

RIMS-1754

**On the a posteriori estimates for inverse operators of  
linear parabolic equations with applications to  
the numerical enclosure of solutions for nonlinear problems**

By

Takehiko KINOSHITA, Takuma KIMURA,  
and Mitsuhiro T. NAKAO

July 2012



**京都大学 数理解析研究所**

RESEARCH INSTITUTE FOR MATHEMATICAL SCIENCES

KYOTO UNIVERSITY, Kyoto, Japan

# On the a posteriori estimates for inverse operators of linear parabolic equations with applications to the numerical enclosure of solutions for nonlinear problems

Takehiko Kinoshita · Takuma Kimura ·  
Mitsuhiro T. Nakao

July 24, 2012

**Abstract** We consider the guaranteed a posteriori estimates for the inverse parabolic operators with homogeneous initial-boundary conditions. Our estimation technique uses a full-discrete numerical scheme, which is based on the Galerkin method with an interpolation in time by using the fundamental solution for semidiscretization in space. In our technique, the constructive a priori error estimates for a full discretization of solutions for the heat equation play an essential role. Combining these estimates with an argument for the discretized inverse operator and a contraction property of the Newton-type formulation, we derive an a posteriori estimate of the norm for the infinite-dimensional operator. In numerical examples, we show that the proposed method should be more efficient than the existing method. Moreover, as an application, we give some prototype results for numerical verification of solutions of nonlinear parabolic problems, which confirm the actual usefulness of our technique.

**Keywords** Parabolic PDEs · Galerkin methods · A posteriori estimates · Numerical verification methods

**Mathematics Subject Classification (2000)** 35K20 · 65M15 · 65M60

## 1 Introduction

Setting  $\mathcal{L}_t := \frac{\partial}{\partial t} - v\Delta + b \cdot \nabla + c$ , for  $f \in L^2(J; L^2(\Omega))$ , consider the following linear parabolic partial differential equations (PDEs) with homogeneous initial and bound-

---

Takehiko Kinoshita  
Research Institute for Mathematical Sciences, Kyoto University, Kyoto 606-8502, Japan  
E-mail: kinosita@kurims.kyoto-u.ac.jp

Takuma Kimura  
JST CREST / Faculty of Science and Engineering, Waseda University, Tokyo 169-8555, Japan  
E-mail: tkimura@aoni.waseda.jp

Mitsuhiro T. Nakao  
Sasebo National College of Technology, Nagasaki 857-1193, Japan  
E-mail: mtakao@post.cc.sasebo.ac.jp

ary conditions:

$$\begin{cases} \mathcal{L}_t u = f, & \text{in } \Omega \times J, & (1a) \\ u(x, t) = 0, & \text{on } \partial\Omega \times J, & (1b) \\ u(x, 0) = 0, & \text{in } \Omega, & (1c) \end{cases}$$

where  $\Omega \subset \mathbb{R}^d$ , ( $d \in \{1, 2, 3\}$ ) is a bounded polygonal or polyhedral domain,  $J := (0, T) \subset \mathbb{R}$ , ( $T < \infty$ ) is a bounded interval,  $\nu$  is a positive constant,  $b \in L^\infty(J; L^\infty(\Omega))^d$ , and  $c \in L^\infty(J; L^\infty(\Omega))$ . As is well known, for any  $f \in L^2(J; L^2(\Omega))$ , there exists a unique weak solution  $u \in L^2(J; H_0^1(\Omega))$  to the problem (1). Denoting the solution operator of (1) by  $\mathcal{L}_t^{-1}$ , it is a bounded linear operator from  $L^2(J; L^2(\Omega))$  to  $L^2(J; H_0^1(\Omega))$ .

The main aim of this paper is to obtain the concrete value  $C_{L^2L^2, L^2H_0^1} > 0$  satisfying the following estimates:

$$\|\mathcal{L}_t^{-1}\|_{\mathcal{L}(L^2(J; L^2(\Omega)), L^2(J; H_0^1(\Omega)))} \leq C_{L^2L^2, L^2H_0^1}. \quad (2)$$

The constant  $C_{L^2L^2, L^2H_0^1}$  plays an important role in the verification of solutions for the initial-boundary-value problems for the nonlinear parabolic PDEs, and we usually need to estimate it as small as possible. The concrete value  $C_{L^2L^2, L^2H_0^1} > 0$  satisfying (2) can be calculated by the Gronwall inequality or other theoretical considerations (e.g., [16]), which we call the “*a priori estimates*.” However, in general,  $C_{L^2L^2, L^2H_0^1}$  obtained by such a priori estimates is exponentially dependent on the length of the time interval  $J$  unless the corresponding elliptic part of the operator  $\mathcal{L}_t$  is coercive [4, 5]. Thus a priori estimates often lead to an overestimate for the norm of  $\mathcal{L}_t^{-1}$ , which yields worse results for some purposes.

In order to overcome this difficulty, we proposed a method to calculate  $C_{L^2L^2, L^2H_0^1}$  by numerical computation with guaranteed accuracy in [10], which we called “*a posteriori estimates*.” The method is based on combining the a priori error estimates for a semidiscretization with the a priori estimates for the ordinary differential equations (ODEs) in time. It has proven to be more efficient than the existing a priori method; some numerical examples show that this a posteriori method can remove the exponential dependency on the time interval  $J$ . However, it has a very large computational cost, because the semidiscretization of (1) causes stiff ODEs that require a very small step size. Also, it is not clear what time-space ratio to use in the discretization process.

In this paper, we propose a new a posteriori method with a fully discretized Newton-type operator, which uses the Galerkin approximation in the space direction and the Lagrange-type interpolation in the time direction. In the case of the simple heat equations, some fundamental properties (e.g., the stability and a priori error estimates) for this full-discretization scheme have already been obtained in [11]. In the desired estimation of the inverse operator norm  $\|\mathcal{L}_t^{-1}\|_{\mathcal{L}(L^2(J; L^2(\Omega)), L^2(J; H_0^1(\Omega)))}$ , the matrix norm estimates corresponding to the discretized inverse operator and the constructive error analysis for the simple heat equations are important and essential. By constructive analysis, we can also guess an appropriate time-space ratio prior to the actual computation. Moreover, by using numerical examples, we will show that

the proposed method succeeds in obtaining a posteriori estimates with less computational cost than the previous method in [10]. This means that the present method is very robust compared with the previous one.

The contents of this paper are as follows: In section 2, we introduce some function spaces, operators, and other notation. In section 3, we introduce the results of stability and a priori error estimates for the full-discretization scheme for the simple heat equations, which were obtained in [11]. In section 4, we consider the approximate quasi-Newton operator that corresponds to the full-discretization scheme for problem (1). In section 5, we derive the new a posteriori estimates of (2) by combining the results in section 3 with the property of the approximate quasi-Newton operator defined in the previous section. In section 6, we compare the computed values for  $C_{L^2L^2, L^2H_0^1}$  by three methods, namely, the a priori method, the a posteriori estimates in [10], and the new a posteriori method obtained in section 5. In this section we also show some prototype results of the numerical enclosure of solutions for nonlinear parabolic problems as an application of our method.

## 2 Notation

In this section, we introduce some function spaces, operators, and other notation. Let  $L^2(\Omega)$  and  $H^1(\Omega)$  be the usual Lebesgue and Sobolev spaces on  $\Omega$ , respectively, and define the natural inner product of  $u, v$  in  $L^2(\Omega)$  by  $(u, v)_{L^2(\Omega)} := \int_{\Omega} u(x)v(x) dx$ . Also, let  $H_0^1(\Omega)$  be a Sobolev space defined by  $H_0^1(\Omega) := \{u \in H^1(\Omega) ; u = 0 \text{ on } \partial\Omega\}$  with inner product  $(u, v)_{H_0^1(\Omega)} := (\nabla u, \nabla v)_{L^2(\Omega)^d}$ . We will sometimes refer to the following Sobolev inequality on  $H_0^1(\Omega)$ . Namely, for a suitable constant  $p \geq 1$ , which is dependent on the dimension of  $\Omega$ , there exists a constant  $C_{s,p} > 0$  such that

$$\|u\|_{L^p(\Omega)} \leq C_{s,p} \|u\|_{H_0^1(\Omega)}, \quad \forall u \in H_0^1(\Omega). \quad (3)$$

When  $p = 2$ , (3) is called the Poincaré inequality.

Let  $\Delta : L^2(\Omega) \rightarrow L^2(\Omega)$  be the Laplace operator that is self-adjoint on the domain  $D(\Delta) := \{u \in H_0^1(\Omega) ; \Delta u \in L^2(\Omega)\}$ . Let  $V^1(J)$  be a subspace of  $H^1(J)$  defined by  $V^1(J) := \{u \in H^1(J) ; u(0) = 0\}$ . Then,  $V^1(J)$  is a Hilbert space with inner product  $(u, v)_{V^1(J)} := (u', v')_{L^2(J)}$ . The time-dependent Lebesgue space  $L^2(J; L^2(\Omega))$  is defined as a space of square-integrable  $L^2(\Omega)$ -valued functions on  $J$ . Then,  $L^2(J; L^2(\Omega))$  is a Hilbert space with inner product  $(u, v)_{L^2(J; L^2(\Omega))} := \int_J \int_{\Omega} u(x, t)v(x, t) dx dt$ . We denote the function space  $L^2(J; L^2(\Omega))$  as  $L^2L^2$ , for short. Let  $L^2(J; H_0^1(\Omega))$  be a subspace of  $L^2L^2$  defined by

$$L^2(J; H_0^1(\Omega)) := \left\{ u \in L^2L^2 ; \nabla u \in L^2(J; L^2(\Omega))^d, u(\cdot, t) = 0 \text{ on } \partial\Omega, \text{ a.e. } t \in J \right\}.$$

Then,  $L^2H_0^1 \equiv L^2(J; H_0^1(\Omega))$  is a Hilbert space with inner product  $(u, v)_{L^2H_0^1} := (\nabla u, \nabla v)_{(L^2L^2)^d}$ . Let  $V^1(J; L^2(\Omega))$  be a subspace of  $L^2L^2$  defined by

$$V^1(J; L^2(\Omega)) := \left\{ u \in L^2(J; L^2(\Omega)) ; \frac{\partial u}{\partial t} \in L^2(J; L^2(\Omega)), u(\cdot, 0) = 0 \text{ in } L^2(\Omega) \right\}.$$

Then,  $V^1L^2 \equiv V^1(J; L^2(\Omega))$  is a Hilbert space with inner product  $(u, v)_{V^1L^2} := \left( \frac{\partial u}{\partial t}, \frac{\partial v}{\partial t} \right)_{L^2L^2}$ .

We define the Hilbert space  $V := V^1L^2 \cap L^2H_0^1$  with inner product  $(u, v)_V := (u, v)_{V^1L^2} + (u, v)_{L^2H_0^1} = \left( \frac{\partial u}{\partial t}, \frac{\partial v}{\partial t} \right)_{L^2L^2} + (\nabla u, \nabla v)_{(L^2L^2)^d}$ . Moreover, we define the partial differential operator  $\Delta_t : L^2L^2 \rightarrow L^2L^2$  by  $\Delta_t := \frac{\partial}{\partial t} - v\Delta$  on the domain  $D(\Delta_t) := V^1L^2 \cap L^2(J; D(\Delta))$ . Then, the inverse of  $\Delta_t$  exists (e.g., [3]), and we denote it by  $\Delta_t^{-1} \in \mathcal{L}(L^2L^2)$ . Notably, the range of  $\Delta_t^{-1}$  satisfies  $R(\Delta_t^{-1}) = D(\Delta_t)$ . From the compactness of the embedding  $I_e : D(\Delta_t) \hookrightarrow L^2H_0^1$ , the bounded linear operator  $I_e\Delta_t^{-1} \in \mathcal{L}(L^2L^2, L^2H_0^1)$  is also compact.

Let  $S_h(\Omega)$  be a finite-dimensional subspace of  $H_0^1(\Omega)$  dependent on the discretization parameter  $h$ . For example,  $S_h(\Omega)$  is considered to be a finite element space with mesh size  $h$ . Let  $n$  be the number of degrees of freedom of  $S_h(\Omega)$ , and let  $\{\phi_i\}_{i=1}^n \subset H_0^1(\Omega)$  be the basis functions of  $S_h(\Omega)$ . Moreover, we denote a vector of the basis functions of  $S_h(\Omega)$  by  $\phi := (\phi_1, \dots, \phi_n)^T$ . We also assume the inverse estimates on  $S_h(\Omega)$  like as follows:

**Assumption 2.1** *There exists a positive constant  $C_{inv}(h)$  satisfying*

$$\|u_h\|_{H_0^1(\Omega)} \leq C_{inv}(h) \|u_h\|_{L^2(\Omega)}, \quad \forall u_h \in S_h(\Omega). \quad (4)$$

For example, if  $\Omega$  is a bounded open interval in  $\mathbb{R}$ , and  $S_h(\Omega)$  is the P1 finite element space, then Assumption 2.1 is realized with  $C_{inv}(h) = \frac{\sqrt{12}}{h_{\min}}$ , where  $h_{\min}$  is the minimum mesh size in the division of  $\Omega$  (see e.g., [15, Theorem 1.5]).

Let  $P_h^1 : H_0^1(\Omega) \rightarrow S_h(\Omega)$  be an  $H_0^1$ -projection. Namely, for an arbitrary element  $u \in H_0^1(\Omega)$ ,  $P_h^1 u \in S_h(\Omega)$  satisfies the following variational equation:

$$(\nabla(u - P_h^1 u), \nabla v_h)_{L^2(\Omega)^d} = 0, \quad \forall v_h \in S_h(\Omega). \quad (5)$$

We need the following assumptions as the a priori error estimates for  $P_h^1$ .

**Assumption 2.2** *There exists a positive constant  $C_\Omega(h)$  satisfying*

$$\|u - P_h^1 u\|_{H_0^1(\Omega)} \leq C_\Omega(h) \|\Delta u\|_{L^2(\Omega)}, \quad \forall u \in D(\Delta), \quad (6)$$

$$\|u - P_h^1 u\|_{L^2(\Omega)} \leq C_\Omega(h) \|u - P_h^1 u\|_{H_0^1(\Omega)}, \quad \forall u \in H_0^1(\Omega). \quad (7)$$

For example, if  $\Omega$  is a bounded open interval in  $\mathbb{R}$ , and  $S_h(\Omega)$  is the P1 finite element space, then Assumption 2.2 is realized as  $C_\Omega(h) = \frac{h}{\pi}$ , where  $h$  is the mesh size (see e.g., [1, 7]).

Let  $V_k^1(J)$  be a finite-dimensional subspace of  $V^1(J)$  dependent on the discretization parameter  $k$ . For example,  $V_k^1(J)$  is considered to be a finite element space with mesh size (time step size)  $k$ . Let  $m$  be the number of degrees of freedom for  $V_k^1(J)$ , and let  $\{\psi_i\}_{i=1}^m \subset V^1(J)$  be the basis functions of  $V_k^1(J)$ . Moreover, we denote a vector of the basis functions of  $V_k^1(J)$  by  $\psi := (\psi_1, \dots, \psi_m)^T$ .

We assume that  $\Pi_k : V^1(J) \rightarrow V_k^1(J)$  is a Lagrange interpolation operator. Namely, if the mesh points on  $J$  are taken as  $0 = t_0 < t_1 < \dots < t_m = T$ , for any element  $u \in V^1(J)$ ,  $\Pi_k u \in V_k^1(J)$  satisfies

$$u(t_i) = (\Pi_k u)(t_i), \quad \forall i \in \{1, \dots, m\}. \quad (8)$$

We need the following assumption as the a priori error estimate for  $\Pi_k$ .

**Assumption 2.3** *There exists a positive constant  $C_J(k)$  satisfying*

$$\|u - \Pi_k u\|_{L^2(J)} \leq C_J(k) \|u\|_{V^1(J)}, \quad \forall u \in V^1(J). \quad (9)$$

For example, if  $V_k^1(J)$  is the P1 finite element space, then Assumption 2.3 is realized by  $C_J(k) = \frac{k}{\pi}$  (see e.g., [15, Theorem 2.4]).

Let  $V^1(J; S_h(\Omega))$  and  $V_k^1(J; S_h(\Omega))$  be the semidiscretization and the full-discretization subspaces of  $V$ , respectively. We now define the semidiscretization operator  $P_h : V \rightarrow V^1(J; S_h(\Omega))$  by the following weak form for any  $u \in V$

$$\left( \frac{\partial}{\partial t} (u - P_h u)(t), v_h \right)_{L^2(\Omega)} + v (\nabla(u - P_h u)(t), \nabla v_h)_{L^2(\Omega)^d} = 0, \quad \forall v_h \in S_h(\Omega), \text{ a.e. } t \in J. \quad (10)$$

Then the full-discretization operator  $P_{h,k} : V \rightarrow V_k^1(J; S_h(\Omega))$  is defined as the composition of  $P_h$  and  $\Pi_k$ , that is, by  $P_{h,k} := \Pi_k P_h$ .

### 3 Constructive a priori error estimates

In this section, we introduce some results for the stability of, and a priori error estimates for, the full-discretization operator  $P_{h,k}$ . Since the results of this section are given in [11], we omit the proofs.

**Theorem 3.1 ([11, Lemma 5.3 & Theorem 5.4])** *Under Assumption 2.1 and Assumption 2.3, the following constructive a priori estimate holds,*

$$\|P_{h,k} u\|_{L^2(J; H_0^1(\Omega))} \leq \left( \frac{C_{s,2}}{v} + C_{inv}(h) C_J(k) \right) \left\| \frac{\partial u}{\partial t} - v \Delta u \right\|_{L^2 L^2}, \quad \forall u \in D(\Delta_t). \quad (11)$$

Moreover, if  $V_k^1(J)$  is the P1 finite element space then we have the following estimates:

$$\|P_{h,k} u\|_{V^1(J; L^2(\Omega))} \leq 2 \left\| \frac{\partial u}{\partial t} - v \Delta u \right\|_{L^2(J; L^2(\Omega))}, \quad \forall u \in D(\Delta_t). \quad (12)$$

Since the full-discretization scheme proposed in [6, 9] has no  $V^1 L^2$  stability, we can say that the present full-discretized approximation has better properties, in an analytical and practical sense.

Finally, we introduce the constructive a priori error estimates for  $P_{h,k}$ .

**Theorem 3.2 ([11, Theorem 5.5 & Theorem 5.6])** *Under the assumptions 2.1- 2.3, we have the following constructive a priori error estimates:*

$$\|u - P_{h,k} u\|_{L^2(J; H_0^1(\Omega))} \leq C_1(h, k) \left\| \frac{\partial u}{\partial t} - v \Delta u \right\|_{L^2(J; L^2(\Omega))}, \quad \forall u \in D(\Delta_t), \quad (13)$$

$$\|u - P_{h,k} u\|_{L^2(J; L^2(\Omega))} \leq C_0(h, k) \left\| \frac{\partial u}{\partial t} - v \Delta u \right\|_{L^2(J; L^2(\Omega))}, \quad \forall u \in D(\Delta_t), \quad (14)$$

where  $C_1(h, k) := \frac{2}{v} C_\Omega(h) + C_{inv}(h) C_J(k)$  and  $C_0(h, k) = \frac{8}{v} C_\Omega(h)^2 + C_J(k)$ .

#### 4 Discretized quasi-Newton scheme

In this section, we consider a full-discretized approximation scheme for solutions of (1) by using a quasi-Newton operator. Since the full-discretization scheme in this paper uses interpolation in time, its computational method is somewhat complicated. However, it enables us to get an efficient and accurate estimation of the inverse operator norm in (2), as well as the verified computation of solutions to nonlinear problems.

We first describe an easy, but an important operation of matrix-vector multiplication.

**Definition 4.1** *Let  $M$  be an  $m_1$ -by- $m_2$  matrix. Then, we define the  $m_1 m_2$  vector  $\text{vec}(M)$  as follows:*

$$\text{vec}(M) := (M_{1,1}, M_{1,2}, \dots, M_{1,m_2}, M_{2,1}, \dots, M_{m_1,m_2})^T. \quad (15)$$

We call this transformation a “row-major matrix-vector transformation”.

**Definition 4.2** *Let  $M$  be an  $m_1$ -by- $m_2$  matrix. Then, we define the block diagonal matrix  $(I_n \otimes M)$  as follows:*

$$(I_n \otimes M) := \underbrace{\begin{pmatrix} M & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & M \end{pmatrix}}_n. \quad (16)$$

Here,  $I_n$  is the  $n$ -by- $n$  identity matrix, and the operator  $\otimes$  denotes the Kronecker product.

From these definitions, we have the following lemma.

**Lemma 4.3** *For an arbitrary  $n$ -by- $m$  matrix  $M$  and  $m$ -dimensional vector  $x$ , the following equality holds:*

$$Mx = (I_n \otimes x^T) \text{vec}(M). \quad (17)$$

**Proof.** — The elements of  $Mx$  are calculated by

$$Mx = \begin{pmatrix} M_{1,1} & \cdots & M_{1,m} \\ \vdots & \ddots & \vdots \\ M_{n,1} & \cdots & M_{n,m} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix} = \begin{pmatrix} M_{1,1}x_1 + \cdots + M_{1,m}x_m \\ \vdots \\ M_{n,1}x_1 + \cdots + M_{n,m}x_m \end{pmatrix}.$$

On the other hand, the elements of  $(I_n \otimes x^T) \text{vec}(M)$  are calculated by

$$(I_n \otimes x^T) \text{vec}(M) = \begin{pmatrix} x^T & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & x^T \end{pmatrix} \begin{pmatrix} M_{1,1} \\ \vdots \\ M_{n,m} \end{pmatrix} = \begin{pmatrix} M_{1,1}x_1 + \cdots + M_{1,m}x_m \\ \vdots \\ M_{n,1}x_1 + \cdots + M_{n,m}x_m \end{pmatrix}.$$

Therefore, the corresponding components coincide with each other.  $\square$

Next, we consider the quasi-Newton operator of (1) and its full-discretization. Let  $A$  be an integral operator defined by  $A := -I_e \Delta_t^{-1} (b \cdot \nabla + c) : L^2(J; H_0^1(\Omega)) \rightarrow L^2(J; H_0^1(\Omega))$ . Since the domain of  $\Delta_t$  is  $D(\Delta_t)$ , denoting the range of  $A$  by  $R(A)$ , it holds that  $R(A) \subset D(\Delta_t) = V^1 L^2 \cap L^2 D(\Delta)$ . Then, the differential operator of the left-hand side of (1a) can be represented as  $\mathcal{L}_t = \Delta_t (I - A)$ , where  $I$  denotes the identity operator on  $D(\Delta_t)$ . We define the quasi-Newton operator as the inverse of  $I - A$ , i.e.,  $(I - A)^{-1} : L^2 H_0^1 \rightarrow L^2 H_0^1$ .

We now define the symmetric and positive definite matrices  $L_\phi$  and  $D_\phi \in \mathbb{R}^{n \times n}$  by

$$L_{\phi,i,j} := (\phi_j, \phi_i)_{L^2(\Omega)}, \quad D_{\phi,i,j} := (\nabla \phi_j, \nabla \phi_i)_{L^2(\Omega)^d}, \quad \forall i, j \in \{1, \dots, n\}.$$

Let  $L_\phi^{1/2}$  and  $D_\phi^{1/2}$  be the Cholesky factors of  $L_\phi$  and  $D_\phi$ , respectively, i.e., the following equalities hold

$$L_\phi = L_\phi^{1/2} L_\phi^{T/2}, \quad D_\phi = D_\phi^{1/2} D_\phi^{T/2},$$

where  $L_\phi^{1/2}$  and  $D_\phi^{1/2}$  are lower triangular matrices, and  $L_\phi^{T/2}$  and  $D_\phi^{T/2}$  are those matrices transposed. Let  $L_\psi \in \mathbb{R}^{m \times m}$  be the symmetric and positive-definite matrix whose elements are defined by  $L_{\psi,i,j} := (\psi_j, \psi_i)_{L^2(J)}$ . We define  $Z_\phi \in L^\infty(J)^{n \times n}$  as the matrix function on  $J$  whose elements are defined by

$$Z_{\phi,i,j} := ((b \cdot \nabla) \phi_j + c \phi_j, \phi_i)_{L^2(\Omega)}, \quad \forall i, j \in \{1, \dots, n\}.$$

For any  $i \in \{1, \dots, m\}$ , we define the matrices  $\tilde{G}_{\phi,\psi}^{(i)} \in \mathbb{R}^{n \times nm}$  and  $\tilde{G}_{\phi,\psi} \in \mathbb{R}^{nm \times nm}$  by

$$\tilde{G}_{\phi,\psi}^{(i)} := \int_0^{t_i} \exp\left((s - t_i) \nu L_\phi^{-1} D_\phi\right) L_\phi^{-1} Z_\phi(s) (I_n \otimes \psi(s)^T) ds, \quad \tilde{G}_{\phi,\psi} := \begin{pmatrix} \tilde{G}_{\phi,\psi}^{(1)} \\ \vdots \\ \tilde{G}_{\phi,\psi}^{(m)} \end{pmatrix}. \quad (18)$$

Moreover, we define  $G_{\phi,\psi} \in \mathbb{R}^{nm \times nm}$  as  $G_{\phi,\psi} := I_{nm} - \tilde{G}_{\phi,\psi}$ .

We obtain the Theorem 4.4 as a full-discretization scheme of the quasi-Newton operator.

**Theorem 4.4** *Let  $V_k^1(J)$  be a finite element space constituted by the Lagrange elements. For a function  $f_{h,k} \in V_k^1(J; S_h(\Omega))$ , let  $u_{h,k} \in V_k^1(J; S_h(\Omega))$  be a solution of the following equation*

$$u_{h,k} - P_{h,k} A u_{h,k} = f_{h,k}. \quad (19)$$

*Then, the unique existence of a solution  $u_{h,k}$  of (19) is equivalent to the nonsingularity of  $G_{\phi,\psi}$ .*



**Proof.** — First, we consider  $P_{h,k}Au_{h,k}$ . For an arbitrary  $u_{h,k} \in V_k^1(J; S_k(\Omega))$ , there exists a matrix  $U \in \mathbb{R}^{n \times m}$  such that  $u_{h,k}(x, t) = \phi(x)^T U \psi(t)$ . Let  $w_h := P_h Au_{h,k}$ . Similarly, from  $w_h \in V^1(J; S_h(\Omega))$ , there exists a vector function  $\mathfrak{w} \in V^1(J)^n$  such that

$$w_h(x, t) = \phi(x)^T \mathfrak{w}(t) = \sum_{i=1}^n \phi_i(x) \mathfrak{w}_i(t).$$

For each  $v_h \in S_h(\Omega)$ , and almost everywhere  $t \in J$ , from the definition of  $P_h$  and the operator  $A$ , we have

$$\begin{aligned} \left( \frac{\partial w_h}{\partial t}(t), v_h \right)_{L^2(\Omega)} + \mathbf{v}(\nabla w_h(t), \nabla v_h)_{L^2(\Omega)^d} \\ = \left( \frac{\partial Au_{h,k}}{\partial t}(t), v_h \right)_{L^2(\Omega)} + \mathbf{v}(\nabla Au_{h,k}(t), \nabla v_h)_{L^2(\Omega)^d}, \\ = ((b(t) \cdot \nabla) u_{h,k}(t) + c(t) u_{h,k}(t), v_h)_{L^2(\Omega)}. \end{aligned} \quad (20)$$

From the arbitrariness of  $v_h \in S_h(\Omega)$ , the variational equation (20) is equivalent to the following system of first-order linear ODEs with homogeneous initial conditions

$$\left( L_\phi \frac{d}{dt} + \mathbf{v} D_\phi \right) \mathfrak{w} = Z_\phi U \psi. \quad (21)$$

Since (21) is an initial-value problem for an ODE system with constant coefficients, by using its fundamental matrix,  $\mathfrak{w}$  can be presented as

$$\mathfrak{w}(t) = \int_0^t \exp\left((s-t)\mathbf{v}L_\phi^{-1}D_\phi\right) L_\phi^{-1}Z_\phi(s)U\psi(s)ds \quad (22)$$

$$= \left( \int_0^t \exp\left((s-t)\mathbf{v}L_\phi^{-1}D_\phi\right) L_\phi^{-1}Z_\phi(s) (I_n \otimes \psi(s)^T) ds \right) \text{vec}(U), \quad (23)$$

where we have used (17) to make the deformation from (22) to (23). And, from (18), we have

$$\mathfrak{w}(t_i) = \tilde{G}_{\phi, \psi}^{(i)} \text{vec}(U) \in \mathbb{R}^n, \quad \forall i \in \{1, \dots, m\}. \quad (24)$$

Thus, from (23), we obtain the following relation between  $U$  and  $\mathfrak{w}$ :

$$\left( \mathfrak{w}(t_1)^T, \dots, \mathfrak{w}(t_m)^T \right)^T = \tilde{G}_{\phi, \psi} \text{vec}(U).$$

Now, we prove that if (19) is solvable for each  $f_{h,k} \in V_k^1(J; S_h(\Omega))$ , then  $G_{\phi, \psi}$  is nonsingular. For an  $f_{h,k} \in V_k^1(J; S_h(\Omega))$ , we denote the solution of (19) as  $u_{h,k} \in V_k^1(J; S_h(\Omega))$ . From the fact that  $f_{h,k} \in V_k^1 S_h$ , there exists an  $F \in \mathbb{R}^{n \times m}$  such that  $f_{h,k}(x, t) = \phi(x)^T F \psi(t)$ . Note that, for any nodal points  $t_i$ , we have  $(P_{h,k}Au_{h,k})(x, t_i) =$

$(\Pi_k w_h)(x, t_i) = \phi(x)^T \mathfrak{w}(t_i)$  by the definition of  $\Pi_k$ . Therefore, from (19) and (24), we have

$$\begin{aligned} u_{h,k}(x, t_i) - f_{h,k}(x, t_i) &= (P_{h,k} A u_{h,k})(x, t_i), \quad \forall x \in \Omega, \forall i \in \{1, \dots, m\}, \\ &= (\Pi_k w_h)(x, t_i), \end{aligned}$$

which implies

$$\begin{aligned} \phi(x)^T (U - F) \Psi(t_i) &= \phi(x)^T \mathfrak{w}(t_i) \\ &= \phi(x)^T \tilde{G}_{\phi, \Psi}^{(i)} \text{vec}(U). \end{aligned} \quad (25)$$

Since we assume that  $V_k^1(J)$  is the finite element space constituted by the Lagrange elements,  $\Psi_j(t_i) = \delta_{j,i}$  is satisfied, where  $\delta_{j,i}$  denotes the Kronecker delta. Therefore, we get

$$(U - F) \Psi(t_i) = \begin{pmatrix} U_{1,1} - F_{1,1} & \cdots & U_{1,m} - F_{1,m} \\ \vdots & \ddots & \vdots \\ U_{n,1} - F_{n,1} & \cdots & U_{n,m} - F_{n,m} \end{pmatrix} \begin{pmatrix} \Psi_1(t_i) \\ \vdots \\ \Psi_m(t_i) \end{pmatrix} = \begin{pmatrix} U_{1,i} - F_{1,i} \\ \vdots \\ U_{n,i} - F_{n,i} \end{pmatrix}.$$

From the arbitrariness of  $x$  and  $i$ , the variational equation (25) is equivalent to the following simultaneous linear equations:

$$\text{vec}(U - F) = \tilde{G}_{\phi, \Psi} \text{vec}(U).$$

Namely, we have

$$(I_{nm} - \tilde{G}_{\phi, \Psi}) \text{vec}(U) = \text{vec}(F).$$

Therefore, from the arbitrariness of  $f_{h,k}$ , the nonsingularity of  $I_{nm} - \tilde{G}_{\phi, \Psi}$  follows. The converse of this proposition is easily obtained by reversing the discussion.  $\square$

When we apply the proposed a posteriori estimates, it is necessary to confirm that  $G_{\phi, \Psi}$  is nonsingular, which will be able to verify by validated computations such as [14]. Therefore, in what follows, we always assume the nonsingularity of  $G_{\phi, \Psi}$ . Moreover, we define the linear operator  $[I - A]_{h,k}^{-1} : V_k^1(J; S_h(\Omega)) \rightarrow V_k^1(J; S_h(\Omega))$  by the solution of (19). We call this operator a “fully discretized quasi-Newton operator”.

## 5 A posteriori estimates

In this section, we derive a new a posteriori estimate to obtain  $C_{L^2 L^2, L^2 H_0^1}$ , which satisfies (2) by using the fully discretized quasi-Newton operator.

First, we describe a method to calculate the norm of the elements in the full-discretization space. Let  $K_{\phi, \Psi}$  be a matrix in  $\mathbb{R}^{nm \times nm}$  defined by

$$K_{\phi, \Psi} := D_{\phi} \otimes L_{\Psi} = \begin{pmatrix} D_{\phi, 1, 1} L_{\Psi} & \cdots & D_{\phi, 1, n} L_{\Psi} \\ \vdots & \ddots & \vdots \\ D_{\phi, n, 1} L_{\Psi} & \cdots & D_{\phi, n, n} L_{\Psi} \end{pmatrix}. \quad (26)$$

From the symmetric positive definiteness of  $D_\phi$  and  $L_\psi$ , it is readily seen that  $K_{\phi,\psi}$  is also symmetric positive definite. Therefore,  $K_{\phi,\psi}$  is Cholesky decomposable such that  $K_{\phi,\psi} = K_{\phi,\psi}^{1/2} K_{\phi,\psi}^{T/2}$ . Similarly, we define the matrix  $L_{\phi,\psi}$  in  $\mathbb{R}^{nm \times nm}$  as  $L_{\phi,\psi} := L_\phi \otimes L_\psi$ .

**Lemma 5.1** *For an element  $u_{h,k} \in V_k^1(J; S_h(\Omega))$ , taking  $U \in \mathbb{R}^{n \times m}$  such that  $u_{h,k} = \phi^T U \psi$ , then the following equalities hold*

$$\|u_{h,k}\|_{L^2(J; L^2(\Omega))} = \left| L_{\phi,\psi}^{T/2} \text{vec}(U) \right|, \quad (27)$$

$$\|u_{h,k}\|_{L^2(J; H_0^1(\Omega))} = \left| K_{\phi,\psi}^{T/2} \text{vec}(U) \right|, \quad (28)$$

where  $|\cdot|$  denotes the Euclidean norm of a vector.

**Proof.** — Since the proofs of (27) and (28) are almost the same, we will prove only (27). From (17), we have

$$\begin{aligned} \|u_{h,k}\|_{L^2(J; L^2(\Omega))}^2 &= \int_J \int_\Omega \psi(t)^T U^T \phi(x) \phi(x)^T U \psi(t) dx dt \\ &= \int_J \int_\Omega \text{vec}(U)^T (I_n \otimes \psi(t))^T \phi(x) \phi(x)^T (I_n \otimes \psi(t)^T) \text{vec}(U) dx dt \\ &= \text{vec}(U)^T \int_J (I_n \otimes \psi(t))^T L_\phi (I_n \otimes \psi(t)^T) dt \text{vec}(U) \\ &= \text{vec}(U)^T \int_J \begin{pmatrix} L_{\phi,1,1} \psi(t) \psi(t)^T \cdots L_{\phi,1,n} \psi(t) \psi(t)^T \\ \vdots \quad \ddots \quad \vdots \\ L_{\phi,n,1} \psi(t) \psi(t)^T \cdots L_{\phi,n,n} \psi(t) \psi(t)^T \end{pmatrix} dt \text{vec}(U) \\ &= \text{vec}(U)^T \begin{pmatrix} L_{\phi,1,1} L_\psi \cdots L_{\phi,1,n} L_\psi \\ \vdots \quad \ddots \quad \vdots \\ L_{\phi,n,1} L_\psi \cdots L_{\phi,n,n} L_\psi \end{pmatrix} \text{vec}(U) \\ &= \left( L_{\phi,\psi}^{T/2} \text{vec}(U) \right)^T \left( L_{\phi,\psi}^{T/2} \text{vec}(U) \right), \end{aligned}$$

which proves equation (27).  $\square$

Let  $M_{\phi,\psi}(h,k)$  be a nonnegative constant defined by  $M_{\phi,\psi}(h,k) := \left\| K_{\phi,\psi}^{T/2} G_{\phi,\psi}^{-1} K_{\phi,\psi}^{-T/2} \right\|_2$ , where  $\|\cdot\|_2$  denotes the matrix two-norm. The following theorem for  $M_{\phi,\psi}$  holds.

**Theorem 5.2** *It holds that*

$$\left\| [I - A]_{h,k}^{-1} f_{h,k} \right\|_{L^2(J; H_0^1(\Omega))} \leq M_{\phi,\psi} \|f_{h,k}\|_{L^2(J; H_0^1(\Omega))}, \quad \forall f_{h,k} \in V_k^1 S_h. \quad (29)$$

**Proof.** — For any  $f_{h,k} \in V_k^1 S_h$ , we set  $u_{h,k} := [I - A]_{h,k}^{-1} f_{h,k} \in V_k^1 S_h$ . Since  $f_{h,k}$  and  $u_{h,k}$  are the elements of  $V_k^1(J; S_h(\Omega))$ , there exist matrices  $F$  and  $U$  in  $\mathbb{R}^{n \times m}$  such that  $f_{h,k} = \phi^T F \psi$  and  $u_{h,k} = \phi^T U \psi$ , respectively. Moreover, from the proof of Theorem 4.4, it follows that  $\text{vec}(U) = G_{\phi,\psi}^{-1} \text{vec}(F)$ . Therefore, we have the following

estimates

$$\begin{aligned} \|u_{h,k}\|_{L^2(J;H_0^1(\Omega))}^2 &= \text{vec}(U)^T K_{\phi,\psi} \text{vec}(U) \\ &= \left(K_{\phi,\psi}^{T/2} \text{vec}(U)\right)^T \left(K_{\phi,\psi}^{T/2} G_{\phi,\psi}^{-1} K_{\phi,\psi}^{-T/2}\right) \left(K_{\phi,\psi}^{T/2} \text{vec}(F)\right) \\ &\leq \|u_{h,k}\|_{L^2(J;H_0^1(\Omega))} \left\|K_{\phi,\psi}^{T/2} G_{\phi,\psi}^{-1} K_{\phi,\psi}^{-T/2}\right\|_2 \|f_{h,k}\|_{L^2(J;H_0^1(\Omega))}. \end{aligned}$$

This completes the proof.  $\square$

Let  $\mathcal{C}_0$  and  $\mathcal{C}_1$  be the nonnegative constants defined by

$$\mathcal{C}_0 := M_{\phi,\psi} \left( \frac{C_{s,2}}{\nu} + C_{\text{inv}}(h) C_J(k) \right), \quad \mathcal{C}_1 := \|b\|_{L^\infty L^\infty} + C_{s,2} \|c\|_{L^\infty L^\infty},$$

respectively. Moreover, we define the constant  $\kappa_{\phi,\psi}$  as follows:

$$\kappa_{\phi,\psi} := \frac{\|b\|_{L^\infty L^\infty} (1 + \mathcal{C}_0 \mathcal{C}_1) C_1(h,k) + \mathcal{C}_0 \mathcal{C}_1 C_0(h,k) \|c\|_{L^\infty L^\infty}}{1 - C_0(h,k) \|c\|_{L^\infty L^\infty}}, \quad (30)$$

provided that  $1 - C_0(h,k) \|c\|_{L^\infty L^\infty} \neq 0$ .

**Theorem 5.3** *Assume that*

$$0 \leq \kappa_{\phi,\psi} < 1. \quad (31)$$

*Then under the same assumptions as in Theorem 3.2, we have the following constructive a posteriori estimates*

$$\|\mathcal{L}_t^{-1}\|_{\mathcal{L}(L^2(J;L^2(\Omega)), L^2(J;H_0^1(\Omega)))} \leq \frac{1}{1 - \kappa_{\phi,\psi}} \frac{\mathcal{C}_0 + (1 + \mathcal{C}_0 \mathcal{C}_1) C_1(h,k)}{1 - C_0(h,k) \|c\|_{L^\infty L^\infty}}. \quad (32)$$

**Proof.** — For any  $f \in L^2(J;L^2(\Omega))$ , we set  $u := \mathcal{L}_t^{-1} f \in D(\Delta_t)$ . Then we make the following decomposition of (1) into two parts, e.g., the finite- and infinite-dimensional parts, using the projection  $P_{h,k}$ . Namely, in the space  $L^2(J;H_0^1(\Omega))$ , using the following equivalency

$$\begin{aligned} \frac{\partial u}{\partial t} - \nu \Delta u + (b \cdot \nabla) u + cu &= f \\ \iff u &= I_e \Delta_t^{-1} (-(b \cdot \nabla) u - cu + f), \end{aligned} \quad (33)$$

we have the decomposition:

$$\iff \begin{cases} P_{h,k} u = P_{h,k} I_e \Delta_t^{-1} (-(b \cdot \nabla) u - cu + f), & (34a) \\ (I - P_{h,k}) u = (I - P_{h,k}) I_e \Delta_t^{-1} (-(b \cdot \nabla) u - cu + f). & (34b) \end{cases}$$

We set  $u_\perp := u - P_{h,k} u$  for short. From (34a), using the definition of the operator  $A$ , we have

$$P_{h,k} u = P_{h,k} (A(P_{h,k} u + u_\perp) + I_e \Delta_t^{-1} f),$$

by the definition of the operator  $[I - A]_{h,k}^{-1}$ , which implies

$$P_{h,k}u = [I - A]_{h,k}^{-1}P_{h,k}(Au_{\perp} + I_e\Delta_t^{-1}f).$$

Therefore, from (29) and (11), we have the following estimates,

$$\begin{aligned} \|P_{h,k}u\|_{L^2H_0^1} &\leq M_{\phi,\psi} \left( \|P_{h,k}Au_{\perp}\|_{L^2H_0^1} + \|P_{h,k}I_e\Delta_t^{-1}f\|_{L^2H_0^1} \right) \\ &\leq M_{\phi,\psi} \left( \frac{C_{s,2}}{\mathbf{v}} + C_{\text{inv}}(h)C_J(k) \right) (\|(b \cdot \nabla + c)u_{\perp}\|_{L^2L^2} + \|f\|_{L^2L^2}). \end{aligned}$$

From the definition of  $\mathcal{C}_0$ , we have

$$\begin{aligned} \|P_{h,k}u\|_{L^2H_0^1} &\leq \mathcal{C}_0 \|(b \cdot \nabla)u_{\perp} + cu_{\perp}\|_{L^2L^2} + \mathcal{C}_0 \|f\|_{L^2L^2} \\ &\leq \mathcal{C}_0 \|b\|_{L^{\infty}L^{\infty}} \|u_{\perp}\|_{L^2H_0^1} + \mathcal{C}_0 \|c\|_{L^{\infty}L^{\infty}} \|u_{\perp}\|_{L^2L^2} + \mathcal{C}_0 \|f\|_{L^2L^2}. \end{aligned} \quad (35)$$

By calculating the  $L^2L^2$  norm of (34b) using (14), we have

$$\begin{aligned} \|u_{\perp}\|_{L^2L^2} &\leq C_0(h,k) \|-(b \cdot \nabla)u - cu + f\|_{L^2L^2} \\ &\leq C_0(h,k) (\|b\|_{L^{\infty}L^{\infty}} \|u\|_{L^2H_0^1} + \|c\|_{L^{\infty}L^{\infty}} \|u\|_{L^2L^2} + \|f\|_{L^2L^2}), \end{aligned}$$

which yields

$$\begin{aligned} (1 - C_0(h,k) \|c\|_{L^{\infty}L^{\infty}}) \|u_{\perp}\|_{L^2L^2} \\ \leq C_0(h,k) (\|b\|_{L^{\infty}L^{\infty}} \|u\|_{L^2H_0^1} + \|c\|_{L^{\infty}L^{\infty}} \|P_{h,k}u\|_{L^2L^2} + \|f\|_{L^2L^2}). \end{aligned}$$

From (31),  $1 - C_0(h,k) \|c\|_{L^{\infty}L^{\infty}} > 0$  is satisfied. Therefore, we obtain

$$\begin{aligned} \|u_{\perp}\|_{L^2L^2} &\leq \frac{C_0(h,k)}{1 - C_0(h,k) \|c\|_{L^{\infty}L^{\infty}}} \left( \|b\|_{L^{\infty}L^{\infty}} \|P_{h,k}u + u_{\perp}\|_{L^2H_0^1} \right. \\ &\quad \left. + C_{s,2} \|c\|_{L^{\infty}L^{\infty}} \|P_{h,k}u\|_{L^2H_0^1} + \|f\|_{L^2L^2} \right) \\ &\leq \frac{C_0(h,k)}{1 - C_0(h,k) \|c\|_{L^{\infty}L^{\infty}}} (\mathcal{C}_1 \|P_{h,k}u\|_{L^2H_0^1} + \|b\|_{L^{\infty}L^{\infty}} \|u_{\perp}\|_{L^2H_0^1} + \|f\|_{L^2L^2}). \end{aligned} \quad (36)$$

Thus (35) is estimated as

$$\begin{aligned} \|P_{h,k}u\|_{L^2H_0^1} &\leq \mathcal{C}_0 \|b\|_{L^{\infty}L^{\infty}} \|u_{\perp}\|_{L^2H_0^1} + \mathcal{C}_0 \|f\|_{L^2L^2} \\ &\quad + \mathcal{C}_0 \frac{C_0(h,k) \|c\|_{L^{\infty}L^{\infty}}}{1 - C_0(h,k) \|c\|_{L^{\infty}L^{\infty}}} \left( \mathcal{C}_1 \|P_{h,k}u\|_{L^2H_0^1} + \|b\|_{L^{\infty}L^{\infty}} \|u_{\perp}\|_{L^2H_0^1} + \|f\|_{L^2L^2} \right). \end{aligned} \quad (37)$$

Setting nonnegative constants  $R_{1,1}$ ,  $R_{1,2}$ , and  $b_1$  as follows:

$$\begin{aligned} R_{1,1} &:= 1 - \mathcal{C}_0 \mathcal{C}_1 \frac{C_0(h,k) \|c\|_{L^{\infty}L^{\infty}}}{1 - C_0(h,k) \|c\|_{L^{\infty}L^{\infty}}}, & R_{1,2} &:= \frac{\mathcal{C}_0 \|b\|_{L^{\infty}L^{\infty}}}{1 - C_0(h,k) \|c\|_{L^{\infty}L^{\infty}}}, \\ b_1 &:= \frac{\mathcal{C}_0}{1 - C_0(h,k) \|c\|_{L^{\infty}L^{\infty}}}, \end{aligned}$$

(37) is rewritten as

$$R_{1,1} \|P_{h,k}u\|_{L^2(J;H_0^1(\Omega))} - R_{1,2} \|u_\perp\|_{L^2(J;H_0^1(\Omega))} \leq b_1 \|f\|_{L^2(J;L^2(\Omega))}. \quad (38)$$

On the other hand, by considering the  $L^2H_0^1$  norm of (34b), from (13) we have

$$\begin{aligned} \|u_\perp\|_{L^2H_0^1} &\leq C_1(h,k) \|-(b \cdot \nabla)u - cu + f\|_{L^2L^2} \\ &\leq C_1(h,k) \left( \|b\|_{L^\infty L^\infty} \|P_{h,k}u + u_\perp\|_{L^2H_0^1} + \|c\|_{L^\infty L^\infty} \|P_{h,k}u + u_\perp\|_{L^2L^2} + \|f\|_{L^2L^2} \right) \\ &\leq C_1(h,k) \left( \mathcal{C}_1 \|P_{h,k}u\|_{L^2H_0^1} + \|b\|_{L^\infty L^\infty} \|u_\perp\|_{L^2H_0^1} + \|c\|_{L^\infty L^\infty} \|u_\perp\|_{L^2L^2} + \|f\|_{L^2L^2} \right). \end{aligned}$$

From (36), we obtain

$$\begin{aligned} \|u_\perp\|_{L^2H_0^1} &\leq C_1(h,k) \mathcal{C}_1 \|P_{h,k}u\|_{L^2H_0^1} + C_1(h,k) \|b\|_{L^\infty L^\infty} \|u_\perp\|_{L^2H_0^1} + C_1(h,k) \|f\|_{L^2L^2} \\ &\quad + C_1(h,k) \frac{C_0(h,k) \|c\|_{L^\infty L^\infty}}{1 - C_0(h,k) \|c\|_{L^\infty L^\infty}} \left( \mathcal{C}_1 \|P_{h,k}u\|_{L^2H_0^1} + \|b\|_{L^\infty L^\infty} \|u_\perp\|_{L^2H_0^1} + \|f\|_{L^2L^2} \right). \end{aligned} \quad (39)$$

We set nonnegative constants  $R_{2,1}$ ,  $R_{2,2}$ , and  $b_2$  as follows:

$$\begin{aligned} R_{2,1} &:= \frac{\mathcal{C}_1 C_1(h,k)}{1 - C_0(h,k) \|c\|_{L^\infty L^\infty}}, & R_{2,2} &:= 1 - \frac{\|b\|_{L^\infty L^\infty} C_1(h,k)}{1 - C_0(h,k) \|c\|_{L^\infty L^\infty}}, \\ b_2 &:= \frac{C_1(h,k)}{1 - C_0(h,k) \|c\|_{L^\infty L^\infty}}, \end{aligned}$$

where we note that the positivity of  $R_{2,2}$  follows by the condition (31). Thus (39) can be rewritten as

$$-R_{2,1} \|P_{h,k}u\|_{L^2(J;H_0^1(\Omega))} + R_{2,2} \|u_\perp\|_{L^2(J;H_0^1(\Omega))} \leq b_2 \|f\|_{L^2(J;L^2(\Omega))}. \quad (40)$$

From (38) and (40), we have the following simultaneous inequalities,

$$\begin{pmatrix} R_{1,1} & -R_{1,2} \\ -R_{2,1} & R_{2,2} \end{pmatrix} \begin{pmatrix} \|P_{h,k}u\|_{L^2(J;H_0^1(\Omega))} \\ \|u_\perp\|_{L^2(J;H_0^1(\Omega))} \end{pmatrix} \leq \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \|f\|_{L^2(J;L^2(\Omega))}.$$

By assumption (31), we obtain

$$\det \begin{pmatrix} R_{1,1} & -R_{1,2} \\ -R_{2,1} & R_{2,2} \end{pmatrix} = 1 - \kappa_{\phi,\psi} > 0.$$

Therefore, the simultaneous inequalities can be solved as follows:

$$\begin{pmatrix} \|P_{h,k}u\|_{L^2(J;H_0^1(\Omega))} \\ \|u_\perp\|_{L^2(J;H_0^1(\Omega))} \end{pmatrix} \leq \frac{1}{1 - \kappa_{\phi,\psi}} \begin{pmatrix} R_{2,2} & R_{1,2} \\ R_{2,1} & R_{1,1} \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \|f\|_{L^2(J;L^2(\Omega))}. \quad (41)$$

Finally, from (41), we have

$$\begin{aligned} \|u\|_{L^2(J;H_0^1(\Omega))} &\leq \|P_{h,k}u\|_{L^2(J;H_0^1(\Omega))} + \|u_\perp\|_{L^2(J;H_0^1(\Omega))} \\ &\leq \frac{R_{2,2}b_1 + R_{1,2}b_2 + R_{2,1}b_1 + R_{1,1}b_2}{1 - \kappa_{\phi,\psi}} \|f\|_{L^2(J;L^2(\Omega))}, \end{aligned}$$

which proves the desired estimates.  $\square$

## 6 Numerical examples

In this section, we show several rigorous numerical results for  $C_{L^2L^2,L^2H_0^1}$  satisfying (2) for test problems by three kinds of methods, namely, a priori estimates (the Gronwall inequality), a posteriori estimates proposed in [10], and the new method in Theorem 5.3. Moreover, we also show several rigorous error bounds of the numerical solutions for the nonlinear parabolic equations as an application of the estimates of (2).

We considered the norm estimates for an inverse operator of the following  $\mathcal{L}_t$ :

$$\mathcal{L}_t := \frac{\partial}{\partial t} - \mathbf{v} \Delta - 2u_h^k, \quad (42)$$

that is,  $b = 0$  and  $c = -2u_h^k$  in (2). Here,  $u_h^k$  is assumed to be an approximate solution of the following nonlinear parabolic problem:

$$\begin{cases} \frac{\partial u}{\partial t} - \mathbf{v} \Delta u = u^2 + f, & \text{in } \Omega \times J, & (43a) \\ u(x, t) = 0, & \text{on } \partial\Omega \times J, & (43b) \\ u(x, 0) = 0, & \text{in } \Omega. & (43c) \end{cases}$$

Therefore, (42) becomes a linearized operator of (43) at  $u_h^k$ . We only considered one-space-dimensional case ( $d = 1$ ) with  $\Omega = (0, 1)$ . Furthermore, the function  $f$  was chosen so that the problem (43) had the following exact solutions:

- $u(x, t) = 0.5t \sin(\pi x)$ ,  $\mathbf{v} = 0.1$ , (Example 1.1);
- $u(x, t) = 0.5t \sin(\pi x)$ ,  $\mathbf{v} = 1.0$ , (Example 1.2);
- $u(x, t) = \sin(\pi t) \sin(\pi x)$ ,  $\mathbf{v} = 0.1$ , (Example 2.1);
- $u(x, t) = \sin(\pi t) \sin(\pi x)$ ,  $\mathbf{v} = 1.0$ , (Example 2.2).

Note that Example 1.1 and Example 2.1 are studied in [10]. In each example, the function  $u_h^k$  was computed as an approximation of the corresponding  $u$  by using a piecewise-cubic Hermite interpolation in the space direction with a piecewise-linear interpolation in the time direction. Therefore,  $u_h^k$  belongs to  $V^1(J; H_0^1(\Omega) \cap H^2(\Omega))$ .

We used the finite-dimensional spaces  $S_h(\Omega)$  and  $V_k^1(J)$ , spanned by piecewise linear functions with uniform mesh size  $h$  and  $k$ , respectively, so that they satisfied  $k = h^2$ . Then, it was seen that the constants in previous sections could be taken as  $C_\Omega(h) = h/\pi$ ,  $C_{\text{inv}}(h) = \sqrt{12}/h$ ,  $C_J(k) = k/\pi = h^2/\pi$ , and  $C_{s,2} = 1/\pi$ , respectively. Moreover, we had

$$\|c\|_{L^\infty(J; L^\infty(\Omega))} = 2 \left\| u_h^k \right\|_{L^\infty(J; L^\infty(\Omega))} \leq \begin{cases} T & \text{(Example 1.1 and 1.2)} \\ 2 & \text{(Example 2.1 and 2.2)}. \end{cases}$$

### 6.1 A posteriori estimates of the inverse parabolic operator

We now present the results computed for  $C_{L^2L^2,L^2H_0^1}$  by using three kinds of method. Here, we used the following a priori estimate, which comes from the Gronwall in-

equality

$$\|\mathcal{L}_T^{-1}\|_{\mathcal{L}(L^2(J;L^2(\Omega)),L^2(J;H_0^1(\Omega)))} \leq \exp(\gamma T) \frac{C_{s,2}}{\nu}, \quad \gamma := \max \left\{ \sup_{\Omega \times J} (-c), 0 \right\}.$$

In this subsection, we will refer to the a posteriori estimates studied in [10] and our present estimates (32) as “a posteriori estimate I” and “a posteriori estimate II,” respectively. To compute a posteriori estimate I, we used the same parameters as in [10], i.e.,  $(n, m) = (5, 700 \cdot T^2)$  for Example 1.1,  $(n, m) = (5, 100 \cdot 4^T)$  for Example 2.1, where we note that  $h = 1/(n+1)$  and  $k = 1/m$ . For a posteriori estimate II, we used  $h = 1/8$  and  $h = 1/16$ , with  $k = h^2$ .

*Example 1.1 and 1.2:*  $u_h^k(x, t) \approx -t \sin(\pi x)$

Figures 1-2 show the values of  $C_{L^2L^2, L^2H_0^1}$  for Example 1.1-1.2, plotted out on log-linear coordinates. For  $T > 1$ , the values of the proposed estimates are smaller than

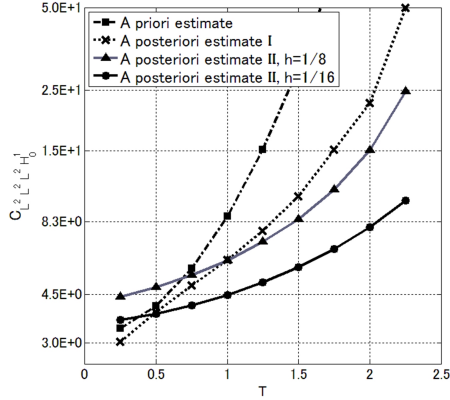


Fig. 1  $\nu = 0.1$

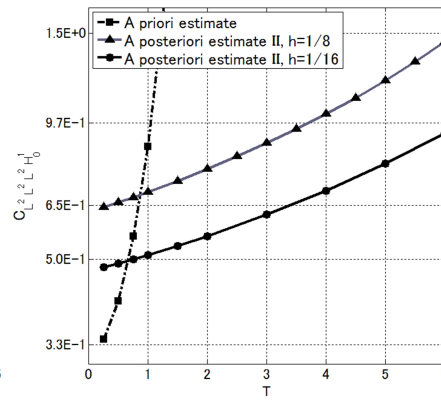


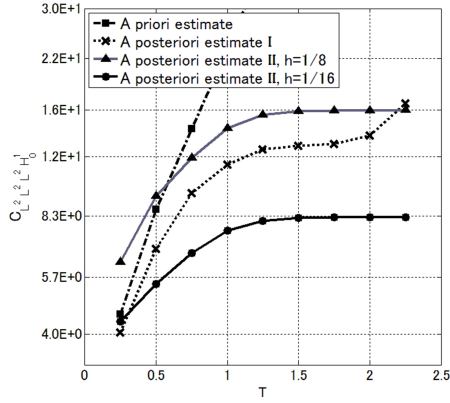
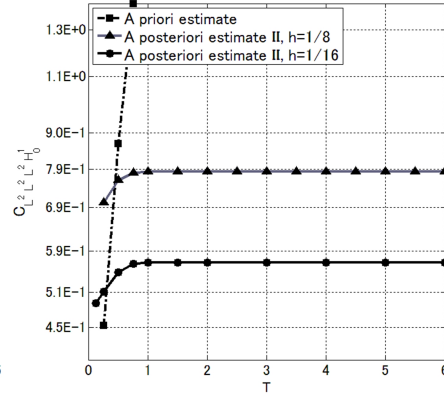
Fig. 2  $\nu = 1$

the other estimates. The two kinds of a posteriori estimates require the validated upper bound for the matrix two-norm of the corresponding unsymmetric dense matrices (e.g.,  $M_{\phi, \psi}$ ), and most of the computational costs is due to this task. In Example 1.1, for  $T = 2$ , a posteriori estimate I a matrix of size 14000, but in a posteriori estimate II, we can attain our purpose with a matrix of size 896 for  $h = 1/8$ , and 7680 for  $h = 1/16$ . This fact shows, in the case of a posteriori estimate II, that it is not necessary to take special account of the stiff property of the ODEs coming from the semidiscretization.

*Example 2.1 and 2.2:*  $u_h^k(x, t) \approx -2 \sin(\pi t) \sin(\pi x)$

Figures 3-4 show the values of  $C_{L^2L^2, L^2H_0^1}$  for Example 2.1-2.2 (log-linear coordinates). For  $T > 1/2$ , the values of the proposed estimates with  $h = 1/16$  are smaller than the other estimates. In Example 2.1, for  $T = 2$ , a posteriori estimate I requires



Fig. 3  $\nu = 0.1$ Fig. 4  $\nu = 1$ 

a matrix of size 8000, but a posteriori estimate II requires only one of size 896 for  $h = 1/8$  and size 7680 for  $h = 1/16$ . It is notable that the results of the proposed estimates show no exponential dependency for  $T$ . On the other hand, due to the stiffness of the corresponding ODEs, we were not successful in computing the inverse operator of a posteriori estimate I, except for the case where  $T$  was very small.

## 6.2 Verification results for solutions of nonlinear parabolic equations

Applying the estimates (2), we implemented a numerical verification method to prove the existence of solutions for the nonlinear parabolic problems. As a prototype application, we considered the nonlinear parabolic initial-boundary-value problems of the form (43). In a similar way as in [8] for the elliptic case, we defined the fixed-point equation for a compact operator, which is equivalent to (43) with the Newton-type residual form, and derived a verification condition by applying the Schauder fixed-point theorem.

First, we considered the following residual equation for (43):

$$\begin{cases} \frac{\partial w}{\partial t} - \nu \Delta w - 2u_h^k w = g(w), & \text{in } \Omega \times J, & (44a) \\ w(x, t) = 0, & \text{on } \partial\Omega \times J, & (44b) \\ w(x, 0) = 0, & \text{in } \Omega, & (44c) \end{cases}$$

where

$$g(w) = w^2 + \varepsilon, \quad \varepsilon = (u_h^k)^2 + f - \left( \frac{\partial u_h^k}{\partial t} - \nu \Delta u_h^k \right).$$

Note that if the approximate solution  $u_h^k$  is close to the exact solution of (43), then  $w \approx 0$ ,  $\varepsilon \approx 0$ , and  $g(w) \approx 0$ . Thus (44) can be rewritten as the following fixed-point equation of the compact map  $F$ :

$$w = \mathcal{L}_t^{-1} g(w) =: F(w). \quad (45)$$

Next, for any positive constants  $\alpha$  and  $\beta$ , we define the candidate set  $W_{\alpha,\beta}$  as

$$W_{\alpha,\beta} := \left\{ w \in V ; \|w\|_{L^2H_0^1} \leq \alpha, \|w\|_{V^1L^2} \leq \beta \right\}.$$

From the Schauder fixed-point theorem, noting that the continuity of the map  $F$  in the space  $L^2H_0^1$ , if the set  $W_{\alpha,\beta}$  satisfies

$$F(W_{\alpha,\beta}) \subset W_{\alpha,\beta}, \quad (46)$$

then a fixed point of (45) exists in the set  $\overline{W_{\alpha,\beta}}$ , where  $\overline{W_{\alpha,\beta}}$  stands for the closure of the set  $W_{\alpha,\beta}$  in  $L^2H_0^1$ .

Now, by some simple calculations using the Sobolev embedding theorem and the Poincaré inequality, it is easily seen that the following inequalities hold for any  $w \in W_{\alpha,\beta}$ :

$$\begin{aligned} \|F(w)\|_{L^2H_0^1} &\leq C_{L^2L^2,L^2H_0^1} \left( \alpha\beta\sqrt{\frac{T}{8}} + \|\varepsilon\|_{L^2L^2} \right), \\ \|F(w)\|_{V^1L^2} &\leq \left( \frac{2}{\pi} C_{L^2L^2,L^2H_0^1} \|u_h^k\|_{L^\infty L^\infty} + 1 \right) \left( \alpha\beta\sqrt{\frac{T}{8}} + \|\varepsilon\|_{L^2L^2} \right). \end{aligned}$$

From these inequalities, we have the following sufficient condition for (46):

$$\begin{cases} C_{L^2L^2,L^2H_0^1} \left( \alpha\beta\sqrt{\frac{T}{8}} + \|\varepsilon\|_{L^2L^2} \right) \leq \alpha, \\ \left( \frac{2}{\pi} C_{L^2L^2,L^2H_0^1} \|u_h^k\|_{L^\infty L^\infty} + 1 \right) \left( \alpha\beta\sqrt{\frac{T}{8}} + \|\varepsilon\|_{L^2L^2} \right) \leq \beta. \end{cases}$$

Thus, we obtain the verification condition for the existence of the solutions of (44). Since it holds that  $w = u - u_h^k$ , by solving the above simultaneous algebraic inequalities in  $\alpha$  and  $\beta$ , we have error bounds of the form  $\|u - u_h^k\|_{L^2H_0^1} \leq \alpha$ .

We now present the verification results for the solutions of (44), namely,  $\alpha$ ,  $\beta$ , and the residual norm,  $\|\varepsilon\|_{L^2L^2}$ . In Figure 5, we chose the function  $f$  so that (43) has the exact solution  $u(x,t) = 0.5t \sin(\pi x)$ , with  $\nu = 0.1$  and  $\nu = 1.0$ , which correspond to Example 1.1 and Example 1.2, respectively, in the previous subsection. We show more results in Figure 6, in which the function  $f$  is chosen so that (43) has the exact solution  $u(x,t) = \sin(\pi t) \sin(\pi x)$ , with  $\nu = 0.1$  and  $\nu = 1.0$ , which correspond to Example 2.1 and Example 2.2, respectively.

From these figures, it is seen that the error bounds increase in proportion to the residual norms. This property should be expected in our verification conditions. Namely, the validated accuracy of the present method is essentially dependent on the residual norm of the approximate solutions.

We see in the left side of Figure 6, for the case  $\nu = 0.1$ , that our verification method failed for  $T > 0.5$  with  $h = 1/8$ . On the other hand, since  $C_\Omega(h)$  and the residual norm  $\|\varepsilon\|_{L^2L^2}$  for the mesh size  $h = 1/16$  are smaller than in case of  $h = 1/8$ , we succeeded in verification up to  $T \leq 1.25$ . This fact shows that a smaller  $h$  yields better verification, which should also be quite expected.

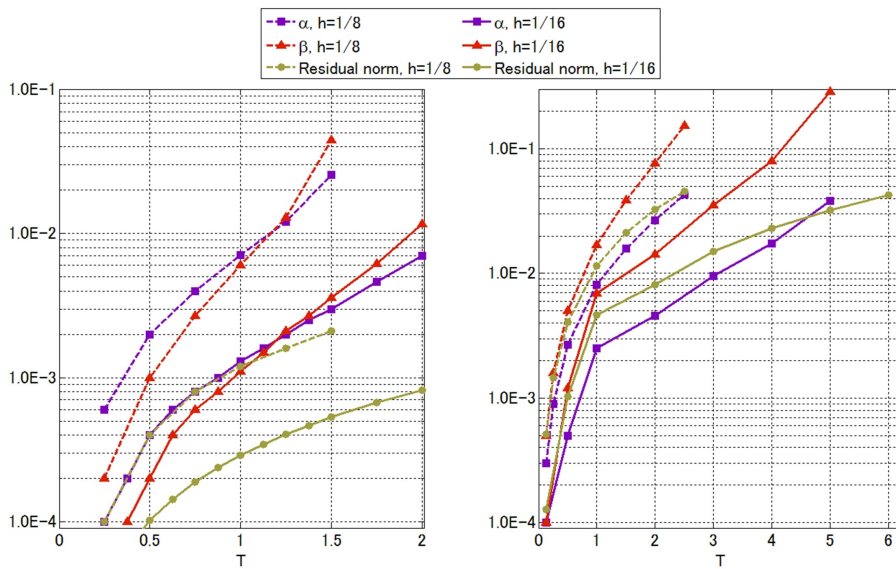


Fig. 5 Verification results: Exact solution for  $u(x,t) = 0.5t \sin(\pi x)$ , with  $v = 0.1$ (left),  $v = 1.0$ (right)

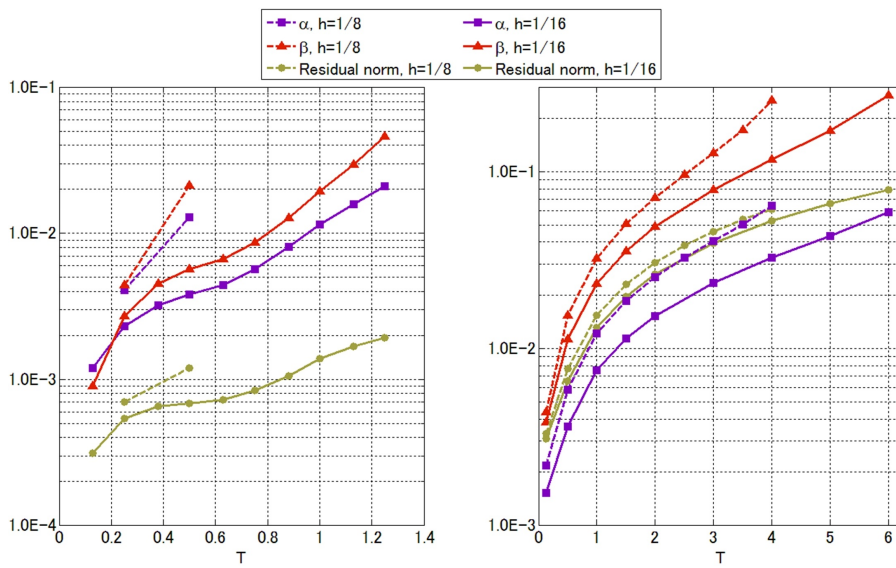


Fig. 6 Verification results: Exact solution for  $u(x,t) = \sin(\pi t) \sin(\pi x)$ , with  $v = 0.1$ (left),  $v = 1.0$ (right)

**Remark 6.1 (Computer environment)** *All computations were carried out on a Dell Precision T7500 (Intel Xeon x5680, 72 GB of memory) with MATLAB R2010b. The computation errors have been taken into account by using INTLAB 6.0, a toolbox for self-validating algorithms, developed by Rump [14].*

## 7 Conclusions

We propose a method to compute constructive a posteriori estimates of the inverse operators for parabolic initial-boundary-value problems. This method is based on the full-discretization quasi-Newton operator, as well as the constructive a priori error estimates for the Galerkin method, with an interpolation in time by effectively using the fundamental solution for the spatial semidiscretization. Our proposed new estimates (32) seem to be better and more robust than the previous estimates [10], as illustrated in the test problems. Moreover, by applying the method to some prototype examples, we illustrated that our method can be used to enclose solutions for nonlinear parabolic problems.

**Acknowledgements** This work was supported by Grants-in-Aid from the Ministry of Education, Culture, Sports, Science and Technology of Japan (No. 20224001 and No. 23740074) and supported by the Global Center of Excellence (GCOE) program, “Fostering top leaders in mathematics”, Kyoto University.

## References

1. Kimura, S., Yamamoto, N.: On explicit bounds in the error for the  $H_0^1$ -projection into piecewise polynomial spaces. *Bull. Inform. Cybernet.* 31 no. 2, 109–115 (1999)
2. Kinoshita, T., Kimura, T., Nakao, M. T.: A posteriori estimates of inverse operators for initial value problems in linear ordinary differential equations. *J. Comput. Appl. Math.* 236 no. 6, 1622–1636 (2011)
3. Ladyženskaja, O. A., Solonnikov, V. A., Ural'ceva, N. N.: *Linear and Quasi-linear Equations of Parabolic Type*. AMS, Rhode Island (1968)
4. Luskin, M., Rannacher, R.: On the smoothing property of the Galerkin method for parabolic equations. *SIAM J. Numer. Anal.* 19 no. 1, 93–113 (1982)
5. Minamoto, T.: Numerical Existence and Uniqueness Proof for Solutions of Semilinear Parabolic Equations. *Appl. Math. Lett.* 14 no. 6, 707–714 (2001)
6. Nakao, M. T.: Solving nonlinear parabolic problems with result verification. Part I: One-space-dimensional case. *J. Comput. Appl. Math.* 38, 323–334 (1991)
7. Nakao, M. T., Yamamoto, N., Kimura, S.: On the best constant in the error bound for the  $H_0^1$ -projection into piecewise polynomial spaces. *J. Approx. Theory* 93, 491–500 (1998)
8. Nakao, M. T., Hashimoto, K., Watanabe, Y.: A numerical method to verify the invertibility of linear elliptic operators with applications to nonlinear problems. *Computing* 75, 1–14 (2005)
9. Nakao, M. T., Hashimoto, K.: A numerical verification method for solutions of nonlinear parabolic problems. *Journal of Math-for-Industry* 1, 69–72 (2009)
10. Nakao, M. T., Kinoshita, T., Kimura, T.: On a posteriori estimates of inverse operators for linear parabolic initial-boundary value problems. *Computing* 94 no. 2, 151–162 (2012)
11. Nakao, M. T., Kimura, T., Kinoshita, T.: Constructive a priori error estimates for a full discrete approximation of the heat equation. *RIMS Preprints RIMS-1745* (2012)
12. Plum, M.: Computer-Assisted Existence Proofs for Two-Point Boundary Value Problems. *Computing* 46 no. 1, 19–34 (1991)
13. Plum, M.: Computer-Assisted Proofs for Semilinear Elliptic Boundary Value Problems. *Japan J. Indust. Appl. Math.* 26 no. 2-3, 419–442 (2009)

14. Rump, S. M.: INTLAB–INTERVAL LABORATORY. In: Csendes, T., (ed.) *Developments in Reliable Computing*, pp. 77–104. Kluwer, Dordrecht (1999) <http://www.ti3.tu-harburg.de/rump/intlab/>
15. Schultz, M. H.: *Spline Analysis*. Prentice-Hall, Englewood Cliffs, N.J. (1973)
16. Zeidler, E.: *Nonlinear Functional Analysis and its Applications II/A*. Springer-Verlag, New York (1990)