and inverse inequalities of finite element space [19], we get

$$\|\Delta_h v\|_{-2} = \sup_{0 \neq \varphi \in H^2} \frac{(\Delta_h v, \varphi)}{\|\varphi\|_2} = \sup_{0 \neq \varphi \in H^2} \frac{(\Delta_h v, \varphi - R_h \varphi)}{\|\varphi\|_2} + \sup_{0 \neq \varphi \in H^2} \frac{(\Delta_h v, R_h \varphi)}{\|\varphi\|_2}$$
$$\leq Ch\|\Delta_h v\|_{-1} + \sup_{0 \neq \varphi \in H^2} -\frac{(\nabla v, \nabla R_h \varphi)}{\|\varphi\|_2} \leq Ch\|\nabla v\| + \sup_{0 \neq \varphi \in H^2} \frac{(v, \Delta_h R_h \varphi)}{\|\varphi\|_2}$$
$$\leq C\|v\| + \sup_{0 \neq \varphi \in H^2} \frac{(v, P_h \Delta \varphi)}{\|\varphi\|_2} \leq C\|v\|,$$

which is the desired result.
$\square$

Then we derive a stronger stability result, that is the difference quotient for the numerical solution $u_h^n$ is also uniformly bounded, which is needed to get the improved error estimates (see the term $Q_{12}$ in Theorem 2). The derivation process is divided into two steps, Lemmas 10 and 11 respectively. In fact, from (6) and (17), it follows that $\|D_t u_h^n\|_{-1} \leq \|\nabla \mu_h^n\|$. Together with (9), we have

$$\sum_{n=0}^{N-1} \Delta t \|D_t u_h^{n+1}\|_{-1}^2 \leq C_0. \tag{54}$$

In Lemma 10 we shall derive the boundedness for $\sum_{n=0}^{N-1} \Delta t \|D_t u_h^{n+1}\|^2$ via (54) when $\lambda > 0$. Then it follows in Lemma 11 the desired stability, that is $\|D_t u_h^n\|$ is uniformly bounded.

**Lemma 10** *For the $H^{-1}$ gradient flow and all $N \leq T/\Delta t$, there exists a positive constant $C$ such that*

$$\lambda \sum_{n=0}^{N-1} \Delta t \|D_t u_h^{n+1}\|^2 \leq C, \tag{55}$$

*where we define $\|\Delta_h^{-1/2} u_h\| := \|\nabla \Delta_h^{-1} u_h\|$.*

*Proof* By replacing $u_h$ by $\Delta_h^{-1} u_h$, we derive from (10) that

$$(\nabla \Delta_h^{-1} u_h, \nabla v) = -(u_h, v), \quad \forall v \in V_h. \tag{56}$$

Taking the difference of (6)–(7) at the time $t^{n+1}$ and $t^n$, we have

$$(D_t u_h^{n+1} - D_t u_h^n, v) = -(\nabla(\mu_h^{n+1} - \mu_h^n), \nabla v), \quad \forall v \in V_h, \tag{57}$$

$$(\mu_h^{n+1} - \mu_h^n, \tau) = (\nabla(u_h^{n+1} - u_h^n), \nabla \tau) + \lambda(u_h^{n+1} - u_h^n, \tau) + \frac{r_h^{n+1} - r_h^n}{\sqrt{E_1(u_h^n)}}(g(u_h^n), \tau)$$

$$+ r_h^n \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u_h^{n-1})}{\sqrt{E_1(u_h^{n-1})}}, \tau \right), \quad \forall \tau \in V_h. \tag{58}$$

We choose $(v, \tau) = (-\Delta_h^{-1} D_t u_h^{n+1}, D_t u_h^{n+1})$ in (57)–(58) and using Young's inequality arrive at

$$\frac{1}{2}(\|\Delta_h^{-1/2} D_t u_h^{n+1}\|^2 - \|\Delta_h^{-1/2} D_t u_h^n\|^2) + \Delta t \|\nabla D_t u_h^{n+1}\|^2 + \lambda \Delta t \|D_t u_h^{n+1}\|^2$$

$$\leq \frac{\Delta t}{2 E_1(u_h^n)}(g(u_h^n), D_t u_h^{n+1})^2 + |r_h^n| \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u_h^{n-1})}{\sqrt{E_1(u_h^{n-1})}}, D_t u_h^{n+1} \right)$$

$$\leq C\Delta t (\|D_t u_h^n\|^2 + \|D_t u_h^{n+1}\|^2)$$

$$\leq \frac{\Delta t}{2}(\|\nabla D_t u_h^n\|^2 + \|\nabla D_t u_h^{n+1}\|^2) + C\Delta t (\|D_t u_h^n\|_{-1}^2 + \|D_t u_h^{n+1}\|_{-1}^2),$$

where we also used (8), (9), (41), (49) and (56). Noticing (54) and $\|D_t u_h^0\|_1 \leq C$, we thus from summing the last inequality over $n$ from 0 to $N-1$ obtain the desired result (55). □

**Lemma 11** *For the $H^{-1}$ gradient flow and all $N \leq T/\Delta t$, if $\lambda > 0$ there exists a positive constant $C$ such that*

$$\|D_t u_h^n\| \leq C, \quad \forall 0 \leq n \leq N. \tag{59}$$

**Proof** Taking $(v, \tau) = (D_t u_h^{n+1}, -\Delta_h D_t u_h^{n+1})$ in (57)–(58) leads to

$$
\frac{1}{2}(\|D_t u_h^{n+1}\|^2 - \|D_t u_h^n\|^2) + \Delta t \|\Delta_h D_t u_h^{n+1}\|^2 + \lambda \Delta t \|\nabla D_t u_h^{n+1}\|^2
$$
$$
\leq \frac{\Delta t}{2 E_1(u_h^n)} (g(u_h^n), D_t u_h^{n+1})(g(u_h^n), \Delta_h D_t u_h^{n+1})
$$
$$
+ |r_h^n| \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u_h^{n-1})}{\sqrt{E_1(u_h^{n-1})}}, \Delta_h D_t u_h^{n+1} \right)
$$
$$
\leq C \Delta t (\|D_t u_h^{n+1}\| + \|D_t u_h^n\|) \|\Delta_h D_t u_h^{n+1}\|
$$
$$
\leq C \Delta t \left( \|D_t u_h^n\|^2 + \|D_t u_h^{n+1}\|^2 \right) + \frac{1}{2} \Delta t \|\Delta_h D_t u_h^{n+1}\|^2,
$$

where we also used (8), (9), (41) and (49). After summing the last inequality over $n$ from 0 to $N - 1$, using (55) we complete the proof. $\qquad\square$

From the proof of Theorem 1, we can see the "trouble" term is $\left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}}, D_t \theta_u^{n+1} \right)$ in (38). Due to the fact $g(u_h^n) \notin H^2$, it implies that there appears $\|D_t \theta_u^{n+1}\|_{-1}$ on the right-hand side of (40), and so by (36) $\|\nabla \theta_\mu^{n+1}\|$ arises with no doubt. Therefore, we can only obtain the error estimate with order $h^k$. In order to get the optimal error estimate in $L^2$ norm, we need to rewrite (38) as follows

$$
D_t e^{n+1} - \frac{1}{2} \left( \frac{g(u^n)}{\sqrt{E_1(u^n)}}, D_t \theta_u^{n+1} \right) = \frac{1}{2} \left( \frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}, D_t u_h^{n+1} \right)
$$
$$
- \frac{1}{2} \left( \frac{g(u^n)}{\sqrt{E_1(u^n)}}, D_t \rho_u^{n+1} \right) - E_r^n. \quad (60)
$$

Then with the help of (59), instead we have $\|D_t \theta_u^{n+1}\|_{-2}$ and $\|\frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}\|$ on the right-hand side of the following equality (61) such that the following theorem holds.

**Theorem 2** *Under the same assumptions as in Lemma 8, we have*

$$
\|u - u_h\|_{l^\infty(L^2)} + \|\mu - \mu_h\|_{l^2(L^2)} + |r - r_h|_{l^\infty}
$$
$$
\leq C h^{k+1} \left( \|u_t\|_{L^2(H^{\max(2,k)})} + \|u\|_{L^2(H^{k+3})} + \|u\|_{L^\infty(H^{k+1})} \right) + C \Delta t.
$$

**Proof** We substitute $v = \theta_u^{n+1}$ into (36) respectively, take $\tau = \theta_\mu^{n+1}$ in (37), and multiply (60) by $2 \Delta t e^{n+1}$. Then we add the resultant equalities to get

$$
\frac{1}{2}(\|\theta_u^{n+1}\|^2 - \|\theta_u^n\|^2) + \Delta t \|\theta_\mu^{n+1}\|^2 + (e^{n+1})^2 - (e^n)^2
$$

$$\leq \Delta t(D_t\rho_u^{n+1}, \theta_u^{n+1}) - \Delta t(E_u^n, \theta_u^{n+1}) + \lambda\Delta t(\theta_u^{n+1}, \theta_\mu^{n+1})$$

$$+\Delta t(\rho_\mu^{n+1}, \theta_\mu^{n+1}) - \lambda\Delta t(\rho_u^{n+1}, \theta_\mu^{n+1})$$

$$+\frac{e^{n+1}\Delta t}{\sqrt{E_1(u^n)}}(g(u^n), \theta_\mu^{n+1}) + \Delta t r_h^{n+1}\left(\frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}, \theta_\mu^{n+1}\right) - \Delta t(E_\mu^n, \theta_\mu^{n+1})$$

$$+e^{n+1}\Delta t\left(\frac{g(u^n)}{\sqrt{E_1(u^n)}}, D_t\theta_u^{n+1}\right) + e^{n+1}\Delta t\left(\frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}, D_t u_h^{n+1}\right)$$

$$-e^{n+1}\Delta t\left(\frac{g(u^n)}{\sqrt{E_1(u^n)}}, D_t\rho_u^{n+1}\right) - 2e^{n+1}\Delta t E_r^n$$

$$=\sum_{i=1}^{14} Q_i. \tag{61}$$

Since $Q_i$, $1 \leq i \leq 10$ have been already appeared on the right-hand side of (40), we only need to estimate the remaining terms $Q_i$, $11 \leq i \leq 14$ as follows.

Since by (42), (53) and the fact that $\|\cdot\|_{-2} \leq \|\cdot\|_{-1}$,

$$\|D_t\theta_u^{n+1}\|_{-2} \leq \|\Delta_h\theta_\mu^{n+1}\|_{-2} + \|D_t\rho_u^{n+1}\|_{-1} + \|E_u^n\|_{-1}$$

$$\leq \|\theta_\mu^{n+1}\| + Ch^{k+1}(\Delta t)^{-1}\int_{t^n}^{t^{n+1}} \|u_t\|_{\max(2,k)}dt$$

$$+\int_{t^n}^{t^{n+1}} \|u_{tt}\|_{-1}dt,$$

it follows that

$$Q_{11} \leq C\Delta t(e^{n+1})^2 + \frac{1}{16}\Delta t\|D_t\theta_u^{n+1}\|_{-2}^2$$

$$\leq C\Delta t(e^{n+1})^2 + \frac{1}{8}\Delta t\|\theta_\mu^{n+1}\|^2 + Ch^{2k+2}\int_{t^n}^{t^{n+1}} \|u_t\|_{\max(2,k)}^2 dt$$

$$+C\Delta t^2\int_{t^n}^{t^{n+1}} \|u_{tt}\|_{-1}^2 dt.$$

As a consequence of (49) and (59),

$$Q_{12} \leq Ce^{n+1}\Delta t\left\|\frac{g(u_h^n)}{\sqrt{E_1(u_h^n)}} - \frac{g(u^n)}{\sqrt{E_1(u^n)}}\right\|$$

$$\leq C\Delta t(e^{n+1})^2 + C\Delta t\|\theta_u^n\|^2 + C\Delta t h^{2k+2}\|u^n\|_{k+1}^2.$$

Then from Lemma 8, we have

$$Q_{13} \leq C\Delta t(e^{n+1})^2 + Ch^{2k+2} \int_{t^n}^{t^{n+1}} \|u_t\|^2_{\max(2,k)} dt,$$

$$Q_{14} \leq \Delta t(e^{n+1})^2 + C\Delta t^2 \int_{t^n}^{t^{n+1}} |r_{tt}(s)|^2 + \|u_{tt}(s)\|^2 ds.$$

Combing the estimates for $Q_i$, $1 \leq i \leq 14$, we have

$$\frac{1}{2}(\|\theta_u^{n+1}\|^2 - \|\theta_u^n\|^2) + \frac{1}{16}\Delta t\|\theta_\mu^{n+1}\|^2 + (e^{n+1})^2 - (e^n)^2$$

$$\leq Ch^{2k+2} \int_{t^n}^{t^{n+1}} \|u_t\|^2_{\max(2,k)} dt + C\Delta t h^{2k+2} (\|\mu^{n+1}\|^2_{k+1} + \|u^{n+1}\|^2_{k+1} + \|u^n\|^2_{k+1})$$

$$+ C\Delta t(\|\theta_u^n\|^2 + \|\theta_u^{n+1}\|^2 + (e^{n+1})^2) + \Delta t^2 \int_{t^n}^{t^{n+1}} |r_{tt}(s)|^2 + \|u_{tt}\|^2 ds + C\Delta t^3.$$

Add the last inequality for $n = 0, \ldots, N - 1$, then by discrete Gronwall inequality we have

$$\|\theta_u\|_{l^\infty(L^2)} + \|\theta_\mu\|_{l^2(L^2)} + |e|_{l^\infty} \leq Ch^{k+1} \left(\|u_t\|_{L^2(\max(2,k))} + \|u\|_{L^2(H^{k+3})}\right) + C\Delta t,$$

from which we obtain the desired result by (21) and the triangle inequality. □

## 5 Error estimates for a second-order fully discrete scheme

A main advantage of the SAV approach (as well as the IEQ approach [22]) is that linear second- or even higher-order energy stable schemes can be easily constructed. In this section, we construct a finite element algorithm based on Crank–Nicolson scheme, then analyze the corresponding error similarly as we did in the last section. First, we state below the finite element approximations for the Crank–Nicolson SAV scheme to the $H^{-1}$ gradient flow, given by: Find $(u_h^{n+1}, \mu_h^{n+1/2}, r_h^{n+1}) \in [V_h]^2 \times R$ such that

$$(u_h^{n+1} - u_h^n, v) = -\Delta t(\nabla \mu_h^{n+1/2}, \nabla v), \quad \forall v \in V_h, \tag{62}$$

$$(\mu_h^{n+1/2}, \tau) = (\nabla u_h^{n+1/2}, \nabla \tau) + \lambda(u_h^{n+1/2}, \tau) + \frac{r_h^{n+1/2}}{\sqrt{E_1(\bar{u}_h^n)}}(g(\bar{u}_h^n), \tau), \quad \forall \tau \in V_h, \tag{63}$$

$$r_h^{n+1} - r_h^n = \frac{1}{2\sqrt{E_1(\bar{u}_h^n)}}(g(\bar{u}_h^n), u_h^{n+1} - u_h^n), \tag{64}$$

with

$$u_h^{n+1/2} = \frac{1}{2}(u_h^{n+1} + u_h^n), \quad r_h^{n+1/2} = \frac{1}{2}(r_h^{n+1} + r_h^n), \quad \bar{u}_h^n = \frac{1}{2}(3u_h^n - u_h^{n-1}). \tag{65}$$

Note that for $n$ fixed these equations employs three time levels rather than the two of our previous methods. We therefore have to restrict its use to $n \geq 1$. With $u_h^0$ given, we

then also need to define $u_h^1$ in some way, e.g. by employing one step of the backward Euler method. Just as in the first-order scheme, one can eliminate $\mu_h^{n+1/2}$ and $r_h^{n+1}$ from (62)–(65) to obtain a linear equation for $u_h^{n+1}$ similar to (11), so it can be solved by using the Sherman–Morrison–Woodbury formula (12) which only involves two linear equations with constant coefficients of the form (14). At the following we shall establish error estimates for the scheme (62)–(65) using a similar procedure.

Similarly as the backward Euler scheme, we can derive the unconditional stability for the scheme (62)–(65).

**Lemma 12** *Let* $(u_h^{n+1}, \mu_h^{n+1/2}, r_h^{n+1})$ *be the solution of the the scheme* (62)–(65)*, then we have*

$$
\max_{0 \leq n \leq N} \frac{\lambda}{2} \|u_h^n\|^2 + \frac{1}{2} \|\nabla u_h^n\|^2 + (r_h^n)^2 + \sum_{n=0}^{N-1} \Delta t \|\nabla \mu_h^{n+1/2}\|^2
$$

$$
\leq \frac{\lambda}{2} \|u_h^0\|^2 + \frac{1}{2} \|\nabla u_h^0\|^2 + (r_h^0)^2. \tag{66}
$$

**Proof** Taking the inner product of (62) and (63) with $\mu_h^{n+1/2}$ and $u_h^{n+1} - u_h^n$ respectively, and multiplying (64) by $r_h^{n+1} + r_h^n$, we derive the following discrete energy law:

$$
\frac{\lambda}{2} \|u_h^{n+1}\|^2 + \frac{1}{2} \|\nabla u_h^{n+1}\|^2 + (r_h^{n+1})^2 - \frac{\lambda}{2} \|u_h^n\|^2 - \frac{1}{2} \|\nabla u_h^n\|^2 - (r_h^n)^2
$$
$$
+ \Delta t \|\nabla \mu_h^{n+1/2}\|^2 = 0,
$$

from which (66) follows immediately. □

We shall first analyze the truncation errors due to the time discretization. Then we will present error estimates for the finite element approximations of the Crank–Nicolson SAV scheme with the help of the stability result (66) and the mathematical induction on $L^\infty$ boundedness derived in Theorem 3.

**Lemma 13** *In addition to the assumptions in Lemma* 8*, we assume the following regularity holds:*

$$
u_{tt} \in L^\infty((0, T); H^1), \quad u_{ttt} \in L^2((0, T); L^2). \tag{67}
$$

*Then the truncation errors due to the time discretization satisfy that*

$$
\max_{1 \leq n \leq N} \|\hat{E}_u^n\| + \|\hat{E}_\mu^n\|_1 + |\hat{E}_r^n| \leq C \Delta t^2. \tag{68}
$$

**Proof** First, like (28), $\|r_{ttt}\|$ can be bounded by $\|\nabla u\|_{L^\infty(L^\infty)}$, $\|u_{tt}\|_{L^2(H^1)}$ and $\|u_{ttt}\|_{L^2(H^{-1})}$ through direct calculation. Then similarly as Lemma 8, we define $v^{n+1/2} = v(t^{n+1/2}) = v((t^n + t^{n+1})/2)$ for any continuous function $v$. Note that the exact solution $(u, \mu, r)$ and $\bar{u}^n = \frac{1}{2}(3u^n - u^{n-1})$ satisfy the equations

$$
(u^{n+1} - u^n, v) = -\Delta t (\nabla \mu^{n+1/2}, \nabla v) + \Delta t (\hat{E}_u^n, v), \tag{69}
$$

$$(\mu^{n+1/2}, \tau) = (\nabla u^{n+1/2}, \nabla \tau) + \lambda(u^{n+1/2}, \tau) + \frac{r^{n+1/2}}{\sqrt{E_1(\bar{u}^n)}}(g(\bar{u}^n), \tau) + (\hat{E}_\mu^n, \tau),$$
(70)

$$r^{n+1} - r^n = \frac{1}{2\sqrt{E_1[\bar{u}^n]}} \int_\Omega g(\bar{u}^n)(u^{n+1} - u^n)dx + \Delta t \hat{E}_r^n,$$
(71)

where

$$|\hat{E}_u^n| = |D_t u^{n+1} - \partial_t u^{n+1/2}|$$

$$= \left| \frac{1}{2\Delta t} \left( \int_{t^n}^{t^{n+1/2}} (s - t^n)^2 u_{ttt}(s)ds + \int_{t^{n+1/2}}^{t^{n+1}} (s - t^{n+1})^2 u_{ttt}(s)ds \right) \right|$$

$$\leq C\Delta t \int_{t^n}^{t^{n+1}} |u_{ttt}(\cdot, s)|ds,$$
(72)

$$|\hat{E}_r^n| = \left| D_t r^{n+1} - \partial_t r^{n+1/2} - \frac{1}{2} \int_\Omega \frac{g(u^{n+1/2})}{\sqrt{E_1[u^{n+1/2}]}}(D_t u^{n+1} - \partial_t u^{n+1/2})dx \right.$$

$$\left. + \frac{1}{2} \int_\Omega \left( \frac{g(u^{n+1/2})}{\sqrt{E_1(u^{n+1/2})}} - \frac{g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}} \right) D_t u^{n+1}dx \right|$$

$$\leq C(\|u\|_{L^\infty(L^\infty)}, \|u_t\|_{L^\infty(L^2)}, \|u_{tt}\|_{L^\infty(L^2)})\left( \Delta t \int_{t^n}^{t^{n+1}} |r_{ttt}(s)|ds \right.$$

$$\left. + \Delta t \int_{t^n}^{t^{n+1}} \int_\Omega |u_{ttt}(x, s)|dsdx + \Delta t^2 \right),$$
(73)

and by (15) we have

$$\|\hat{E}_\mu^n\|_s = |r^{n+1/2}| \| \frac{g(u^{n+1/2})}{\sqrt{E_1(u^{n+1/2})}} - \frac{g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}} \|_s,$$

$$\leq \sup_{t \in [0,T]} |r(t)| \left( \|g(u^{n+1/2})\|_s \frac{|E_1(\bar{u}^n) - E_1(u^{n+1/2})|}{\sqrt{E_1(u^{n+1/2})E_1(\bar{u}^n)}(E_1(\bar{u}^n) + E_1(u^{n+1/2}))} \right.$$

$$\left. + \frac{\|g(u^{n+1/2}) - g(\bar{u}^n)\|_s}{\sqrt{E_1(\bar{u}^n)}} \right)$$

$$\leq C(\|r\|_{L^\infty}, \|u_{tt}\|_{L^\infty(H^s)}, \|u\|_{L^\infty(L^\infty)})\Delta t^2.$$

□

**Theorem 3** *In addition to the same assumptions as Lemma 13, we assume that $u_{tt} \in L^2((0, T); H^3)$. Let $u, \mu, r$ and $u_h, \mu_h, r_h$ be the solution of (3)–(5) and (62)–(64) respectively. Taking $(u_h^0, r_h^0) = (R_h u^0, \sqrt{E_1(u^0)})$ and $u_h^1$ by employing one step of*

*the backward Euler method, then we have the following error estimates*

$$\|u - u_h\|_{l^\infty(L^2)} + \|\nabla(u - u_h)\|_{l^\infty(L^2)} + \|\nabla(\mu - \mu_h)\|_{\hat{l}^2(L^2)} + |r - r_h|_{l^\infty}$$
$$\leq Ch^k \left( \|u\|_{L^\infty(H^{k+1})} + \|u\|_{H^1(H^{\max(k-1,1)})} + \|u\|_{L^2(H^{k+3})} \right) + C\Delta t^2$$
$$+ C\Delta t \left( \int_0^{\Delta t} \|u_{tt}(s)\|_{-1}^2 ds \right)^{1/2},$$

*where the last term on the right-hand side is the error produced by the first step and we denote*

$$\|\nabla(\mu - \mu_h)\|_{\hat{l}^2(L^2)} := \left( \sum_{n=1}^{N-1} \Delta t \|\nabla(\mu^{n+1/2} - \mu_h^{n+1/2})\|^2 \right)^{1/2}.$$

**Proof** With $\rho^n$ bounded as above, we only need to consider $\theta^n$. Substracting (69)–(71) from (62)–(64) gives that

$$(D_t \theta_u^{n+1}, v) + (\nabla \theta_\mu^{n+1/2}, \nabla v) = (D_t \rho_u^{n+1}, v) - (\hat{E}_u^n, v), \tag{74}$$

$$(\theta_\mu^{n+1/2}, \tau) - (\nabla \theta_u^{n+1/2}, \nabla \tau) - \lambda(\theta_u^{n+1/2}, \tau) = (\rho_\mu^{n+1/2}, \tau) - \lambda(\rho_u^{n+1/2}, \tau)$$
$$+ \frac{e^{n+1/2}}{\sqrt{E_1(\bar{u}^n)}} (g(\bar{u}^n), \tau) + r_h^{n+1/2} \left( \frac{g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}} - \frac{g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}}, \tau \right) - (\hat{E}_\mu^n, \tau) \tag{75}$$

$$D_t e^{n+1} - \frac{1}{2} \left( \frac{g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}}, D_t \theta_u^{n+1} \right) = \frac{1}{2} \left( \frac{g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}} - \frac{g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}}, D_t u^{n+1} \right)$$
$$- \frac{1}{2} \left( \frac{g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}}, D_t \rho_u^{n+1} \right) - \hat{E}_r^n \tag{76}$$

for all $v, \tau \in V_h$ where by (65)

$$\theta_\mu^{n+1/2} = \mu_h^{n+1/2} - R_h \mu^{n+1/2}, \ \rho_\mu^{n+1/2} = \mu^{n+1/2} - R_h \mu^{n+1/2}, \ \rho_u^{n+1/2}$$
$$= u^{n+1/2} - R_h u^{n+1/2},$$
$$\theta_u^{n+1/2} = u_h^{n+1/2} - R_h u^{n+1/2} = \frac{1}{2}(\theta_u^{n+1} + \theta_u^n) + R_h \left( \frac{u^{n+1} + u^n}{2} - u^{n+1/2} \right):$$
$$= \frac{1}{2}(\theta_u^{n+1} + \theta_u^n) + w^n,$$
$$e^{n+1/2} = r_h^{n+1/2} - r^{n+1/2} = \frac{1}{2}(e^{n+1} + e^n) + \frac{1}{2}(r^{n+1} + r^n) - r^{n+1/2}:$$
$$= \frac{1}{2}(e^{n+1} + e^n) + y^n.$$

From $(u_h^0, r_h^0) = (R_h u^0, \sqrt{E_1(u^0)})$, it implies

$$\theta_u^0 = e^0 = 0. \tag{77}$$

As for the backward Euler method, from (50) there exist positive constants $h_1$, $\Delta t_1$ such that when $h < h_1$ and $\Delta t < \Delta t_1$ we have $\ell_h \le 1/h$ and

$$\|\theta_u^1\|_1^2 + \Delta t\|\Delta_h \theta_u^1\|^2 + |e^1|^2 \le C_1 \Delta t^3 + C_1 \Delta t h^{2k} + C_1 \Delta t^2 \left( \int_0^{\Delta t} \|u_{tt}(s)\|_{-1}^2 ds \right). \tag{78}$$

Setting $\tau = \Delta_h \theta_u^{n+1}$, we derive from (75) that

$$\frac{\Delta t}{4} \left( \|\Delta_h \theta_u^{n+1}\|^2 - \|\Delta_h \theta_u^n\|^2 + \|\Delta_h(\theta_u^{n+1} + \theta_u^n)\|^2 + \lambda\|\nabla\theta_u^{n+1}\|^2 - \lambda\|\nabla\theta_u^n\|^2 \right.$$
$$\left. + \lambda\|\nabla(\theta_u^{n+1} + \theta_u^n)\|^2 \right)$$
$$= \Delta t \left( -\lambda(\nabla w^n, \nabla\theta_u^{n+1}) + (\nabla\theta_\mu^{n+1/2}, \nabla\theta_u^{n+1}) - (\nabla\rho_\mu^{n+1/2}, \nabla\theta_u^{n+1}) + \lambda(\nabla\rho_\mu^{n+1/2}, \nabla\theta_u^{n+1}) \right)$$
$$- \Delta t \left( \frac{e^{n+1/2}}{\sqrt{E_1(\bar{u}^n)}} (\nabla P_h g(\bar{u}^n), \nabla\theta_u^{n+1}) + r_h^{n+1/2} \left( \frac{\nabla g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}} - \frac{\nabla g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}}, \nabla\theta_u^{n+1} \right) \right)$$
$$+ \Delta t \left( (\nabla\hat{E}_\mu^n, \nabla\theta_u^{n+1}) + (\nabla\Delta_h w^n, \nabla\theta_u^{n+1}) \right) := M_1. \tag{79}$$

Then we substitute $v = \theta_\mu^{n+1/2}$, $\theta_u^{n+1/2}$ into (74) and $\tau = D_t\theta_u^{n+1}$, $\theta_u^{n+1/2}$ into (75) respectively. Multiplying (76) with $e^{n+1} + e^n$ and adding the resultant equalities with (79), we get

$$\frac{1}{2}(\|\theta_u^{n+1}\|_1^2 - \|\theta_u^n\|_1^2) + \Delta t\|\theta_\mu^{n+1/2}\|_1^2 + (e^{n+1})^2 - (e^n)^2$$
$$+ \frac{\Delta t}{4} \left( \|\Delta_h \theta_u^{n+1}\|^2 - \|\Delta_h \theta_u^n\|^2 + \|\Delta_h(\theta_u^{n+1} + \theta_u^n)\|^2 \right) \le \sum_{i=1}^5 M_i, \tag{80}$$

with

$$M_2 = \Delta t \left( (D_t\rho_u^{n+1}, \theta_\mu^{n+1/2}) - (\hat{E}_u^n, \theta_\mu^{n+1/2}) + (D_t\rho_u^{n+1}, \theta_u^{n+1/2}) - (\hat{E}_u^n, \theta_u^{n+1/2}) \right)$$
$$+ \Delta t \left( \lambda(\theta_u^{n+1/2}, \theta_\mu^{n+1/2}) + (\rho_\mu^{n+1/2}, \theta_\mu^{n+1/2}) - \lambda(\rho_u^{n+1/2}, \theta_u^{n+1/2}) - (\hat{E}_\mu^n, \theta_\mu^{n+1/2}) \right),$$

$$M_3 = \Delta t(\lambda\rho_u^{n+1/2} - (\lambda+1)w^n - \rho_\mu^{n+1/2}, D_t\theta_u^{n+1}) - \Delta t(\nabla w^n, \nabla D_t\theta_u^{n+1})$$
$$- \Delta t y^n \left( \frac{g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}}, D_t\theta_u^{n+1} \right) + \Delta t(\hat{E}_\mu^n, D_t\theta_u^{n+1}),$$

$$M_4 = -\frac{r^{n+1} + r^n}{2}\Delta t \left( \frac{g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}} - \frac{g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}}, D_t\theta_u^{n+1} \right) + \frac{e^{n+1} + e^n}{2}\Delta t \cdot$$
$$\left( \frac{g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}} - \frac{g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}}, D_t u^{n+1} \right) + \Delta t r_h^{n+1/2} \left( \frac{g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}} - \frac{g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}}, \theta_\mu^{n+1/2} \right),$$

$$M_5 = -\frac{e^{n+1} + e^n}{2} \Delta t \left( \frac{g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}}, D_t \rho_u^{n+1} \right) - (e^{n+1} + e^n) \Delta t \hat{E}_r^n + \Delta t \frac{e^{n+1/2}}{\sqrt{E_1(\bar{u}^n)}} (g(\bar{u}^n), \theta_\mu^{n+1/2}).$$

At this moment, we invoke a mathematical induction on (41). Since from (19), (22), (77) and (78) we have

$$\|u_h^n - u^n\|_{L^\infty} \le \|\theta_u^n\|_{L^\infty} + \|u^n - R_h u^n\|_{L^\infty} \le C(h^k + \Delta t^{1/2}), \quad n = 0, 1,$$

it implies that when $n = 0, 1$ there exist positive constants $h_2$, $\Delta t_2$ such that (41) holds provided $h < h_2$ and $\Delta t < \Delta t_2$. In the following, we present estimates of the finite element solution by assuming that (41) holds for $0 \le n \le m$, for some integer $m \ge 1$. We shall see that if (41) holds for $0 \le n \le m$, then it also holds for $n = m + 1$.

Now we analyze each term on the righthand side of (80). By (68) and (72),

$$M_2 \le Ch^{2k} \int_{t^n}^{t^{n+1}} \|u_t\|_{\max(1,k-1)}^2 dt + C\Delta t h^{2k+2} \left( \|\mu^{n+1/2}\|_{k+1}^2 + \|u^{n+1/2}\|_{k+1}^2 \right)$$

$$+ C\Delta t^4 \int_{t^n}^{t^{n+1}} \|u_{ttt}\|_{-1}^2 + \|u_{tt}\|_1^2 dt + C\Delta t (\|\theta_u^{n+1}\|_1^2 + \|\theta_u^n\|_1^2) + C\Delta t^5$$

$$+ \frac{\Delta t}{16} \|\theta_\mu^{n+1/2}\|_1^2,$$

where we have used the fact that

$$\|w^n\|_s^2 = \left\| R_h \left( \frac{u^{n+1} + u^n}{2} - u^{n+1/2} \right) \right\|_s^2$$

$$\le \left\| \frac{u^{n+1} + u^n}{2} - u^{n+1/2} \right\|_s^2 \le C\Delta t^3 \int_{t^n}^{t^{n+1}} \|u_{tt}\|_s^2 dt. \qquad (81)$$

Since it follows from (74) that

$$D_t \theta_u^{n+1} = \Delta_h \theta_\mu^{n+1/2} + P_h \left( D_t \rho_u^{n+1} - \hat{E}_u^n \right)$$

together with (17), (25), (72) and (23) we can derive

$$\|D_t \theta_u^{n+1}\|_{-1} \le \|\nabla \theta_\mu^{n+1/2}\| + \|D_t \rho_u^{n+1}\|_{-1} + \|\hat{E}_u^n\|_{-1}$$

$$\le \|\nabla \theta_\mu^{n+1/2}\| + Ch^k (\Delta t)^{-1} \int_{t^n}^{t^{n+1}} \|u_t\|_{\max(1,k-1)} dt$$

$$+ C\Delta t \int_{t^n}^{t^{n+1}} \|u_{ttt}\|_{-1} dt. \qquad (82)$$

Then by (24), (18), (81), (68) and (82) we have

$$M_3 \leq C\Delta t(\|\rho_u^{n+1/2}\|_1^2 + \|w^n\|_1^2 + \|\rho_\mu^{n+1/2}\|_1^2 + \|\Delta w^n\|_1^2 + |y^n|^2 + \|\hat{E}_\mu^n\|_1^2)$$

$$+ \frac{\Delta t}{16}\|D_t\theta_u^{n+1}\|_{-1}^2$$

$$\leq C\Delta t h^{2k}(\|u^{n+1/2}\|_{k+1}^2 + \|\mu^{n+1/2}\|_{k+1}^2) + C\Delta t^5$$

$$+ C\Delta t^4 \int_{t^n}^{t^{n+1}} (r_{tt}^2 + \|u_{tt}\|_3^2)dt + X^n,$$

where like (81) we have used that

$$|y^n|^2 = \left|\frac{r^{n+1} + r^n}{2} - r^{n+1/2}\right|^2 \leq C\Delta t^3 \int_{t^n}^{t^{n+1}} r_{tt}^2 dt, \tag{83}$$

and denoted

$$X^n := \frac{\Delta t}{8}\left\|\theta_\mu^{n+1/2}\|_1^2 + Ch^{2k} \int_{t^n}^{t^{n+1}} \right\|u_t\|_{\max(1,k-1)}^2 dt$$

$$+ C\Delta t^4 \int_{t^n}^{t^{n+1}} \|u_{ttt}\|_{-1}^2 dt. \tag{84}$$

Now we estimate

$$\frac{\nabla g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}} - \frac{\nabla g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}} = \nabla g(\bar{u}^n)\frac{E_1(\bar{u}^n) - E_1(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)E_1(\bar{u}^n)}(E_1(\bar{u}^n) + E_1(\bar{u}_h^n))}$$

$$+ \frac{\nabla g(\bar{u}_h^n) - \nabla g(\bar{u}^n)}{\sqrt{E_1(\bar{u}_h^n)}} = B_1 + B_2. \tag{85}$$

The two terms on the righthand side are bounded by

$$\|B_1\| \leq C\|g\|_{W^{1,\infty}}\|\bar{u}^n - \bar{u}_h^n\| \leq C(\|\theta_u^n\| + \|\rho_u^n\| + \|\theta_u^{n-1}\| + \|\rho_u^{n-1}\|), \tag{86}$$

and

$$\|B_2\| \leq C\|\nabla g(\bar{u}_h^n) - \nabla g(\bar{u}^n)\|$$

$$\leq C\|(g'(\bar{u}_h^n) - g'(\bar{u}^n))\nabla\bar{u}^n\| + |g'(\bar{u}_h^n)|(\|\nabla\rho_u^n\| + \|\nabla\theta_u^n\| + \|\nabla\rho_u^{n-1}\|$$

$$+ \|\nabla\theta_u^{n-1}\|) \leq C(\|\rho_u^n\|_1 + \|\theta_u^n\|_1 + \|\rho_u^{n-1}\|_1 + \|\theta_u^{n-1}\|_1) \tag{87}$$

where we have used the fact that $E_1[u_h^{n-1}]$, $E_1[u_h^n] \geq C_0 > 0$ and $g'(u_h^{n-1})$, $g''(u_h^{n-1})$, $g'(u_h^n)$, $g''(u_h^n)$ have a uniform upper bound by (15), (16) and (41). Therefore, we combine (85), (86) and (87) to get

$$\left\| \frac{\nabla g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}} - \frac{\nabla g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}} \right\| \le C(\|\rho_u^n\|_1 + \|\theta_u^n\|_1 + \|\rho_u^{n-1}\|_1 + \|\theta_u^{n-1}\|_1)$$

$$\le Ch^k(\|u^n\|_{k+1} + \|u^{n-1}\|_{k+1}) + C(\|\theta_u^n\|_1 + \|\theta_u^{n-1}\|_1). \tag{88}$$

Hence, together with (66) we derive that

$$M_4 \le C\Delta t((e^{n+1})^2 + (e^n)^2) + C\Delta t \left\| \frac{g(\bar{u}_h^n)}{\sqrt{E_1(\bar{u}_h^n)}} - \frac{g(\bar{u}^n)}{\sqrt{E_1(\bar{u}^n)}} \right\|_1^2 + X^n$$

$$\le C\Delta t((e^{n+1})^2 + (e^n)^2 + \|\theta_u^n\|_1^2 + \|\theta_u^{n-1}\|_1^2)$$

$$+ C\Delta t h^{2k}(\|u^n\|_{k+1}^2 + \|u^{n-1}\|_{k+1}^2) + X^n.$$

Thanks to (66), (41), (73) and (83), we have

$$M_5 \le C\Delta t((e^{n+1})^2 + (e^n)^2)$$

$$+ C\left( \Delta t^4 \int_{t^n}^{t^{n+1}} (r_{tt}^2 + r_{ttt}^2)ds + \Delta t^4 \int_{t^n}^{t^{n+1}} \int_\Omega u_{ttt}^2 dxds + \Delta t^5 \right)$$

$$+ Ch^{2k} \int_{t^n}^{t^{n+1}} \|u_t\|_{\max(1,k-1)}^2 dt + \frac{\Delta t}{16} \|\theta_\mu^{n+1/2}\|_1^2.$$

Then it implies from (18), (81), (66), (68) and (88) that

$$M_1 \le \frac{\Delta t}{16} \|\theta_\mu^{n+1/2}\|_1^2 + C\Delta t \left( \|\theta_u^n\|_1^2 + \|\theta_u^{n+1}\|_1^2 + \|\theta_u^{n-1}\|_1^2 + (e^{n+1})^2 \right.$$

$$\left. + (e^n)^2 \right) + C\Delta t^4 \int_{t^n}^{t^{n+1}} (r_{tt}^2 + \|u_{tt}\|_3^2)dt$$

$$+ C\Delta t h^{2k}(\|\mu^{n+1/2}\|_{k+1}^2 + \|u^{n+1/2}\|_{k+1}^2 + \|u^n\|_{k+1}^2 + \|u^{n-1}\|_{k+1}^2).$$

We combine the estimate for each term $M_i$, $1 \le i \le 5$ and (84) to get

$$\frac{1}{2}(\|\theta_u^{n+1}\|_1^2 - \|\theta_u^n\|_1^2) + \frac{\Delta t}{16} \|\theta_\mu^{n+1/2}\|_1^2 + (e^{n+1})^2 - (e^n)^2$$

$$+ \frac{\Delta t}{4} \left( \|\Delta_h \theta_u^{n+1}\|^2 - \|\Delta_h \theta_u^n\|^2 + \|\Delta_h(\theta_u^{n+1} + \theta_u^n)\|^2 \right)$$

$$\le C\Delta t \left( \|\theta_u^n\|_1^2 + \|\theta_u^{n+1}\|_1^2 + \|\theta_u^{n-1}\|_1^2 + (e^{n+1})^2 + (e^n)^2 \right)$$

$$+ Ch^{2k} \int_{t^n}^{t^{n+1}} \|u_t\|_{\max(1,k-1)}^2 dt$$

$$+C\Delta t h^{2k}\left(\|\mu^{n+1/2}\|_{k+1}^2+\|u^n\|_{k+1}^2+\|u^{n-1}\|_{k+1}^2+\|u^{n+1/2}\|_{k+1}^2\right)$$

$$+C\Delta t^4\int_{t^n}^{t^{n+1}}r_{tt}^2+r_{ttt}^2+\|u_{tt}\|_3^2+\|u_{ttt}\|^2ds+C\Delta t^5.$$

By applying Gronwall's inequality and using (78), we have

$$\max_{1\le n\le m}\left(\|\theta_u^{n+1}\|_1^2+\Delta t\|\Delta_h\theta_u^{n+1}\|^2+(e^{n+1})^2\right)$$

$$+\Delta t\sum_{n=1}^m\left(\|\theta_\mu^{n+1/2}\|_1^2+\|\Delta_h(\theta_u^{n+1}+\theta_u^n)\|^2\right)$$

$$\le C_1\Delta t^3+C_1\Delta t h^{2k}+\underline{C_1\Delta t^2\left(\int_0^{\Delta t}\|u_{tt}(s)\|_{-1}^2ds\right)+C_2\Delta t^4}$$

$$\underline{+C_2h^{2k}\left(\|u\|_{H^1(H^{\max(1,k-1)})}+\|\mu\|_{L^2(H^{k+1})}+\|u\|_{L^2(H^{k+1})}\right)^2}\qquad(89)$$

where the underlined part is the error by the first step.

Hence, when $h<h_2$ and $\Delta t<\Delta t_2$ and if $\Delta t\le h$, from the inverse estimate (c.f. [19, Lemma 6.4]) we have

$$\|\theta_u^{m+1}\|_{L^\infty}^2\le C\log(1/h)\|\theta_u^{m+1}\|_1^2\le Ch^{-1}(h^{2k}+\Delta t^2)\le C(h^{2k-1}+\Delta t).$$

On the other hand, if $\Delta t\ge h$, then we have

$$\|\Delta_h\theta_u^{m+1}\|^2\le C\left(\frac{h^{2k}}{\Delta t}+\Delta t\right)\le C(h^{2k-1}+\Delta t).$$

Overall, together with (19) we have $\|\theta_u^{m+1}\|_{L^\infty}\le C(h^{k-1/2}+\Delta t^{1/2})$, and so

$$\|u_h^{m+1}-u^{m+1}\|_{L^\infty}\le\|\theta_u^{m+1}\|_{L^\infty}+\|u^{m+1}-R_hu^{m+1}\|_{L^\infty}\le C(h^{k-1/2}+\Delta t^{1/2}),$$

where $C$ is independent of $m$. Thus there exists positive constants $h_3,\Delta t_3$ such that when $h<h_3$ and $\Delta t<\Delta t_3$ we have

$$\|u_h^{m+1}-u^{m+1}\|_{L^\infty}\le 1,$$

and this completes the mathematical induction on (41) in the case that $h<h_3$ and $\Delta t<\Delta t_3$. Thus (89) holds for $m=N-1$ with the same constant $C_1,C_2$, provided $h<h_3$ and $\Delta t<\Delta t_3$.

If $h\ge h_3$ or $\Delta t\ge\Delta t_3$, from (66) we see that

$$\max_{0 \le n \le N-1} \left( \|\theta_u^{n+1}\|_1^2 + (e^{n+1})^2 \right) + \Delta t \sum_{n=0}^{N-1} \|\nabla \theta_\mu^{n+1/2}\|^2$$

$$\le C_3 \le C_3(\Delta t_3^{-4} + h_3^{-2k}) \left( \Delta t^4 + h^{2k} + C_1 \Delta t^2 \left( \int_0^{\Delta t} \|u_{tt}(s)\|_{-1}^2 ds \right) \right)$$

$$(90)$$

for some positive constant $C_3$. From (89) and (90) we obtain for any $\Delta t$ and $h$,

$$\max_{0 \le n \le N-1} \left( \|\theta_u^{n+1}\|_1^2 + (e^{n+1})^2 \right) + \Delta t \sum_{n=0}^{N-1} \|\nabla \theta_\mu^{n+1/2}\|^2$$

$$\le (C_1 + C_2 + C_3(\Delta t_3^{-4} + h_3^{-2k})) \left( \Delta t^4 + h^{2k} + \Delta t^2 \left( \int_0^{\Delta t} \|u_{tt}(s)\|_{-1}^2 ds \right) \right)$$

from which we obtain the desired result by (21) and the triangle inequality. □

Thanks to the fact that (66) and (41) also hold for this second order scheme, following a similar procedure as in the proof of (59) and Theorem 2, we can also get the following optimal error estimates in $L^2$ norm for the finite element approximation.

**Theorem 4** *Under the same assumptions as in Theorem 3, we have*

$$\|u - u_h\|_{l^\infty(L^2)} + \|\mu - \mu_h\|_{\hat{l}^2(L^2)} + |r - r_h|_{l^\infty}$$

$$\le Ch^{k+1} \left( \|u\|_{H^1(H^{\max(k,2)})} + \|u\|_{L^2(H^{k+3})} + \|u\|_{L^\infty(H^{k+1})} \right)$$

$$+ C\Delta t^2 + C\Delta t \left( \int_0^{\Delta t} \|u_{tt}(s)\|_{-1}^2 ds \right)^{1/2}.$$

For other second-order SAV schemes, such as the one based on BDF2 method, one can derive similar error estimates as we did in this section.

## 6 Numerical experiments

In this section, we present several numerical experiments to validate our theoretical estimates. Since the time discretization errors of the SAV approach have been examined previously in [15], we shall concentrate on the spacial discretization errors.

We consider the Eq. (1) with boundary condition (2). First, we choose $\Omega = (0, 2\pi) \times (0, 2\pi)$ and take the initial condition to be

$$u_0(x, y) = 0.05 \cos x \cos y. \tag{91}$$

Since the exact solution is unknown, we take the numerical solution computed by the cubic elements on the finest grid $128 \times 128$ as the reference.

First, we consider the first-order scheme (6)–(8), and use Lagrange elements of degree $r$ for both $u$ and $\mu$. Then we compute the example by taking $\Delta t = 10^{-4}$, $T = 1$,

**Table 1** Compute with $P^1$ on square by backward Euler method at final time ($h = 3.93E{-}01$)

| Mesh size | $\|u - u_h\|$ | Rate | $\|\nabla(u - u_h)\|$ | Rate |
|-----------|---------------|------|------------------------|------|
| h | 4.90E−03 | 1.95 | 4.33E−02 | 0.97 |
| h/2 | 1.24E−03 | 1.98 | 2.18E−02 | 0.99 |
| h/4 | 3,12E−04 | 1.99 | 1.09E−02 | 1.00 |
| h/8 | 7.78E−05 | 2.00 | 5.45E−03 | 1.00 |

**Table 2** Compute with $P^2$ on square by backward Euler method at final time ($h = 3.93E - 01$)

| Mesh size | $\|u - u_h\|$ | Rate | $\|\nabla(u - u_h)\|$ | Rate |
|-----------|---------------|------|------------------------|------|
| h | 1.56E−04 | 2.75 | 3.42E−03 | 1.97 |
| h/2 | 2.09E−05 | 2.90 | 8.46E−04 | 2.01 |
| h/4 | 2.68E−06 | 2.96 | 2.11E−04 | 2.01 |
| h/8 | 3.37E−07 | 2.99 | 5.27E−05 | 2.00 |

**Table 3** Compute with $P^1$ on square by C–N method at final time ($h = 3.93E - 01$)

| Mesh size | $\|u - u_h\|$ | Rate | $\|\nabla(u - u_h)\|$ | Rate |
|-----------|---------------|------|------------------------|------|
| h | 4.96E−03 | 1.95 | 4.33E−02 | 0.96 |
| h/2 | 1.25E−03 | 1.99 | 2.18E−02 | 0.99 |
| h/4 | 3.14E−04 | 1.99 | 1.09E−02 | 1.00 |
| h/8 | 7.91E−05 | 1.99 | 5.47E−03 | 1.00 |

**Table 4** Compute with $P^2$ on square by C–N method at final time ($h = 3.93E - 01$)

| Mesh size | $\|u - u_h\|$ | Rate | $\|\nabla(u - u_h)\|$ | Rate |
|-----------|---------------|------|------------------------|------|
| h | 1.56E−04 | 2.75 | 3.42E−03 | 1.97 |
| h/2 | 2.09E−05 | 2.90 | 8.47E−04 | 2.01 |
| h/4 | 2.68E−06 | 2.96 | 2.11E−04 | 2.01 |
| h/8 | 3.37E−07 | 2.99 | 5.27E−05 | 2.00 |

where the error due to the time discretization should be much smaller than the one due to the spatial discretization since the time step is small enough. Tables 1 and 2 show the computation results for $r = 1$ and $r = 2$ respectively, from which we can find that the convergence rate in each norm is as predicted by the theory.

Next, we compute the same example by the SAV method based on Crank–Nicolson scheme, that is (62)–(65). Here we take $\Delta t = 10^{-4}$, $T = 1$ and compute with $u$, $\mu$ discretized by $P^1$ (or $P^2$) Lagrange element spaces. The corresponding numerical results are shown in Tables 3 and 4 respectively, which also show the convergence rate in each norm is just as predicted by the theory. Again, we observe the predicted convergence rates.

# References

1. Allen, S.M., Cahn, J.W.: A microscopic theory for antiphase boundary motion and its application to antiphase domain coarsening. Acta Metall. **27**(6), 1085–1095 (1979)
2. Baskaran, A., Lowengrub, J.S., Wang, C., Wise, S.M.: Convergence analysis of a second order convex splitting scheme for the modified phase field crystal equation. SIAM J. Numer. Anal. **51**(5), 2851–2873 (2013)
3. Brenner, S., Scott, R.: The Mathematical Theory of Finite Element Methods. Texts in Applied Mathematics, 3rd edn. Springer, New York (2010)
4. Cahn, J.W., Hilliard, J.E.: Free energy of a nonuniform system. I: interfacial free energy. J. Chem. Phys. **28**, 258 (1958)
5. Ciarlet, P.G.: The Finite Element Method for Elliptic Problems, Studies in Mathematics and its Applications, vol. 4. North-Holland Publishing Co., Amsterdam (1978)
6. Condette, N., Melcher, C., Süli, E.: Spectral approximation of patternforming nonlinear evolution equations with double-well potentials of quadratic growth. Math. Comput. **80**(273), 205–223 (2011)
7. Du, Q., Nicolaides, R.A.: Numerical analysis of a continuum model of phase transition. SIAM J. Numer. Anal. **28**(5), 1310–1322 (1991)
8. Elliott, C.M., Stuart, A.M.: The global dynamics of discrete semilinear parabolic equations. SIAM J. Numer. Anal. **30**(6), 1622–1663 (1993)
9. Feng, X., He, Y., Liu, C.: Analysis of finite element approximations of a phase field model for two-phase fluids. Math. Comput. **76**, 539–571 (2007)
10. Feng, X., Prohl, A.: Error analysis of a mixed finite element method for the Cahn–Hilliard equation. Numerische Mathematik **99**(1), 47–84 (2004)
11. Kay, D., Styles, V., Süli, E.: Discontinuous galerkin finite element approximation of the Cahn–Hilliard equation with convection. SIAM J. Numer. Anal. **47**(4), 2660–2685 (2009)
12. Kessler, D., Nochetto, R.H., Schmidt, A.: A posteriori error control for the Allen–Cahn problem: circumventing Gronwall's inequality. ESAIM: Math. Modell. Numer. Anal. **38**(1), 129–142 (2004)
13. Liu, Y., Chen, W., Wang, C., Wise, S.M.: Error analysis of a mixed finite element method for a Cahn–Hilliard–Hele–Shaw system. Numer. Math. **135**(3), 679–709 (2017)
14. Quarteroni, A., Valli, A.: Numerical Approximation of Partial Differential Equations. Springer, Berlin (2008)
15. Shen, J., Xu, J.: Convergence and error analysis for the scalar auxiliary variable (SAV) schemes to gradient flows. SIAM J. Numer. Anal. **56**(5), 2895–2912 (2018)
16. Shen, J., Jie, X., Yang, J.: The scalar auxiliary variable (SAV) approach for gradient flows. J. Comput. Phys. **353**, 407–416 (2018)
17. Shen, J., Jie, X., Yang, J.: A new class of efficient and robust energy stable schemes for gradient flows. SIAM Rev. **61**, 474–506 (2019)
18. Temam, R.: Infinite-Dimensional Dynamical Systems in Mechanics and Physics, 2nd edn. Springer, Berlin (1997)
19. Thomée, V.: Galerkin Finite Element Methods for Parabolic Problems. Springer Series in Computational Mathematics, vol. 25, 2nd edn. Springer, Berlin (2006)
20. Yang, X.: Linear, first and second-order, unconditionally energy stable numerical schemes for the phase field model of homopolymer blends. J. Comput. Phys. **327**, 294–316 (2016)
21. Yang, X., Ju, L.: Linear and unconditionally energy stable schemes for the binary fluid-surfactant phase field model. Comput. Methods Appl. Mech. Eng. **318**, 1005–1029 (2017)
22. Yang, X., Zhao, J., Wang, Q., Shen, J.: Numerical approximations for a three components Cahn–Hilliard phase-field model based on the invariant energy quadratization method. Math. Models Methods Appl. Sci. **27**(11), 1993–2030 (2017)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.