



# Stability Analysis and Error Estimate of the Explicit Single-Step Time-Marching Discontinuous Galerkin Methods with Stage-Dependent Numerical Flux Parameters for a Linear Hyperbolic Equation in One Dimension

Yuan Xu<sup>1</sup> · Chi-Wang Shu<sup>2</sup> · Qiang Zhang<sup>3</sup> 

Received: 14 March 2024 / Revised: 27 June 2024 / Accepted: 5 July 2024 /

Published online: 13 July 2024

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2024

## Abstract

In this paper, we present the  $L^2$ -norm stability analysis and error estimate for the explicit single-step time-marching discontinuous Galerkin (DG) methods with stage-dependent numerical flux parameters, when solving a linear constant-coefficient hyperbolic equation in one dimension. Two well-known examples of this method include the Runge–Kutta DG method with the downwind treatment for the negative time marching coefficients, as well as the Lax–Wendroff DG method with arbitrary numerical flux parameters to deal with the auxiliary variables. The stability analysis framework is an extension and an application of the matrix transferring process based on the temporal differences of stage solutions, and a new concept, named as the averaged numerical flux parameter, is proposed to reveal the essential upwind mechanism in the fully discrete status. Distinguished from the traditional analysis, we have to present a novel way to obtain the optimal error estimate in both space and time. The main tool is a series of space–time approximation functions for a given spatial function, which preserve the local structure of the fully discrete schemes and the balance of exact evolution under the control of the partial differential equation. Finally some numerical experiments are given to validate the theoretical results proposed in this paper.

**Keywords** Discontinuous Galerkin method · Explicit single step time marching · Stage-dependent numerical flux parameters · Hyperbolic equation · Stability analysis and error estimate

**Mathematics Subject Classification** 65M12 · 65M15

---

✉ Qiang Zhang  
qzh@nju.edu.cn

Yuan Xu  
yuanxu@nju.edu.cn

Chi-Wang Shu  
chi-wang\_shu@brown.edu

<sup>1</sup> School of Mathematical Sciences, Nanjing Normal University, Nanjing 210023, Jiangsu, China

<sup>2</sup> Division of Applied Mathematics, Brown University, Providence, RI 02912, USA

<sup>3</sup> Department of Mathematics, Nanjing University, Nanjing 210093, Jiangsu, China

## 1 Introduction

In this paper we would like to present the  $L^2$ -norm stability analysis and error estimate for the explicit single-step time-marching discontinuous Galerkin (ESTDG) methods in a more general application of numerical fluxes. Two well-known examples include the RKDG method and the LWDG method to solve hyperbolic equations, which respectively employ the Runge–Kutta time marching [5–9], and the Lax–Wendroff time marching [13, 23] to solve the semidiscrete discontinuous Galerkin (DG) method. Many applications have shown that these methods are good at solving nonlinear conservation laws, due to good stability, high order accuracy and the ability for capturing shocks sharply. For more details, we refer to the review papers [10, 15, 20, 21] and the references therein.

Besides the time marching algorithms, the major concepts in these methods are the numerical fluxes used in the DG spatial discretization. We remark that, in numerical applications, nonlinear limiters are also used to improve the numerical performance when shocks appear. However, in this paper we do not consider the limiters and only pay attention to the interaction between the numerical fluxes and the time discretization. In most numerical experiments, numerical fluxes are often taken as the same type or with the same parameter at every element boundaries and time stages. However, the numerical fluxes are allowed to be changed and this strategy has been actually applied in many numerical simulations. A famous example is the downwind treatment in high order RKDG methods to deal with the negative time-marching coefficients [7, 10], in order to ensure the total variation diminishing in the means (TVDM) property (coupled with a suitable limiter) under the strong-stability-preserving (SSP) framework [11] such that a good numerical performance might be obtained nearby the shock. This downwind treatment is necessary because the Runge–Kutta algorithm for nonlinear problems must have negative time-marching coefficients to achieve fifth or higher orders of time accuracy, as well as in the fourth order with only four stages [12, 17]. We would like to mention that the downwind treatment is also used in many high order numerical methods (for instance, the TVD, ENO and WENO schemes) with the Runge–Kutta algorithms [12, 16, 18, 22] and the multistep algorithms [19]. Another example is the LWDG method, where the DG discretization for those high order Lax–Wendroff expansion is often different to that for the first order (convection) term; see the second example in Sect. 2.2.

To accurately understand the numerical effects of the above treatments, we need to carry out the corresponding theoretical analysis for the ESTDG method with stage-dependent numerical flux functions. However, as far as the authors know, till now there is not any discussion on this topic, even for a simple model equation. To fill in this gap, we would like in this paper to carry out the  $L^2$ -norm stability analysis and establish optimal error estimates of the ESTDG method in a unified framework, for the linear constant-coefficient hyperbolic equation in one dimension

$$\partial_t U + \beta \partial_x U = 0, \quad x \in I = (0, 1), \quad t > 0, \quad (1.1)$$

equipped with the initial condition  $U(x, 0) = U_0(x)$  and the periodic boundary condition. We think that the deep research on this topic provides a starting point to push ahead theoretical studies on the fully discrete DG method that is really used in the practical simulation of nonlinear conservation laws. For simplicity, we further assume that  $\beta$  is a positive constant, and that the numerical flux parameter involved in the numerical flux only changes at different stage time and does not change with respect to the space position; see Sect. 2.2. Different from the special case that numerical flux parameters are the same in the RKDG methods, we have to spend extra effort and propose a new strategy to carefully handle the analysis

difficulties resulted from the perturbation of the numerical flux parameters in the ESTDG methods.

There are two major difficulties to carry out the  $L^2$ -norm stability analysis. On one hand, it is well known [2] that the DG method coupled with the forward Euler time-marching is unstable for any fixed CFL number if the polynomial space is not piecewise constant. That is to say, the  $L^2$ -norm stability of ESTDG methods can not be derived under the so-called SSP framework. We have to set up a facilitating energy equation to carry out energy analysis. However, this is difficult for the high order in time fully discrete DG methods. Recently this trouble is systematically settled by the technique of matrix transferring process based on temporal differences of stage solutions, which can automatically achieve the expected energy equation step by step. This technique has been successfully applied for the RKDG methods when numerical flux parameters are the same; see the references [1, 26–30]. Similar works on this issue can be found in the framework to analyze the stability of the explicit RK methods to solve an ODEs with semi-negative linear operator [24]. On the other hand, in this paper we have to overcome the new difficulty resulting from the stage-dependent numerical flux parameters. As a main highlight of this paper, we make an application and/or an extension of the matrix transferring process and put forward an important quantity, named as the *averaged numerical flux parameter*. This quantity reveals the overall upwind effect in every step time-marching, so it should be greater than one half from the viewpoint of practice. Further, by deep discussions on two detailed examples we find out that, via adjusting the numerical flux parameters (even though the averaged numerical flux are not enlarged), we have a chance to improve the stability performance of the ESTDG method, for example, from the strong stability to the monotonicity stability. For more detailed concepts and statements, see Sect. 3.

Unfortunately, for the ESTDG method with stage-dependent numerical flux parameters, the optimal error estimate becomes very difficult, although the suboptimal error estimate is trivial by traditional treatments. When numerical flux parameters are the same, this purpose has been achieved for the RKDG methods [27, 31, 32] by virtue of the above stability analysis and the generalized Gauss–Radau (GGR) projection with a fixed parameter. However, this proof strategy does not work well for the general case that numerical flux parameters are changed at different occurrence. The main reason is that the element boundary errors at different stages can not be simultaneously eliminated by a fixed GGR projection. To overcome this difficulty, we propose in this paper a new analysis tool, named as *a series of space–time approximation functions* for any given spatial function. They preserve not only the local structure of the fully discrete scheme, but also the local balance of exact evolution under the rule of the considered partial differential equation (PDE). Hence, they are able to provide a group of good reference functions belonging to the finite element space, such that the error accumulation in time of the fully discrete scheme is elaborately scattered over every gap, at the time level  $t^n$ , between the head function (the first one in the series) of  $U(x, t^n)$  and the tail function (the last one in the series) of  $U(x, t^n - \tau)$ . Here  $U(x, t)$  is the exact solution and  $\tau$  is the time step. With the help of the results and the stability conclusions for the nonhomogeneous problem (as a trivial extension of those in Sect. 3.2), the difficulty to obtaining the optimal error estimate is shifted to how to prove the optimal estimate for a series of space–time approximation functions. From our point of view, this analysis line is specifically designed for the fully discrete scheme and is remarkably distinguished to the traditional analysis line, which is often pushed ahead from a semi-discrete scheme in either time or space to the fully discrete scheme.

Because each one in a series of space–time approximation functions cannot be regarded as a traditional projection of the given function, we encounter serious difficulties in proving the optimal approximation property; see Lemma 4.1. Fortunately, this aim can be accomplished

by the aid of those techniques and concepts proposed in the matrix transferring process, for instance, the temporal differences of stage solutions and the evolution identity. Here we would like to mention that the averaged numerical flux parameter still plays an important role in this analysis process. With this special quantity, the  $L^2$ -norms of the specially-defined error function sequences (see (4.25) for details) can be mainly bounded by each other in the forward and reverse directions, respectively; see Lemmas 4.2 and 4.3. In this entire process, the GGR projection and the flux lifting function (see Sect. 4.2) are fully utilized.

It is worthy to emphasize that the averaged numerical flux parameter makes significant contributions throughout the theoretical analysis of this paper. To prove Lemma 3.7, we have to make a deep investigation on the matrix transferring process and make more efforts to establish the subtle relationship among the one-step time marching and the multistep one. This procedure involves many manipulations of matrices, such as the matrix description of matrix transferring process and the Kronecker products of matrices. By tedious and rigorous calculations, we discover the important role of the hidden zero restriction related to the averaged numerical flux parameter, which is stated in Lemma 3.5 with  $m = 1$  or the equivalent identity (7.28). This zero restriction helps us to prove that the concerned submatrix in the multistep spatial matrix is close to a symmetric positive definite (SPD) matrix congruent to the Hilbert matrix such that the distance is reciprocal to the multistep number; see the appendix. Another application of this zero restriction is the proof of Lemma 4.3, where the coefficient in front of the jump term of the head function is successfully eliminated; see Sect. 4.2.2.

The rest of paper is organized as follows. In Sect. 2 we describe the ESTDG method with stage-dependent numerical flux parameters and then present two well-known examples that will be analyzed and numerically tested in this paper. In Sect. 3 we present a framework to derive energy equation and carry out the  $L^2$ -norm stability analysis, where the averaged numerical flux parameter is proposed. Section 4 is devoted to obtaining the optimal error estimate in  $L^2$ -norm, where a series of space-time approximation functions are proposed and analyzed. Some numerical experiments are given in Sect. 5 to verify the theoretical results. The concluding remarks and some technical proofs are respectively presented in Sect. 6 and the appendix.

## 2 The ESTDG Method

In this section we present the detailed definition of the ESTDG methods to solve (1.1) and then show two well-known examples including the RKDG method and the LWDG method.

### 2.1 The Semidiscrete DG Method

Let  $J$  be any positive integer and  $0 = x_{1/2} < x_{3/2} < \dots < x_{J+1/2} = 1$  be a quasi-uniform partition of the spatial interval  $I$ . Each element  $I_j = (x_{j-1/2}, x_{j+1/2})$  has the length  $h_j = x_{j+1/2} - x_{j-1/2}$  for  $j = 1, 2, \dots, J$ . Denote  $h = \max_{1 \leq j \leq J} h_j$ . Then we define the discontinuous finite element space by

$$V_h = \{v \in L^2(I) : v|_{I_j} \in \mathcal{P}^k(I_j), j = 1, 2, \dots, J\}, \quad (2.1)$$

where  $\mathcal{P}^k(I_j)$  is the polynomial space in  $I_j$  of degree at most  $k \geq 0$ . As usual we denote by  $v^+$  and  $v^-$  the limits of  $v$  from two sides.

In this paper,  $I_h$  denotes the partition and  $\Gamma_h$  the element boundaries. The inner product in  $L^2(I_h)$  and  $L^2(\Gamma_h)$  are respectively denoted by  $(\cdot, \cdot)_{I_h}$  and  $\langle \cdot, \cdot \rangle_{\Gamma_h}$ . The associated norms are  $\|\cdot\|_{L^2(I)} = \|\cdot\|_{L^2(I_h)}$  and  $\|\cdot\|_{L^2(\Gamma_h)}$ , respectively. For any  $v \in V_h$  there hold the inverse inequalities [4, 15]:

$$\|\partial_x v\|_{L^2(I)} \leq \mu h^{-1} \|v\|_{L^2(I)}, \quad \|v^\pm\|_{L^2(\Gamma_h)} \leq \mu h^{-\frac{1}{2}} \|v\|_{L^2(I)}, \tag{2.2}$$

where  $\mu > 0$  is the inverse constant independent of  $v$  and  $h$ .

The semidiscrete DG method for the model Eq. (1.1) is defined as follows: find a map  $u(t) : [0, T] \rightarrow V_h$  such that it satisfies

$$\left(\partial_t u, v\right)_{I_h} = \mathcal{H}^\theta(u, v), \quad \forall v \in V_h, \quad t \in (0, T], \tag{2.3}$$

with a well-defined initial solution  $u(0) \in V_h$ . Here  $\mathcal{H}^\theta(u, v)$  is the so-called spatial DG discretization, defined in the form

$$\mathcal{H}^\theta(u, v) = \underbrace{\sum_{1 \leq j \leq J} \int_{I_j} \beta u \partial_x v dx}_{(\beta u, \partial_x v)_{I_h}} + \underbrace{\sum_{1 \leq j \leq J} \beta \{u\}_{j+\frac{1}{2}}^\theta \llbracket v \rrbracket_{j+\frac{1}{2}}}_{(\beta \{u\}^\theta, \llbracket v \rrbracket)_{\Gamma_h}}, \tag{2.4}$$

with the weighted average and the jump at element boundary

$$\{u\}^\theta = \theta u^- + (1 - \theta)u^+, \quad \llbracket v \rrbracket = v^+ - v^-.$$

In this paper,  $\theta$  is called the numerical flux parameter. It is often assumed to be independent of time and greater than 1/2 in order to provide the upwind mechanism and the  $L^2$ -norm stability.

The following properties [29] for the spatial DG discretization (2.4) will be used. Let  $u$  and  $v$  be any piecewise smooth functions. A simple application of integration by parts yields the approximating skew-symmetric property

$$\mathcal{H}^\theta(u, v) + \mathcal{H}^\theta(v, u) = -\beta(2\theta - 1) \langle \llbracket u \rrbracket, \llbracket v \rrbracket \rangle_{\Gamma_h}, \tag{2.5a}$$

which implies the nonpositive property (if  $\theta > 1/2$ )

$$\mathcal{H}^\theta(u, u) = -\frac{1}{2} \beta(2\theta - 1) \|\llbracket u \rrbracket\|_{L^2(\Gamma_h)}^2, \tag{2.5b}$$

to explicitly show the numerical viscosity in the spatial discretization. Moreover, we also have the weak boundedness property (with the parameter  $\theta$ )

$$|\mathcal{H}^\theta(u, v)| \leq C \beta h^{-1} \|u\|_{L^2(I)} \|v\|_{L^2(I)}, \quad \forall u, v \in V_h, \tag{2.5c}$$

where the bounding constant  $C > 0$  depends on  $\theta$  and the inverse constant  $\mu$ .

### 2.2 The Fully Discrete ESTDG Methods

Let  $N > 0$  be any positive integer and  $\{t^n = n\tau : 0 \leq n \leq N\}$  be a uniform partition of the time interval  $[0, T]$ , where  $\tau = T/N$  is the time step. In this paper we would like to seek the numerical solution at every time level  $t^n$ , denoted by  $u^n \in V_h$ , by employing an explicit single-step time-marching algorithm to solve the semidiscrete DG method (2.3).

Suppose that  $u^n$  has been obtained at the current time level, we are able to seek  $u^{n+1}$  at the next time level through  $s$  intermediate (or generalized stage) solutions. The detailed procedure is often described in the Shu–Osher form as follows:

1. Let  $u^{n,0} = u^n$ .
2. For  $\ell = 0, 1, \dots, s - 1$ , successively find the generalized stage solution  $u^{n,\ell+1} \in V_h$  through the variational formula

$$\left(u^{n,\ell+1}, v\right)_{I_h} = \sum_{0 \leq \kappa \leq \ell} \left[ c_{\ell\kappa} \left(u^{n,\kappa}, v\right)_{I_h} + \tau d_{\ell\kappa} \mathcal{H}^{\theta_{\ell\kappa}} \left(u^{n,\kappa}, v\right) \right], \quad \forall v \in V_h. \tag{2.6}$$

Here the time-marching coefficients,  $c_{\ell\kappa}$  and  $d_{\ell\kappa}$ , are inherited from the  $r$ -th order explicit single-step algorithm. In this paper we demand  $d_{\ell\ell} \neq 0$  and  $c_{\ell\kappa} \geq 0$  for any  $\ell$  and  $\kappa$ . Note that  $s \geq r$  in general.

3. Let  $u^{n+1} = u^{n,s}$ .

The initial solution  $u^0 \in V_h$  can be set as any approximation of  $U_0$ . In this paper we define it by the local  $L^2$ -projection  $\mathbb{P}_h$ , namely

$$\left(u^0, v\right)_{I_h} = \left(U_0, v\right)_{I_h}, \quad \forall v \in V_h. \tag{2.7}$$

Till now we have completed the definition of the considered fully discrete method, which is named as the ESTDG( $s, r, k$ ) method in this paper for convenience.

We remark again that the numerical flux parameters in (2.6) are allowed to be changed at every time stage. In this paper we mainly consider two well-known examples and investigate their stability and accuracy order in the  $L^2$ -norm.

**Example 2.1** The first example is the RKDG(4, 4,  $k$ ) method with the downwind treatment [22] to deal with the negative time-marching coefficients in

$$\{c_{\ell\kappa}\} = \begin{pmatrix} 1 & & & & \\ 1/2 & 1/2 & & & \\ 1/9 & 2/9 & 2/3 & & \\ 0 & 1/3 & 1/3 & 1/3 & \end{pmatrix}, \quad \{d_{\ell\kappa}\} = \begin{pmatrix} 1/2 & & & & \\ -1/4 & 1/2 & & & \\ -1/9 & -1/3 & 1 & & \\ 0 & 1/6 & 0 & 1/6 & \end{pmatrix}, \tag{2.8}$$

where  $\ell$  and  $\kappa$  are taken from the set  $\{0, 1, 2, 3\}$  in the natural order.

To be more general than [22], we would like in this paper to take the numerical flux parameters under the following rule:  $\theta_{\ell\kappa} > 1/2$  if  $d_{\ell\kappa} \geq 0$  and  $\theta_{\ell\kappa} < 1/2$  otherwise.

**Example 2.2** The second one is the LWDG( $r, k$ ) method, which adopts the  $r$ th order Lax–Wendroff time marching to solve (2.3). This method has been proposed and analyzed in [13, 23] for  $r = 2, 3$ , with some special numerical flux parameters.

For example, the second order LWDG method [23] is given in the form

$$\begin{aligned} \left(p^n, v\right)_{I_h} &= -\mathcal{H}^{\theta_{00}} \left(u^n, v\right), \\ \left(u^{n+1}, v\right)_{I_h} &= \left(u^n, v\right)_{I_h} + \tau \mathcal{H}^{\theta_{10}} \left(u^n, v\right) - \frac{1}{2} \tau^2 \mathcal{H}^{\theta_{11}} \left(p^n, v\right), \end{aligned} \tag{2.9}$$

where  $p$  is an auxiliary variable to approximate  $\partial_t U = -\beta \partial_x U$ . As for the numerical flux parameters, the authors only take  $\theta_{00} = \theta_{10} = 1$  and  $\theta_{11}$  to be either 0 or 1. Obviously, this method can be written as an ESTDG method by defining the so-called stage solution  $u^{n,1} = -\tau p^n$ .

Actually, every  $LWDG(r, k)$  method can be written as an  $ESTDG(r, r, k)$  method with the contributory (or nonzero) parameters

$$c_{r-1,0} = 1; \quad d_{\ell\ell} = 1, \quad 0 \leq \ell \leq r - 2; \quad d_{r-1,\kappa} = \frac{1}{(\kappa + 1)!}, \quad 0 \leq \kappa \leq r - 1, \quad (2.10)$$

by similar treatments for all auxiliary variables. In this paper we would like to investigate the  $LWDG$  method in a general case and remove the technical limitations that some numerical flux parameters must be the same [23], for example,  $\theta_{00} = \theta_{r-1,0}$ .

To end this section we give a remark on the condition

$$\sum_{0 \leq \kappa \leq \ell} c_{\ell\kappa} \equiv 1, \quad \ell \geq 0,$$

which is often true for the  $RKDG$  method as a condition to ensure the consistency of the Runge–Kutta algorithm. However, (2.10) shows that the  $LWDG$  method does not satisfy this condition. Hence, in this paper we would like to discard this unessential condition for the  $ESTDG$  method and directly employ those results given in [26, 29], provided that this condition is not applied in the proofs.

### 3 Stability Analysis

In this section we devote to analyzing the  $L^2$ -norm stability for the  $ESTDG$  methods with stage-dependent numerical flux parameters. The presented analysis framework can be looked upon as an application and/or an extension of the technique of the matrix transferring process [27, 29] when numerical flux parameters are the same.

#### 3.1 The Matrix Transferring Process

In order to accurately understand the stability performance, we have to investigate the scheme when combining several time steps together in the time-marching. For this purpose, we introduce the generalized notations for stage solutions, as those in [27, 29]. Namely, for any nonnegative integers  $n, i$  and  $j$ , we denote

$$u^{n,si+j} = u^{n+i,j}. \quad (3.1)$$

Remark that this notational rule has been used in the scheme’s description.

In this paper we use an integer  $m \geq 1$  to represent the multistep number. It is evident for the  $ESTDG(s, r, k)$  method that every  $m$ -steps marching with time step  $\tau$  can be regarded as one-step marching of an  $ESTDG(ms, r, k)$  method with time step  $m\tau$ . Namely, for  $0 \leq \ell \leq ms - 1$ , there holds the following variation formula: for any  $v \in V_h$ ,

$$\left(u^{n,\ell+1}, v\right)_{I_h} = \sum_{0 \leq \kappa \leq \ell} \left[ c_{\ell\kappa}(m) \left(u^{n,\kappa}, v\right)_{I_h} + m\tau d_{\ell\kappa}(m) \mathcal{H}^{\theta_{\ell\kappa}(m)}(u^{n,\kappa}, v) \right]. \quad (3.2)$$

Let  $\ell' = \ell \pmod s$  and  $\kappa' = \kappa \pmod s$ . The contributory (or nonzero) parameters in (3.2) only emerge for those  $\ell$  and  $\kappa$  satisfying  $\ell - \ell' = \kappa - \kappa'$ , such that

$$c_{\ell\kappa}(m) = c_{\ell'\kappa'}, \quad d_{\ell\kappa}(m) = \frac{1}{m} d_{\ell'\kappa'}, \quad \theta_{\ell\kappa}(m) = \theta_{\ell'\kappa'}. \quad (3.3)$$

Here  $\ell'$  and  $\kappa'$  are both taken from  $\{0, 1, \dots, s - 1\}$ .

### 3.1.1 Temporal Differences of Stage Solutions and Evolution Identity

For  $1 \leq i \leq ms$ , we would like to define the  $i$ th order temporal difference of stage solutions in the form

$$\mathbb{D}_i(m)u^n = \sum_{0 \leq j \leq i} \sigma_{ij}(m)u^{n,j}, \tag{3.4}$$

where  $\sigma_{ij}(m)$  are the undetermined combination coefficients independent of stage solutions. For convenience, we also denote  $\mathbb{D}_0(m)u^n = u^n$  and  $\sigma_{00}(m) = 1$  throughout this paper.

**Remark 3.1** The above concepts and notations originate from the error estimates [31, 32] and have been systematically studied in [27, 29], for the RKDG methods with the same numerical flux parameters.

The combination coefficients in (3.4) can be inductively defined along the same way as in [27, 29]. Assuming, for a certain integer  $i \geq 0$ , the temporal differences of stage solutions up to the  $i$ th order have been well defined, we would like to define the next one in the form

$$\mathbb{D}_{i+1}(m)u^n = \sum_{0 \leq \ell \leq i} \phi_{i\ell}(m) \left[ u^{n,\ell+1} - \sum_{0 \leq \kappa \leq \ell} c_{\ell\kappa}(m)u^{n,\kappa} \right], \tag{3.5}$$

where the combination coefficients  $\phi_{i\ell}(m)$  will be determined by the following procedure.

Since the above linear combination does not involve any terms about spatial discretization, we can easily define the combination coefficients by the special case that all the numerical flux parameters are the same. Hence we introduce an arbitrary fixed constant, denoted by  $\vartheta$  in this paper. Due to (3.2) and (3.5), after a changing of summation orders we yield

$$\left( \mathbb{D}_{i+1}(m)u^n, v \right)_{I_h} = m\tau \Phi_i(v) + m\tau \Psi_i(v), \tag{3.6}$$

where the two terms on the right hand side show the kernel construct and the perturbation effect, respectively. They read

$$\Phi_i(v) = \sum_{0 \leq \kappa \leq i} \sum_{\kappa \leq \ell \leq i} \phi_{i\ell}(m) d_{\ell\kappa}(m) \mathcal{H}^\vartheta(u^{n,\kappa}, v), \tag{3.7a}$$

$$\Psi_i(v) = \sum_{0 \leq \kappa \leq i} \sum_{\kappa \leq \ell \leq i} \phi_{i\ell}(m) d_{\ell\kappa}(m) \left[ \mathcal{H}^{\vartheta_{\ell\kappa}}(u^{n,\kappa}, v) - \mathcal{H}^\vartheta(u^{n,\kappa}, v) \right]. \tag{3.7b}$$

We call (3.7b) the perturbation term, since  $\Psi_i(v) = 0$  if  $\vartheta_{\ell\kappa} \equiv \vartheta$ .

We want to define (3.5) to ensure a nice structure among the temporal differences of stage solutions, similar as that in [27, 29] for the RKDG method with the same numerical flux parameters. The process is described as follows. Since every diagonal entry  $d_{\kappa\kappa}(m)$  is nonzero, the triangular system of linear equations

$$\sum_{\kappa \leq \ell \leq i} \phi_{i\ell}(m) d_{\ell\kappa}(m) = \sigma_{i\kappa}(m), \quad \kappa = 0, 1, \dots, i \tag{3.8}$$

uniquely determines  $\phi_{i\ell}(m)$  for  $0 \leq \ell \leq i$ . Substituting this into (3.7a), we can achieve the same expression as that in [27, 29]

$$\Phi_i(v) = \mathcal{H}^\vartheta(\mathbb{D}_i(m)u^n, v). \tag{3.9}$$



At this moment, by comparing with the coefficients in the front of  $u^{n,\kappa}$ , on both sides of (3.5), we are able to inductively define

$$\sigma_{i+1,\kappa}(m) = \phi_{i,\kappa-1}(m) - \sum_{\kappa \leq \ell \leq i} \phi_{i,\ell}(m)c_{\ell\kappa}(m), \quad \kappa = 0, 1, \dots, i, \tag{3.10a}$$

with the supplemental notation  $\phi_{i,-1}(m) = 0$ , and

$$\sigma_{i+1,i+1}(m) = \phi_{ii}(m) = \frac{\sigma_{ii}(m)}{d_{ii}(m)} \neq 0. \tag{3.10b}$$

By these data we complete the definition of  $\mathbb{D}_{i+1}(m)u^n$ .

**Remark 3.2** In [26, 27, 29] for the RKDG method, it seems that we have demanded

$$\sum_{0 \leq j \leq i} \sigma_{ij}(m) = 0, \quad i \geq 1.$$

Actually, this condition does not take effect in any analysis therein. In this paper we would like to completely abandon this condition for the ESTDG method, since it does not hold for the LWDG method; see Sect. 3.2.2.

After all temporal differences of stage solutions have been well defined, due to (3.10b), the inversion manipulation yields the linear equivalence of two function sequences  $\{u^{n,0}, u^{n,1}, \dots, u^{n,m_s}\}$  and  $\{\mathbb{D}_0(m)u^n, \mathbb{D}_1(m)u^n, \dots, \mathbb{D}_{m_s}(m)u^n\}$ . Specially, there holds the evolution identity

$$u^{n+m} = \sum_{0 \leq i \leq m_s} \alpha_i(m)\mathbb{D}_i(m)u^n. \tag{3.11}$$

**Remark 3.3** In [27, 29], the left hand side of (3.11) was written as  $\alpha_0(m)u^{n+m}$ , where  $\alpha_0(m) > 0$  is introduced only for scaling. In this paper we always take  $\alpha_0(m) = 1$  for convenience.

Note that the above manipulations only depend on the time-marching coefficients,  $c_{\ell\kappa}$  and  $d_{\ell\kappa}$ , and they are totally independent on the numerical flux parameters. Hence all  $\sigma_{ij}(m)$  and  $\alpha_i(m)$  are the same as those when numerical flux parameters are the same; refer to [26, 27, 29] for more detailed conclusions.

The following lemma [26, Lemma 2.2] will be frequently used in this paper. It can be easily proved by the fact that the used single-step time-marching algorithm has the  $r$ th order in time.

**Lemma 3.1** For any  $m \geq 1$ , there holds  $\alpha_\ell(m) = 1/\ell!$  for  $0 \leq \ell \leq r$ .

### 3.1.2 Relationship Among Temporal Differences of Stage Solutions

In what follows we continue to discuss (3.6) and set up the relationship among temporal differences of stage solutions. A simple manipulation yields

$$\mathcal{H}^{\theta_{\ell\kappa}}(w, v) - \mathcal{H}^\vartheta(w, v) = \beta(\vartheta - \theta_{\ell\kappa}) \left\langle \llbracket w \rrbracket, \llbracket v \rrbracket \right\rangle_{\Gamma_h},$$

so the perturbation term (3.7b) can be written in the form

$$\Psi_i(v) = \beta \sum_{0 \leq \kappa \leq i} \sum_{\kappa \leq \ell \leq i} \phi_{i,\ell}(m)d_{\ell\kappa}(m)(\vartheta - \theta_{\ell\kappa}(m)) \left\langle \llbracket u^{n,\kappa} \rrbracket, \llbracket v \rrbracket \right\rangle_{\Gamma_h}. \tag{3.12}$$

To express the right hand side in terms of temporal differences of stage solutions, we would like to introduce a series of numbers  $q_{i\ell}(m; \vartheta)$ , for  $0 \leq \ell \leq i$ , by the triangular system of linear equations: for  $\kappa = 0, 1, \dots, i$ ,

$$\sum_{\kappa \leq \ell \leq i} q_{i\ell}(m; \vartheta) \sigma_{\ell\kappa}(m) = - \sum_{\kappa \leq \ell \leq i} \phi_{i\ell}(m) d_{\ell\kappa}(m) (\vartheta - \theta_{\ell\kappa}(m)). \tag{3.13}$$

The existence and uniqueness are trivial since every diagonal entry  $\sigma_{\kappa\kappa}(m)$  is nonzero, due to (3.10b). By substituting (3.13) into the previous identity and changing the summary order, we can deduce

$$\begin{aligned} \Psi_i(v) &= -\beta \sum_{0 \leq \kappa \leq i} \sum_{\kappa \leq \ell \leq i} q_{i\ell}(m; \vartheta) \sigma_{\ell\kappa}(m) \left( \llbracket u^{n,\kappa} \rrbracket, \llbracket v \rrbracket \right)_{\Gamma_h} \\ &= -\beta \sum_{0 \leq \ell \leq i} q_{i\ell}(m; \vartheta) \left( \llbracket \mathbb{D}_\ell(m) u^n \rrbracket, \llbracket v \rrbracket \right)_{\Gamma_h}, \end{aligned} \tag{3.14}$$

where the definition of temporal differences of stage solutions, like (3.4), is used at the last step.

Substituting (3.9) and (3.14) into (3.6), we eventually achieve the relationship among the temporal differences of stage solutions: for any  $v \in V_h$ , there holds

$$\begin{aligned} \left( \mathbb{D}_{i+1}(m) u^n, v \right)_{I_h} &= m\tau \mathcal{H}^\vartheta \left( \mathbb{D}_i(m) u^n, v \right) \\ &\quad - m\tau\beta \sum_{0 \leq \ell \leq i} q_{i\ell}(m; \vartheta) \left( \llbracket \mathbb{D}_\ell(m) u^n \rrbracket, \llbracket v \rrbracket \right)_{\Gamma_h}. \end{aligned} \tag{3.15}$$

It is worthy to mention that the right hand side is independent of the choice of  $\vartheta$ , hence its value can be set arbitrarily.

At the end of this subsection, we present the kernel relationship that will be extensively used in the matrix transferring process. For convenience of notations, we would like to denote a series of quantities independent of the choice of  $\vartheta$ , namely

$$\tilde{q}_{i\ell}(m) = q_{i\ell}(m; \vartheta) + \delta_{i\ell} \vartheta, \quad 0 \leq \ell \leq i \text{ and } 0 \leq i \leq ms - 1. \tag{3.16}$$

Throughout this paper  $\delta_{i\ell}$  is a standard Kronecker symbol, being 1 if  $i = \ell$  and otherwise 0. The independence is easily verified, because (3.16) satisfies the triangular system of linear equations that is independent of  $\vartheta$ ,

$$\sum_{\kappa \leq \ell \leq i} \tilde{q}_{i\ell}(m) \sigma_{\ell\kappa}(m) = \sum_{\kappa \leq \ell \leq i} \phi_{i\ell}(m) d_{\ell\kappa}(m) \theta_{\ell\kappa}(m), \quad \kappa = 0, 1, \dots, i,$$

due to (3.13) and (3.8).

We will see later two kinds of fundamental members in the matrix transferring process. One is the joint of two  $L^2(I_h)$ -inner products terms (named as the temporal information terms)

$$\mathcal{J}(i, j) = \left( \mathbb{D}_{i+1}(m) u^n, \mathbb{D}_j(m) u^n \right)_{I_h} + \left( \mathbb{D}_i(m) u^n, \mathbb{D}_{j+1}(m) u^n \right)_{I_h}, \tag{3.17}$$

and the other is the  $L^2(\Gamma_h)$ -inner product term (the essential ingredient of the spatial information)

$$\mathcal{P}(i, j) = \left( \llbracket \mathbb{D}_i(m) u^n \rrbracket, \llbracket \mathbb{D}_j(m) u^n \rrbracket \right)_{\Gamma_h}.$$

The kernel relationship is stated in the following lemma.

**Lemma 3.2** For  $0 \leq i, j \leq ms - 1$ , there holds

$$\mathcal{J}(i, j) = -m\tau\beta \left[ -\mathcal{P}(i, j) + \sum_{0 \leq i' \leq i} \tilde{q}_{ii'}(m)\mathcal{P}(i', j) + \sum_{0 \leq j' \leq j} \tilde{q}_{jj'}(m)\mathcal{P}(i, j') \right].$$

**Proof** This lemma follows from (3.15), (3.16) and an application of (2.5a). □

**Remark 3.4** Suppose all numerical flux parameters are the same, say,  $\theta_{\ell\kappa} \equiv \theta$ . Taking the fixed parameter  $\vartheta = \theta$ , it is easy to see  $q_{\ell\kappa}(m; \theta) = 0$  due to (3.13) and hence  $\tilde{q}_{\ell\kappa}(m) = \theta\delta_{\ell\kappa}$  due to (3.16). Then we have

$$\mathcal{J}(i, j) = -m\tau\beta(2\theta - 1)\mathcal{P}(i, j),$$

from the above lemma. This result is the same as that in [29].

### 3.1.3 Derivation of Energy Equations

Along the same line as that in the previous works [26, 27, 29], we would like to carry out the matrix transferring process to automatically achieve a perfect energy equation for the considered ESTDG method, through a sequence of energy equations

$$\|u^{n+m}\|_{L^2(I)}^2 - \|u^n\|_{L^2(I)}^2 = \text{TM}(\ell; m) + \text{SP}(\ell; m). \tag{3.18}$$

Here  $\ell \geq 0$  stands for the sequence number, and

$$\text{TM}(\ell; m) = \sum_{0 \leq i \leq ms} \sum_{0 \leq j \leq ms} a_{ij}^{(\ell)}(m) \left( \mathbb{D}_i(m)u^n, \mathbb{D}_j(m)u^n \right)_{I_h}, \tag{3.19a}$$

$$\text{SP}(\ell; m) = -m\tau\beta \sum_{0 \leq i \leq ms} \sum_{0 \leq j \leq ms} b_{ij}^{(\ell)}(m) \left\langle \llbracket \mathbb{D}_i(m)u^n \rrbracket, \llbracket \mathbb{D}_j(m)u^n \rrbracket \right\rangle_{\Gamma_h}, \tag{3.19b}$$

respectively contain all temporal information and all spatial information. For convenience, we abbreviate two formulas in (3.19) by two symmetric matrices

$$\mathbb{A}^{(\ell)}(m) = \{a_{ij}^{(\ell)}(m)\}_{0 \leq i, j \leq ms}, \quad \mathbb{B}^{(\ell)}(m) = \{b_{ij}^{(\ell)}(m)\}_{0 \leq i, j \leq ms}. \tag{3.20}$$

**Remark 3.5** It is worthy to mention that (3.19b) is different to that in [26, 27, 29] for the RKDG methods. Actually, the modification in this paper originates from the application of the approximating skew-symmetric property (2.5a); see Lemma 3.2.

As a purpose of matrix transferring process, we expect to dig out the contribution of the spatial discretization as much as possible, by successively transforming the lower order temporal information into the spatial information. To show that, in what follows we give a more detailed description on the matrix transferring process.

The initial energy equation is easily derived by squaring and integrating the evolution identity (3.11). It deduces the initial matrices with

$$a_{ij}^{(0)}(m) = \begin{cases} 0, & i = j = 0, \\ \alpha_i(m)\alpha_j(m), & \text{otherwise;} \end{cases} \quad \text{and} \quad b_{ij}^{(0)}(m) = 0, \tag{3.21}$$

with  $\alpha_0(m) \equiv 1$  as stated in Remark 3.3. Remark that this energy equation does not reflect any contribution of the spatial discretization.

The matrix transferring process is carried out step by step. By induction, for  $\ell \geq 1$ , the  $\ell$ th step matrix transform starts from two obtained matrices

$$\mathbb{A}^{(\ell-1)} = \begin{pmatrix} \textcircled{0} & \textcircled{0} & \textcircled{0} & \cdots \\ \textcircled{0} & a_{\ell-1,\ell-1}^{(\ell-1)} & a_{\ell-1,\ell}^{(\ell-1)} & \cdots \\ \textcircled{0} & a_{\ell,\ell-1}^{(\ell-1)} & a_{\ell,\ell}^{(\ell-1)} & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}, \quad \mathbb{B}^{(\ell-1)} = \begin{pmatrix} \star & \star & \star & \cdots \\ \star & b_{\ell-1,\ell-1}^{(\ell-1)} & b_{\ell-1,\ell}^{(\ell-1)} & \cdots \\ \star & b_{\ell,\ell-1}^{(\ell-1)} & 0 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix},$$

where  $\textcircled{0}$  remarks the zero block and  $\star$  remarks the transformed (nonzero) region. Here and below the notation  $(m)$  is dropped for convenience if there is no confusion.

The next action depends on the leading element  $a_{\ell-1,\ell-1}^{(\ell-1)}(m)$ . If it is equal to zero, we carry out the following manipulations. In this step, we would like to use Lemma 3.2 to eliminate every entry at the  $(\ell - 1)$ th row and column of  $\mathbb{A}^{(\ell-1)}$ . This process generates two new matrices  $\mathbb{A}^{(\ell)}$  and  $\mathbb{B}^{(\ell)}$ .

More specifically, for the temporal matrix  $\mathbb{A}^{(\ell)}$ , the entries at the lower triangular region are given by the following formulas

$$a_{ij}^{(\ell)} = \begin{cases} 0, & \ell - 1 \leq i \leq ms \text{ and } j = \ell - 1, \\ a_{ij}^{(\ell-1)} - 2a_{i+1,j-1}^{(\ell-1)}, & i = \ell \text{ and } j = \ell, \\ a_{ij}^{(\ell-1)} - a_{i+1,j-1}^{(\ell-1)}, & \ell + 1 \leq i \leq ms - 1 \text{ and } j = \ell, \\ a_{ij}^{(\ell-1)}, & \text{otherwise.} \end{cases} \tag{3.22}$$

Since  $\mathbb{A}^{(\ell)}$  is demanded to be symmetric, the upper triangular entry is easily filled in. To understand the above formula in (3.22), we give some comments.

- The difference between the second and the third line results from whether the basic elimination (with respect to one entry) along the row and column is superimposed on the same position.
- The third line is used to eliminate  $a_{i+1,\ell-1}^{(\ell-1)}$  by applying Lemma 3.2 to the joint term  $a_{i+1,\ell-1}^{(\ell-1)}\mathcal{J}(i, \ell - 1)$ , with the help of the neighbor entry  $a_{i,\ell}^{(\ell-1)}$ .

Remark that the order of the left-top zero block is enlarged now.

The above elimination is accompanied by the changing of the spatial matrix. For stage-dependent numerical flux parameters, each basic elimination affects many entries of the spatial matrix. For example, the basic elimination on  $a_{i+1,\ell-1}^{(\ell-1)}$  (corresponding to the third line in (3.22)) has influence on those entries of the spatial matrix at both the left half-line and the top half-line starting from the position  $(i, \ell - 1)$ . As a result, it is not easy to present short and unified formulas for calculating each entry of the new spatial matrix  $\mathbb{B}^{(\ell)}$ . However, the manipulation process can be conveniently expressed in the pseudo-code and summarized as **Algorithm 1**.

**Remark 3.6** Recalling Remark 3.4 for the same numerical flux parameters, it is easy to see that the two sub-loops in Step 2 really execute only for  $i = \kappa$  and  $j = \ell - 1$ , respectively. Note that  $\tilde{q}_{\kappa,\kappa} = \tilde{q}_{\ell-1,\ell-1} = \vartheta = \theta$  in this case.

Otherwise, if  $a_{\ell-1,\ell-1}^{(\ell-1)}(m) \neq 0$ , we stop the entire transform process and name this entry as the *central objective* of temporal matrix. At the same time, we output the *termination index* of time marching

$$\zeta(m) = \ell - 1, \tag{3.23}$$

---

**Algorithm 1.** Generate the spatial matrix  $\mathbb{B}^{(\ell)} = \{b_{ij}^{(\ell)}\}$  for the given  $\ell$ .

---

*Step 1.* Initialization: set  $g_{ij} = 0$  for any  $0 \leq i, j \leq ms$ ;  
*Step 2.* Modification: for  $\kappa = \ell - 1, \dots, ms - 1$ , do  
 if  $\kappa = \ell - 1$  then let  $\nu = 1/2$ ; otherwise,  $\nu = 1$ ;  
 compute  $g_{\kappa, \ell-1} \leftarrow g_{\kappa, \ell-1} - \nu a_{\kappa+1, \ell-1}^{(\ell-1)}$ ;  
 compute  $g_{i, \ell-1} \leftarrow g_{i, \ell-1} + \nu a_{\kappa+1, \ell-1}^{(\ell-1)} \tilde{q}_{\kappa, i}$  for  $i = 0, \dots, \kappa$ ;  
 compute  $g_{\kappa, j} \leftarrow g_{\kappa, j} + \nu a_{\kappa+1, \ell-1}^{(\ell-1)} \tilde{q}_{\ell-1, j}$  for  $j = 0, \dots, \ell - 1$ ;  
*Step 3.* Generation: define  $b_{ij}^{(\ell)} = b_{ij}^{(\ell-1)} + g_{ij} + g_{ji}$  for  $0 \leq i, j \leq ms$ .

---

together with the ultimate temporal matrix and the ultimate spatial matrix, respectively denoted by

$$\begin{aligned} \mathbb{A}(m) &= \mathbb{A}^{(\zeta(m))}(m) = \{a_{ij}(m)\}_{0 \leq i, j \leq ms}, \\ \mathbb{B}(m) &= \mathbb{B}^{(\zeta(m))}(m) = \{b_{ij}(m)\}_{0 \leq i, j \leq ms}. \end{aligned}$$

Motivated by [26, 27, 29], it is important for the ultimate spatial matrix to find the largest order of the SPD sequential principal submatrix, i.e.,

$$\rho(m) = \max \left\{ \kappa : 1 \leq \kappa \leq \zeta \text{ and } \{b_{ij}(m)\}_{0 \leq i, j \leq \kappa-1} \text{ is SPD} \right\}. \tag{3.24}$$

This quantity is also named as the *contribution index* of spatial discretization. If the set in (3.24) is empty, we define  $\rho(m) = 0$  as a supplement.

Till now we have completed the description of the matrix transferring process. The stability performance of the ESTDG method will be determined by three important quantities: two indices  $\zeta(m)$  and  $\rho(m)$ , as well as the sign of central objective. See Sect. 3.2 for details.

### 3.1.4 Discussion on Three Important Quantities

Since the ultimate temporal matrix  $\mathbb{A}(m)$  solely depends on the time-marching coefficients, we have the same conclusions as those in [26] for the RKDG method with the same numerical flux parameters. Below we list some conclusions that will be used in this paper.

**Lemma 3.3** *The termination index  $\zeta(m) \geq 1$  is independent of  $m$ , and hence we denote it by  $\zeta$  throughout this paper.*

**Lemma 3.4** *For any  $m \geq 1$ , the central objective  $a_{\zeta\zeta}(m)$  keeps the same sign.*

The ultimate spatial matrix  $\mathbb{B}(m)$  depends on not only the time-marching coefficients but also the numerical flux parameters. Hence, for the time-dependent numerical flux parameters, the property of  $\rho(m)$  becomes a little complex and the corresponding analysis turns out to be much more difficult.

We begin with the assumption that  $\rho(m)$  is always positive, in view of the practical application. This is equivalent to  $b_{00}(m) > 0$  for any  $m \geq 1$ . When all numerical flux parameters are taken to be  $\theta > 1/2$ , we have found out in [26] for the RKDG method that

$$\frac{1}{2} [b_{00}(m) + 1] = \theta > \frac{1}{2}, \quad m \geq 1. \tag{3.25}$$

As an extension of this conclusion, we would like to propose an important concept for the ESTDG method with stage-dependent numerical flux parameters.

**Definition 3.1** For the ESTDG method, the averaged numerical flux parameter every  $m$ -steps marching is defined by

$$\Theta(m) \equiv \frac{1}{2} [b_{00}(m) + 1], \quad m \geq 1. \tag{3.26}$$

Especially,  $\Theta = \Theta(1)$  is called the averaged numerical flux parameter.

To well understand the above definition, we need to make more detailed discussions on (3.26). From **Algorithm 1**, it is easy to see that

$$b_{00}(m) \equiv b_{00}^{(1)}(m) = -a_{10}^{(0)}(m) + \sum_{0 \leq \ell \leq ms-1} 2a_{\ell+1,0}^{(0)}(m)\tilde{q}_{\ell,0}(m),$$

which is determined at the first step of the matrix transferring process. Noticing (3.21) and Lemma 3.1, we have  $a_{10}^{(0)}(m) = \alpha_0(m)\alpha_1(m) = 1$  and  $a_{\ell+1,0}^{(0)}(m) = \alpha_{\ell+1}(m)$ . Then it follows from (3.26) that

$$\Theta(m) = \sum_{0 \leq \ell \leq ms-1} \alpha_{\ell+1}(m)\tilde{q}_{\ell,0}(m). \tag{3.27}$$

As a direct application of this formula, we can easily find out the hidden zero restriction that will be used to analyze the performance of  $\rho(m)$  as  $m$  goes to infinity and used to obtain the optimal error estimate. This important conclusion is stated in the following lemma.

**Lemma 3.5** *There holds for any  $m \geq 1$  that*

$$\sum_{0 \leq \ell \leq ms-1} \alpha_{\ell+1}(m)q_{\ell,0}(m; \Theta(m)) = 0. \tag{3.28}$$

**Proof** We can prove this lemma by (3.16) and taking  $\vartheta = \Theta(m)$  in (3.27). □

Similar as (3.25) for the fixed numerical flux parameters, we have the following lemma for variant numerical flux parameters. The proof is put aside in the appendix, since it shares many materials in the proof of the next lemma.

**Lemma 3.6**  *$\Theta(m)$  is independent of  $m$ , namely  $\Theta(m) = \Theta$ .*

From our point of view,  $\Theta$  is an essential quantity to accurately describe the upwind attribute for the fully discrete method. Owing to Lemma 3.6, the assumption that  $\rho(m) \geq 1$  holds forever is equivalent to demand

$$\Theta > 1/2, \tag{3.29}$$

which means the upwind mechanism at least in the average sense. Actually, this demand plays a critical role in the whole analysis of this paper.

**Lemma 3.7** *If  $\Theta > 1/2$ , then there is an  $m_\star \geq 1$  such that  $\rho(m) = \zeta$  for  $m \geq m_\star$ .*

The proof line of this lemma is the same as that in [26] for the RKDG method with the same numerical flux parameters. However, the stage-dependent numerical flux parameters cause serious analysis difficulties such that the proof process involves many matrix manipulation and looks much lengthy and technical. We would like to postpone the proof of this lemma to the appendix and only present the key points in the proof.

1. **Algorithm 1** is not convenient to carry out the analysis, and we have to set up a matrix description for the ultimate spatial matrix  $\mathbb{B}(m)$ . In this process, many tricks are used to get some convenient and unified formulas, especially for the  $\zeta$ th order sequential principle submatrix of  $\mathbb{B}(m)$ , which is denoted by  $\tilde{\mathbb{B}}(m)$  for convenience of statement.

2. Roughly speaking, in order to prove Lemma 3.7, we would like to split  $\tilde{\mathbb{B}}(m)$  into two symmetrical matrices for any given parameter  $\vartheta$ . Although this matrix is independent of  $\vartheta$ , finding a good separation with a suitable choice of  $\vartheta$  is important for the theoretical analysis.
  - One matrix is just the same as that for the special case that all numerical flux parameters are taken to be  $\vartheta$ . Provided  $\vartheta > 1/2$ , we can prove similarly as in [26] that this matrix tends to a special SPD matrix as  $m$  goes to infinity.
  - The other matrix results from the perturbation of stage-dependent numerical flux parameters with respect to the fixed  $\vartheta$ . As a trivial purpose, we expect that this perturbation matrix tends to zero as  $m$  goes to infinity. In general, this purpose is not easily accomplished for arbitrary choice of  $\vartheta$ . However, this aim is fortunately addressed with the help of a special choice  $\vartheta = \Theta$ , owing to the hidden zero restriction (3.28) in Lemma 3.5.
3. In order to achieve the second goal in the previous item, we need to reveal the relationship of the perturbation matrix with regard to the multistep number  $m$ . To do that, a large number of matrix manipulations (including Kronecker products of matrix) are executed. This simplification process is long and technical.

To end this subsection, we would like to show how to ensure  $\Theta > 1/2$  by adjusting the numerical flux parameters. This purpose can be implemented by the following two propositions, whose proofs will be given in the appendix.

**Proposition 3.1**  *$\Theta$  is a weighted average of the numerical flux parameters. Moreover, it increases with respect to  $\theta_{\ell\kappa}$  if  $d_{\ell\kappa} > 0$  and decreases otherwise.*

For the RKDG method, the averaged numerical flux parameter often depends on all numerical flux parameters. For example, the RKDG(4, 4,  $k$ ) method (2.8) has

$$\Theta = \frac{37}{108}\theta_{00} - \frac{5}{36}\theta_{10} + \frac{5}{18}\theta_{11} - \frac{1}{27}\theta_{20} - \frac{1}{9}\theta_{21} + \frac{1}{3}\theta_{22} + \frac{1}{6}\theta_{31} + \frac{1}{6}\theta_{33}.$$

However, it is not true for the LWDG method.

**Proposition 3.2** *For the LWDG( $r, k$ ) method we always have  $\Theta = \theta_{r-1,0}$ .*

Proposition 3.2 gives a theoretical support to the upwind requirement  $\theta_{r-1,0} > 1/2$  for the LWDG method, which has been implicitly stressed in [13, 23]. This is to say, only the first order term in the Lax–Wendroff expansion demands the DG discretization with upwind numerical flux, and the other term can be arbitrarily discretized.

### 3.2 Energy Analysis and Stability Conclusions

The matrix transferring process yields the final energy Eq. (3.18) for any  $m \geq 1$ , with the termination index  $\ell = \zeta$ , the sign of the central objective, and the contribution index  $\rho(m)$ . Based on these informations, we are able to easily carry out the energy analysis and conclude the  $L^2$ -norm stability performance along the same line as that in [29].

It is worthy pointing out that the stage-dependent numerical flux parameters do not cause any essential analysis difficulty in this subsection. To shorten the length of this paper, we only present the key steps and conclusions and point out the main modifications in this process.

The increment every  $m$  steps is still bounded in the form

$$\|u^{n+m}\|_{L^2(I)}^2 - \|u^n\|_{L^2(I)}^2 + \mathcal{S}_{\text{bry}} \leq a_{\zeta\zeta}(m) \|\mathbb{D}_{\zeta}(m)u^n\|_{L^2(I)}^2 + \mathcal{S}_{\text{hot}}, \tag{3.30}$$

where

$$\begin{aligned}
 \mathcal{S}_{\text{bry}} &= \varepsilon_*(m)m\tau\beta \sum_{0 \leq \ell < \rho(m)} \|\mathbb{D}_\ell(m)u^n\|_{L^2(\Gamma_h)}^2, \\
 \mathcal{S}_{\text{hot}} &= C(m) \sum_{i,j \geq \zeta \text{ except } i=j=\zeta} \left| \left( \mathbb{D}_i(m)u^n, \mathbb{D}_j(m)u^n \right)_{I_h} \right| \\
 &\quad + C(m) \sum_{\max(i,j) \geq \rho(m)} \left| \left( \mathbb{D}_i(m)u^n, \mathbb{D}_j(m)u^n \right)_{\Gamma_h} \right| \tau.
 \end{aligned}$$

Here  $\varepsilon_*(m)$  is the smallest eigenvalue of the SPD submatrix  $\{b_{ij}(m)\}_{0 \leq i,j < \rho(m)}$ , and  $C(m)$  means the generic constant independent of  $n, h$  and  $\tau$ .

All terms in  $\mathcal{S}_{\text{hot}}$  can be well controlled by the inverse inequality to the jump term of temporal differences

$$(\tau\beta\lambda)^{\frac{1}{2}} \|\mathbb{D}_\kappa(m)u^n\|_{L^2(\Gamma_h)} \leq C\lambda \|\mathbb{D}_\kappa(m)u^n\|_{L^2(I)}, \tag{3.31a}$$

and the relationship among temporal differences of stage solutions

$$\begin{aligned}
 \|\mathbb{D}_{i+1}(m)u^n\|_{L^2(I)} &\leq C\lambda \|\mathbb{D}_i(m)u^n\|_{L^2(I)} \\
 &\quad + C(\tau\beta\lambda)^{\frac{1}{2}} \sum_{0 \leq \ell \leq i} \|\mathbb{D}_\ell(m)u^n\|_{L^2(\Gamma_h)},
 \end{aligned} \tag{3.31b}$$

where  $\lambda = \tau\beta/h$  is the CFL number. The inequality (3.31b) can be easily obtained by taking  $v = \mathbb{D}_{i+1}(m)u^n$  in (3.15) and using (2.5c). The last summation on the right hand side of (3.31b) originates from the perturbation of numerical flux parameters, and thus we inevitably encounter some terms involved the jumps of lower order temporal differences in the estimating process to  $\mathcal{S}_{\text{hot}}$ . This is the only analysis difference when numerical flux parameters are stage-dependent. It is worthy pointing out that we can further diminish the jump norms for those temporal differences of order not less than  $\rho(m)$  in (3.31b) if  $i \geq \rho(m)$ , by an inductive application of two inequalities in (3.31). Namely, we have for  $i \geq \rho(m)$  that

$$\|\mathbb{D}_{i+1}(m)u^n\|_{L^2(I)}^2 \leq C\lambda^2 \|\mathbb{D}_{\rho(m)}(m)u^n\|_{L^2(I)}^2 + C\lambda \mathcal{S}_{\text{bry}}. \tag{3.32}$$

By these treatments, together with some applications of Cauchy–Schwartz inequality, we can easily bound  $\mathcal{S}_{\text{hot}}$  in the form

$$\mathcal{S}_{\text{hot}} \leq C\lambda \|\mathbb{D}_{\rho(m)}(m)u^n\|_{L^2(I)}^2 + \left(\frac{1}{2} + C\lambda\right) \mathcal{S}_{\text{bry}}.$$

As long as  $\lambda$  is small enough, the last term in this inequality can be controlled with the help of  $\mathcal{S}_{\text{bry}}$ . Since the obtaining inequality is almost the same as that in the previous works [29], the stability results can be similarly stated if the polynomials degree  $k$  is not specified.

Along the same line as for (3.32) we can similarly have

$$\|\mathbb{D}_{\rho(m)}(m)u^n\|_{L^2(I)}^2 \leq C\lambda^{2\rho(m)} \|u^n\|_{L^2(I)}^2 + C\lambda \mathcal{S}_{\text{bry}}. \tag{3.33}$$

Note that  $\|\mathbb{D}_\zeta(m)u^n\|_{L^2(I)}^2$  can be bounded by either (3.33) or (3.32), since  $\zeta \geq \rho(m)$  due to their definitions. Summing up the above discussions into (3.30), we have the rough estimate

$$\|u^{n+m}\|_{L^2(I)}^2 \leq \left[ 1 + C\lambda^{\min(2\zeta, 2\rho(m)+1)} \right] \|u^n\|_{L^2(I)}^2,$$

Together with Lemma 3.7 for sufficient large  $m$ , we can easily obtain the weak stability, as stated in the next theorem. This conclusion does not consider the effect of the sign of the central objective.



**Theorem 3.1** *The ESTDG method (2.6) has the weak( $2\zeta$ ) stability at least. Namely, for sufficiently small  $h$  there holds*

$$\|u^n\|_{L^2(I)} \leq C \|u^0\|_{L^2(I)}, \quad n \geq 0, \tag{3.34}$$

*under a severe temporal–spatial condition  $\tau \leq Mh^{\frac{2\zeta}{2\zeta-1}}$ . Here  $M$  is any given positive constant, and the bounding constant  $C = C(T, M)$  is independent of  $n, h$  and  $\tau$ .*

As usual, we pay more attention on the stability under suitable CFL conditions. To this end, we introduce an important quantity

$$n_\star = \min \left\{ m : \rho(m) = \rho(m + 1) = \dots = \rho(2m - 1) = \zeta \right\}, \tag{3.35}$$

which is not larger than  $m_\star$  due to Lemma 3.7. Actually, we have proved in [26, Proposition 3.5] that  $n_\star = m_\star$  holds for many RKDG method with the same numerical flux parameters. Although we can not prove this conclusion for any ESTDG method with stage-dependent numerical flux parameters, numerical experiments proposed in this paper indicate that this statement might be true.

Note that the negativity of the central objective plays a pivotal role in the next theorem.

**Theorem 3.2** *If the central objective is negative, the ESTDG method (2.6) has the strong( $n_\star$ ) stability for any  $k \geq 0$ , namely, there exists a maximal CFL number  $\lambda_{\max}$  such that*

$$\|u^n\|_{L^2(I)} \leq \|u^0\|_{L^2(I)}, \quad n \geq n_\star, \tag{3.36}$$

*holds under the CFL condition  $\lambda \leq \lambda_{\max}$ . Furthermore, if  $n_\star = 1$  is allowed, there holds the monotonicity stability in the sense that*

$$\|u^{n+1}\|_{L^2(I)} \leq \|u^n\|_{L^2(I)}, \quad n \geq 0. \tag{3.37}$$

**Remark 3.7** Actually, the strong stability is obtained by the monotonicity stability for the corresponding ESTDG( $m, r, k$ ) method (3.2), where the multistep number  $m$  goes through  $n_\star, n_\star + 1, \dots, 2n_\star - 1$ . Detailed discussions can be found in [27, 29].

Along the same line as that in [29], we can similarly obtain a nice control among the temporal differences of stage solutions. For instance, the first term on the right hand side of (3.31b) can be replaced by  $C \|(m\tau\beta\partial_x)\mathbb{D}_i(m)u^n\|_{L^2(I)}$ , which helps us to enhance the stability performance for piecewise polynomials of lower degree. Detailed discussions are referred to [29]. The related conclusions are stated in the next theorem.

**Theorem 3.3** *The ESTDG method (2.6) can have a better stability performance for lower degree  $k$ :*

- *the strong( $n_\star$ ) stability for  $k < \zeta$ , if the central objective is positive.*
- *the monotonicity stability for  $k < \rho(1)$  no matter whether the central objective is positive or negative.*

From the last two theorems we are happy to find out an opportunity to enlarge the contribution index  $\rho(m)$  and get better stability performance, by means of suitably adjusting the numerical flux parameters. If so, the quantity  $n_\star$  may become smaller, even to 1 so that the strong stability is improved to the monotonicity stability. In the next subsections we will give some detailed discussions on two examples given in Sect. 2.2.

### 3.2.1 The RKDG Method

Consider the RKDG(4, 4,  $k$ ) method proposed in Example 2.1, and take the numerical flux parameters

$$\left\{ \theta_{\ell k} - \frac{1}{2} \right\} = \varepsilon \begin{pmatrix} 1 & & & & \\ -1 & 1 & & & \\ -1 & -y & z & & \\ & & & 1 & 0 & 1 \end{pmatrix}, \tag{3.38}$$

where  $\varepsilon$ ,  $y$  and  $z$  are positive constants. Three negative entries in the right matrix correspond to the so-called downwind treatment.

Below we take  $z = 1$  and focus on the effect of  $y$ . We begin the stability analysis with  $m = 1$ . The temporal differences of stage solutions are defined as

$$\{\sigma_{ij}(1)\} = \begin{pmatrix} 1 & & & & \\ -2 & 2 & & & \\ 0 & -4 & 4 & & \\ 4 & 0 & -8 & 4 & \\ 8 & 0 & -16 & -16 & 24 \end{pmatrix}, \quad 0 \leq i, j \leq 4,$$

and the numerical flux parameters lead to

$$\left\{ \tilde{q}_{ij}(1) - \frac{1}{2} \delta_{ij} \right\} = \varepsilon \begin{pmatrix} 1 & & & & \\ 2 & & 1 & & \\ -4/9 + 4y/3 & 2/3 + 2y/3 & 1 & & \\ -100/9 - 8y/3 & -4/3 - 4y/3 & 0 & 1 & \end{pmatrix}, \quad 0 \leq i, j \leq 3.$$

The matrix transferring process gives two matrices. The first one is the ultimate temporal matrix

$$\mathbb{A}(1) = \begin{pmatrix} \mathbb{O}_3 & & \\ -1/72 & 1/144 & \\ 1/144 & 1/576 & \end{pmatrix},$$

where  $\mathbb{O}_3$  is the third order zero matrix. This matrix implies that the termination index of time marching is  $\zeta = 3$  and the central objective satisfies  $a_{\zeta\zeta}(1) = -1/72 < 0$ . The second one is the ultimate spatial matrix

$$\mathbb{B}(1) = \varepsilon \begin{pmatrix} 2y/9 + 79/27 & y/9 + 65/54 & 1/3 & y/36 + 17/108 & 0 \\ y/9 + 65/54 & y/18 + 13/18 & 1/4 & y/72 + 7/72 & 0 \\ 1/3 & 1/4 & 1/12 & 1/24 & 0 \\ y/36 + 17/108 & y/72 + 7/72 & 1/24 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

of which the first three leading principle determinants are

$$\varepsilon \left( \frac{2y}{9} + \frac{79}{27} \right), \quad \varepsilon^2 \left( \frac{y}{18} + \frac{1937}{2916} \right), \quad \varepsilon^3 \left( \frac{y}{324} - \frac{125}{17496} \right). \tag{3.39}$$

The first quantity indicates that the averaged numerical flux parameter indeed satisfies Proposition 3.1. Now we can claim the following stability results:

- For  $y > 125/54$ , these three quantities are all positive and hence  $\rho(1) = 3 = \zeta$ . Now we can claim the monotonicity stability for  $k \geq 0$  by Theorem 3.2.

- For  $y < 125/54$ , the stability performance becomes a little weaker. In this case, only the first two quantities in (3.39) are positive, and thus  $\rho(1) = 2$  becomes smaller. A series of matrix transferring process for multisteps time-marching yields  $\rho(2) = \rho(3) = 3 = \zeta$ , as we have predicted in Lemma 3.7. By Theorems 3.2 and 3.3 we can claim the strong(2) stability for any  $k \geq 0$  and the monotonicity stability only for  $k \leq 1$ . The sharpness of this statement will be shown in the numerical experiments.

**Remark 3.8** Consider the same RKDG method with  $y = 1$ , and focus on the effect of  $z$ . The matrix transferring process derives that  $\rho(1) = \zeta = 3$  (i.e., the monotonicity stability) hold only for  $z < (2\sqrt{598} - 37)/27 \approx 0.441$ . If we increase  $z$  out of the above region, although  $\Theta = \varepsilon(67/54 + z/3) + 1/2$  becomes larger, the stability performance is weakened to the strong stability.

From this numerical example (or the LWDG method), we can see that sometimes the stability performance may be not improved by enlarging the averaged numerical flux parameter. This contradicts the commonly accepted concept that the greater numerical viscosity provides the better stability. Hence it seems to show that this quantity should not be simply understood as the numerical viscosity coefficients for the ESTDG methods of stage-dependent numerical flux parameters.

### 3.2.2 The LWDG Method

We now turn to the LWDG( $r, k$ ) method for  $r \leq 5$ ; see Example 2.2. For simplicity, all involved numerical flux parameters are taken to be  $1/2 \pm \varepsilon$ , where  $\varepsilon$  is a positive constant. Due to Proposition 3.2, we must set  $\theta_{r-1,0} = 1/2 + \varepsilon$  to ensure  $\Theta > 1/2$  for all cases.

Let us take the second order ( $r = 2$ ) LWDG method as an example. By the matrix transferring process we can obtain

$$\{\sigma_{ij}(1)\}_{0 \leq i, j \leq 2} = \begin{pmatrix} 1 & & \\ 0 & 1 & \\ -2 & -2 & 2 \end{pmatrix} \quad \text{and} \quad \mathbb{A}(1) = \begin{pmatrix} 0 & & \\ 0 & & \\ & & 1/4 \end{pmatrix},$$

and get  $\zeta = 2$  and  $a_{\zeta}(1) = 1/4$ . Due to Theorem 3.1, we claim that this method at least has the weak(4) stability for  $k \geq 0$ .

Due to Theorem 3.3, we can get the strong stability and/or monotonicity stability for lower degree  $k$ . For different combination of  $\theta_{00}$  and  $\theta_{11}$ , we may achieve different values of  $n_*$  by calculating the contribution index of spatial discretization as  $m$  increases. The detailed conclusions are listed as follows.

- Let  $\theta_{00} = \theta_{11} = 1/2 + \varepsilon$ . We get  $\rho(1) = 2 = \zeta$ , since

$$\{\tilde{q}_{ij}(1)\}_{0 \leq i, j \leq 1} = \begin{pmatrix} 1/2 + \varepsilon & \\ & 1/2 + \varepsilon \end{pmatrix}, \quad \{b_{ij}(1)\}_{0 \leq i, j \leq 1} = \varepsilon \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}.$$

This concludes the monotonicity stability for  $k \leq 1$ .

- Let  $\theta_{00} = 1/2 + \varepsilon$  and  $\theta_{11} = 1/2 - \varepsilon$ . Let  $m = 1$  and we get

$$\{\tilde{q}_{ij}(1)\}_{0 \leq i, j \leq 1} = \begin{pmatrix} 1/2 + \varepsilon & \\ & 1/2 - \varepsilon \end{pmatrix}, \quad \{b_{ij}(1)\}_{0 \leq i, j \leq 1} = \varepsilon \begin{pmatrix} 2 & 0 \\ 0 & -1 \end{pmatrix},$$

which implies  $\rho(1) = 1$  and hence the monotonicity stability for  $k = 0$ .

By carrying out the matrix transferring process for increasing multistep number, we can have  $\rho(3) = \rho(4) = \rho(5) = 2 = \zeta$  and then conclude the strong(3) stability for  $k \leq 1$ .

**Table 1** Stability results for the LWDG(2,  $k$ ) methods

Parameters				$n_*$ : strong ( $n_*$ ) stability	
$\theta_{00}$	$\theta_{10}$	$\theta_{11}$	$k \geq 2$	$k = 1$	$k = 0$
+	+	+	Weak (4)	1	1
+	+	-		3	
-	+	+		3	
-	+	-		4	

**Table 2** Stability results for the LWDG(3,  $k$ ) method

Parameters					$n_*$ : strong( $n_*$ ) stability	
$\theta_{00}$	$\theta_{11}$	$\theta_{20}$	$\theta_{21}$	$\theta_{22}$	$k \geq 1$	$k = 0$
+	$\pm$	+	+	$\pm$	1	1
+	-	+	-	$\pm$	3	
-	-	+	$\pm$	$\pm$	3	
+	+	+	-	$\pm$	4	
-	+	+	$\pm$	$\pm$	4	

**Table 3** Stability results for the LWDG(4,  $k$ ) method

Parameters							$n_*$ : strong( $n_*$ ) stability		
$\theta_{00}$	$\theta_{11}$	$\theta_{22}$	$\theta_{30}$	$\theta_{31}$	$\theta_{32}$	$\theta_{33}$	$k \geq 2$	$k = 1$	$k = 0$
+	+	-	+	+	+	$\pm$	1	1	1
+	$\pm$	+	+	+	+	$\pm$	2	1	
+	-	-	+	+	$\pm$	$\pm$	2	1	
+	+	$\pm$	+	+	-	$\pm$	3	1	
+	-	+	+	+	-	$\pm$	3	1	
+	-	$\pm$	+	-	+	$\pm$	5	3	
+	-	$\pm$	+	-	-	$\pm$	6	3	
-	-	$\pm$	+	+	$\pm$	$\pm$	6	3	
-	-	$\pm$	+	-	$\pm$	$\pm$	7	3	
+	+	$\pm$	+	-	$\pm$	$\pm$	7	3	
-	+	$\pm$	+	+	$\pm$	$\pm$	7	3	
-	+	$\pm$	+	-	$\pm$	$\pm$	8	4	

- The other cases can be studied similarly.

The stability results for the LWDG(2,  $k$ ) method are gathered in Table 1, where  $\pm$  stands for  $1/2 \pm \varepsilon$  here and below.

For  $r = 3$  or  $r = 4$ , we are able to similarly find  $\zeta = r - 1$  and that the central objective is negative. Hence we can claim the strong stability for  $k \geq 0$ , due to Theorem 3.2. The detailed results are collected in Tables 2 and 3.

For  $r = 5$ , we can get  $\zeta = 3$  and the positive central objective, which implies the strong stability for  $k \leq 2$  due to Theorem 3.3 and the weak(6) stability for  $k \geq 3$  due to Theorem 3.1. The detailed results are collected in Table 4.

**Remark 3.9** In the above four tables, the first row gives the numerical flux parameters to have the monotonicity stability for some  $k$ . If  $r \neq 4$ , it is acceptable to take  $\theta_{\ell\kappa} \equiv 1/2 + \varepsilon$  for any  $\ell$  and  $\kappa$ . Otherwise, for  $r = 4$ , we have to take all  $\theta_{\ell\kappa} \equiv 1/2 + \varepsilon$  except  $\theta_{22} = 1/2 - \varepsilon$ .

**Table 4** Stability results for the LWDG(5,  $k$ ) method

Parameters									$k \geq 3$	$n_*$ : strong( $n_*$ ) stability		
$\theta_{00}$	$\theta_{11}$	$\theta_{22}$	$\theta_{33}$	$\theta_{40}$	$\theta_{41}$	$\theta_{42}$	$\theta_{43}$	$\theta_{44}$		$k = 2$	$k = 1$	$k = 0$
+	+	±	±	+	+	+	±	±	Weak(6)	1	1	1
+	-	-	±	+	+	±	±	±		2	1	
+	-	+	±	+	+	+	±	±		2	1	
+	-	+	±	+	+	-	±	±		3	1	
+	+	±	±	+	+	-	±	±		3	1	
+	-	±	±	+	-	+	±	±		5	3	
+	-	±	±	+	-	-	±	±		6	3	
-	-	±	±	+	+	±	±	±		6	3	
+	±	±	±	+	-	±	±	±		7	3	
-	+	±	±	+	+	±	±	±		7	3	
-	+	±	±	+	-	±	±	±		8	4	

**Remark 3.10** Associated with the second row in Table 1 with  $\varepsilon = 1/2$ , we get the LWDG(2, 1) method with  $\theta_{00} = \theta_{10} = 1$  and  $\theta_{11} = 0$ . This LWDG scheme has been proved in [23] to have the stability result ( $u^{n,1} = -\tau p^n$ )

$$\|u^n\|_{L^2(I)}^2 + \|u^{n,1}\|_{L^2(I)}^2 \leq \|u_0\|_{L^2(I)}^2 + \|u^{0,1}\|_{L^2(I)}^2.$$

This result can not yield the strong(3) stability  $\|u^n\|_{L^2(I)} \leq \|u^0\|_{L^2(I)}$  for  $n \geq 3$ , as claimed in this paper.

So does for the LWDG(3, $k$ ) method [23] when the numerical flux parameters are taken from the second row of Table 2 with  $\varepsilon = 1/2$ .

### 4 Optimal Error Estimate

In this section we are devoted to obtaining the optimal  $L^2$ -norm error estimate for the ESTDG method with the stage-dependent numerical flux parameters. This result is stated in the following theorem.

**Theorem 4.1** For the ESTDG( $s, r, k$ ) method (2.6) with the averaged numerical flux parameter  $\Theta > 1/2$ , we have the optimal error estimate

$$\|u^N - U(t^N)\|_{L^2(I)} \leq C \|U_0\|_{H^{\vartheta+1}(I)} (h^{k+1} + \tau^r), \tag{4.1}$$

under the same type of temporal–spatial condition to ensure the  $L^2$ -norm stability, as stated in Theorems 3.1 through 3.3. Here  $\vartheta = \max(k + 1, r)$  and the bounding constant  $C > 0$  is independent of  $h, \tau$  and  $U_0$ .

This theorem has been proved for the RKDG method with the same numerical flux parameters, where the fourth order in time scheme is taken as an example [27]. Besides the stability analysis, the major techniques to prove this theorem are the standard GGR projection with a fixed parameter and the good definition of the reference functions which are related to the local time marching of the exact solutions. However, this strategy does not work well for the ESTDG method with stage-dependent numerical flux parameters, because the GGR

projection with any fixed parameter can not simultaneously eliminate the projection error at boundary endpoints for all time stages. We have to find a new approach to prove this theorem.

### 4.1 Proof of Theorem 4.1

To address the difficulties resulted from the stage-dependent numerical flux parameters, we would like in this paper to propose a new tool, named as *a series of space–time approximation functions* for any given spatial function. They will be used to set up a group of good reference functions and delicately define the stage errors for the fully discrete scheme. Different to the traditional analysis technique, they depend on not only the numerical method but also the considered PDE.

**Definition 4.1** Let  $W(x) \in L^2(I)$  be a given periodic function. Associated with the ESTDG( $s, r, k$ ) method of the time step  $\tau > 0$  and the finite element space  $V_h$ , a series of space–time approximation functions, denoted by

$$W_h^\ell = \mathbb{Q}_{h,\tau}^\ell W(x) \in V_h, \quad \ell = 0, 1, \dots, s, \tag{4.2}$$

are defined by the following conditions:

- Matching the local structure of the fully discrete numerical scheme, namely

$$\left( W_h^{\ell+1}, v \right)_{I_h} = \sum_{0 \leq \kappa \leq \ell} \left[ c_{\ell\kappa} \left( W_h^\kappa, v \right)_{I_h} + \tau d_{\ell\kappa} \mathcal{H}^{\theta_{\ell\kappa}} \left( W_h^\kappa, v \right) \right], \quad \forall v \in V_h, \tag{4.3a}$$

holds for  $0 \leq \ell \leq s - 1$ ;

- Preserving the balance of exact evolution under the control of PDE (1.1), namely

$$\left( W_h^s - W_h^0, v \right)_{I_h} = \left( W(x - \tau\beta) - W(x), v \right)_{I_h}, \quad \forall v \in V_h^*. \tag{4.3b}$$

Here  $V_h^* = \{v \in V_h : (v, 1)_{I_h} = 0\}$  is an orthogonal complementary space.

- Conserving the overall mean for the head function  $W_h^0$ , namely

$$\left( W_h^0, 1 \right)_{I_h} = \left( W(x), 1 \right)_{I_h}. \tag{4.3c}$$

**Remark 4.1** The head function is the most important one in the definition. For convenience of statement,  $W_h^s$  is called the tail function.

In what follows we give some comments to this definition. First of all, we point out that condition (4.3a) can be well understood by making full use of those concepts proposed in the matrix transferring process, for instance, the temporal differences of stage solutions and the evolution identity. That is to say, there holds

$$W_h^s = \sum_{0 \leq \ell \leq s} \alpha_\ell \mathbb{D}_\ell W_h \quad \text{with} \quad \mathbb{D}_\ell W_h = \sum_{0 \leq \kappa \leq \ell} \sigma_{\ell\kappa} W_h^\kappa, \tag{4.4}$$

where  $\alpha_\ell = \alpha_\ell(1)$  and  $\sigma_{\ell\kappa} = \sigma_{\ell\kappa}(1)$  have been defined in (3.11) and (3.4), respectively. Analogously, there also holds for  $0 \leq \ell \leq s - 1$  that

$$\left( \mathbb{D}_{\ell+1} W_h, v \right)_{I_h} = \tau \mathcal{H}^\vartheta \left( \mathbb{D}_\ell W_h, v \right) - \tau\beta \sum_{0 \leq \kappa \leq \ell} q_{\ell,\kappa}(\vartheta) \left( \mathbb{D}_\kappa W_h, v \right)_{I_h}, \tag{4.5}$$

for any  $v \in V_h$ , where  $q_{\ell,\kappa}(\vartheta) = q_{\ell,\kappa}(1; \vartheta)$  has been defined in (3.13).

As for the second condition (4.3b), we point out that it can be extended to the whole finite element space, i.e.,

$$\left(W_h^s - W_h^0, v\right)_{I_h} = \left(W(x - \tau\beta) - W(x), v\right)_{I_h}, \quad \forall v \in V_h. \tag{4.6}$$

This conclusion holds owing to the following facts:

- Since  $\mathcal{H}^\vartheta(\mathbb{D}_\ell W_h, 1) = 0$ , by taking  $v = 1$  in (4.5) we can inductively derive that

$$\left(\mathbb{D}_\ell W_h, 1\right)_{I_h} = 0, \quad \ell \geq 1. \tag{4.7}$$

Together with (4.4), this equality yields  $(W_h^s - W_h^0, 1)_{I_h} = 0$ .

- Since  $W(x)$  is periodic, it is easy to see  $(W(x - \tau\beta) - W(x), 1)_{I_h} = 0$ .

In other words, condition (4.3c) is only used to ensure the uniqueness if the space–time approximation functions are made up of (4.3a) and (4.6).

It is worthy to emphasize that any space–time approximation function given in Definition 4.1 is not a projection of  $W(x)$ , even when the numerical flux parameters are the same. An example is given below. Let  $I_h$  be a uniform mesh with the mesh size  $h$ , and consider the function  $W(x) \in V_h$ , which in each cell is defined by

$$W(x) = L_{j,1}(x) = \frac{1}{h} \left(2x - x_{j-\frac{1}{2}} - x_{j+\frac{1}{2}}\right), \quad x \in I_j.$$

Associated with the classical second order RKDG method (refer to [29] for details) with  $\theta_{\ell\kappa} \equiv 1$ , we can yield (with  $\lambda = |\beta|\tau/h$ )

$$W_h^0 = \frac{\lambda - 1}{3\lambda - 1} W, \quad W_h^1 = (1 - 6\lambda)W_h^0, \quad W_h^2 = (1 - 6\lambda + 18\lambda^2)W_h^0,$$

which are all not equal to  $W$ . This distinct property will cause many difficulties in obtaining the next lemma with respect to the approximation property.

**Lemma 4.1** *For sufficiently small  $\lambda = |\beta|\tau/h$ , a series of space–time approximation functions (4.2) are well defined, and further, if  $W(x) \in H^{\max(k+1, r+1)}(I)$ , the head function  $W_h^0$  satisfies the optimal error estimate*

$$\|W_h^0 - W\|_{L^2(I)} \leq C \left[ h^{k+1} \|W\|_{H^k(I)} + \tau^r \|W\|_{H^r(I)} \right], \tag{4.8}$$

where the bounding constant  $C > 0$  is independent of  $h, \tau$  and  $W$ . Note that the notation  $\natural = \max(k + 1, r)$  has been given in Theorem 4.1.

For ease of reading and understanding, we postpone the lengthy and technical proof of this lemma to the next subsection and come back to prove Theorem 4.1 now. For any  $n \geq 0$ , we utilize Definition 4.1 to define a series of space–time approximation functions with respect to  $U^n(x) = U(x, t^n)$ , namely

$$\chi^{n,\ell} = \mathbb{Q}_{h,\tau}^\ell U^n(x) \in V_h, \quad \ell = 0, 1, \dots, s. \tag{4.9}$$

It is worthy to emphasize that  $\chi^{n+1,0} \neq \chi^{n,s}$  in general, and the accumulation of these gaps at all time level forms the main error of the ESTDG method.

The reference functions are then defined by those functions in (4.9) except  $\ell = s$ . As a result, for any  $n \geq 0$  we denote the stage errors in the finite element space by

$$\xi^{n,\ell} = u^{n,\ell} - \chi^{n,\ell}, \quad \ell = 0, 1, \dots, s - 1. \tag{4.10a}$$

As used in the definition of the ESTDG method, we would like to give a supplementary notation

$$\xi^{n,s} = \xi^{n+1,0} = \xi^{n+1}. \tag{4.10b}$$

Every function  $\chi^{n,\ell}$  in (4.9) satisfies the variation form (4.3a) with  $W_h^\ell = \chi^{n,\ell}$ . Subtracting them from the fully discrete method with the same  $n$  and  $\ell$ , we can obtain a series of error equations as following: for  $\ell = 0, 1, \dots, s - 1$ , there holds

$$\left(\xi^{n,\ell+1}, v\right)_{I_h} = \sum_{0 \leq \kappa \leq \ell} \left[ c_{\ell\kappa} \left(\xi^{n,\kappa}, v\right)_{I_h} + \tau d_{\ell\kappa} \mathcal{H}^{\theta_{\ell\kappa}} \left(\xi^{n,\kappa}, v\right) \right] + \tau \left(F^{n,\ell}, v\right)_{I_h},$$

for any  $v \in V_h$ , where the source term  $F^{n,\ell}$  is equal to zero except the last one

$$F^{n,s-1} = \frac{1}{\tau} \left(\chi^{n,s} - \chi^{n+1,0}\right). \tag{4.11}$$

These error equations have the same form as the nonhomogeneous ESTDG method. Along the similar analysis line as in Sect. 3, we can get

$$\|\xi^N\|_{L^2(I)}^2 \leq C \left[ \|\xi^0\|_{L^2(I)}^2 + \sum_{0 \leq n < N} \|F^{n,s-1}\|_{L^2(I)}^2 \tau \right], \tag{4.12}$$

under the same type of temporal–spatial condition as stated in Theorems 3.1 through 3.3, where the bounding constant  $C > 0$  is independent of  $h$  and  $\tau$ , but may depend on the final time  $T$ .

Below we estimate each term on the right hand side of (4.12). It follows from the initial setting that  $\xi^0 = \mathbb{P}_h U_0 - \mathbb{Q}_{h,\tau}^0 U_0$ . By using the triangle inequality, we have

$$\begin{aligned} \|\xi^0\|_{L^2(I)} &\leq \|U_0 - \mathbb{P}_h U_0\|_{L^2(I)} + \|U_0 - \mathbb{Q}_{h,\tau}^0 U_0\|_{L^2(I)} \\ &\leq C \left[ h^{k+1} \|U_0\|_{H^{\natural}(I)} + \tau^r \|U_0\|_{H^r(I)} \right], \end{aligned} \tag{4.13}$$

where the well-known approximation property of  $\mathbb{P}_h$  ( $L^2$  projection) and Lemma 4.1 are used separately. Since the time step is uniform, definition (4.3) implies that

$$\chi^{n+1,0} - \chi^{n,0} = \mathbb{Q}_{h,\tau}^0 (U^{n+1} - U^n). \tag{4.14}$$

It follows from (4.6) that  $(\chi^{n,s} - \chi^{n,0}, v)_{I_h} = (U^{n+1} - U^n, v)_{I_h}$ . Hence (4.11) implies

$$\left(F^{n,s-1}, v\right)_{I_h} = \left(\frac{U^{n+1} - U^n}{\tau}, v\right)_{I_h} - \left(\mathbb{Q}_{h,\tau}^0 \left(\frac{U^{n+1} - U^n}{\tau}\right), v\right)_{I_h},$$

which, together with Lemma 4.1 again, yields

$$\|F^{n,s-1}\|_{L^2(I)} \leq C \left[ h^{k+1} \left\| \frac{U^{n+1} - U^n}{\tau} \right\|_{H^{\natural}(I)} + \tau^r \left\| \frac{U^{n+1} - U^n}{\tau} \right\|_{H^r(I)} \right].$$

Since  $U(x, t) = U_0(x - \beta t)$  and  $U^{n+1} - U^n = \int_n^{n+1} U_t(x, t') dt'$ , we can obtain from the above inequality that

$$\|F^{n,s-1}\|_{L^2(I)} \leq C \left[ h^{k+1} \|U_0\|_{H^{\natural+1}(I)} + \tau^r \|U_0\|_{H^{r+1}(I)} \right]. \tag{4.15}$$

Since  $\natural = \max(k + 1, r)$ , now we yield  $\|\xi^N\|_{L^2(I)} \leq C(h^{k+1} + \tau^r) \|U_0\|_{H^{\natural+1}(I)}$  by substituting (4.13) and (4.15) into (4.12). It follows from Lemma 4.1 that

$$\|U^N - \chi^{N,0}\|_{L^2(I)} \leq C(h^{k+1} + \tau^r) \|U_0\|_{H^{\natural}(I)}.$$



Since  $u^N - U^N = \xi^N - (U^N - \chi^{N,0})$ , the above two inequalities and the triangle inequality complete the proof of Theorem 4.1.

**Remark 4.2** Due to (4.13), the initial solution is admitted to be any function satisfying  $\|U_0 - u^0\|_{L^2(I)} \leq C \|U_0\|_{H^1(I)} h^{k+1}$ .

**Remark 4.3** In this paper the proof of inequality (4.15) strongly depends on the property (4.14), which only holds for the uniform time step. Till now we have not found a good way to rigorously prove this inequality for nonuniform time step size. How to address this difficulty is left for our further work.

### 4.2 Proof of Lemma 4.1

In this subsection we want to prove Lemma 4.1. Since the total number of the restrictions proposed by Definition 4.1 is equal to the unknowns' degrees of freedom, it is sufficient and necessary to prove the uniqueness and existence by verifying that there is only one trivial solution  $W_h^0 = \dots = W_h^s = 0$  for  $W = 0$ . The proof line of this topic is almost the same as that for the optimal estimate, so we would like to solely present the proof of (4.8) in this subsection.

In the next analysis, we will use the GGR projection and the flux lifting function for any given parameter  $\vartheta \neq 1/2$ . For convenience, we first give the definitions for  $k \geq 1$  and then a remark for  $k = 0$  later on.

**Definition 4.2** Let  $w$  be a periodic function belonging to  $H^1(I_j)$  for  $j = 1, 2, \dots, J$ . The GGR projection, denoted by  $\mathbb{G}_\vartheta w$ , is defined as the unique function in  $V_h$  such that for  $j = 1, 2, \dots, J$ ,

$$\int_{I_j} (\mathbb{G}_\vartheta w) v dx = \int_{I_j} w v dx \quad \forall v \in \mathcal{P}^{k-1}(I_j), \quad \text{and} \quad \{\mathbb{G}_\vartheta w\}_{j+\frac{1}{2}}^\vartheta = \{w\}_{j+\frac{1}{2}}. \quad (4.16)$$

**Definition 4.3** Let  $w^b$  be a single-valued periodic function defined on all element endpoints. The flux lifting function, denoted by  $\mathbb{L}_\vartheta w^b$ , is defined as the unique function in  $V_h$  such that for  $j = 1, 2, \dots, J$ ,

$$\int_{I_j} (\mathbb{L}_\vartheta w^b) v dx = 0 \quad \forall v \in \mathcal{P}^{k-1}(I_j), \quad \text{and} \quad \{\mathbb{L}_\vartheta w^b\}_{j+\frac{1}{2}}^\vartheta = w^b_{j+\frac{1}{2}}. \quad (4.17)$$

It has been proved in [3, Lemma 3.2] that the GGR projection is well-defined and the projection error  $\mathbb{G}_\vartheta^\perp w = w - \mathbb{G}_\vartheta w$  satisfies

$$\|\mathbb{G}_\vartheta^\perp w\|_{L^2(I)} + h^{\frac{1}{2}} \|(\mathbb{G}_\vartheta^\perp w)^\pm\|_{L^2(I_h)} \leq Ch^{\min(\aleph, k+1)} \|w\|_{H^\aleph(I)}, \quad (4.18)$$

where  $\aleph \geq 1$  is the smoothness requirement. Actually, the proof therein has implicitly used  $\mathbb{G}_\vartheta w = \mathbb{P}_h w + \mathbb{L}_\vartheta \{w - \mathbb{P}_h w\}^\vartheta$  and has shown that the flux lifting function is well-defined and satisfies

$$\|\mathbb{L}_\vartheta w^b\|_{L^2(I)} \leq Ch^{\frac{1}{2}} \|w^b\|_{L^2(I_h)}. \quad (4.19)$$

Furthermore, a direct application of Definitions 4.2 and 4.3 yields for any  $v \in V_h$ ,

$$\mathcal{H}^\vartheta(\mathbb{G}_\vartheta^\perp w, v) = 0 \quad \text{and} \quad \mathcal{H}^\vartheta(\mathbb{L}_\vartheta w^b, v) = \beta \left\langle w^b, \llbracket v \rrbracket \right\rangle_{I_h}, \quad (4.20)$$

as well as the property on the overall mean

$$\left(\mathbb{G}_{\vartheta}^{\perp} w, 1\right)_{I_h} = 0 \quad \text{and} \quad \left(\mathbb{L}_{\vartheta} w^b, 1\right)_{I_h} = 0. \tag{4.21}$$

**Remark 4.4** The above two definitions can be extended to  $k = 0$  with some minor modifications. The process is divided into two steps:

- Define two piecewise constant functions by the second condition in (4.16) and (4.17), respectively.
- Respectively subtract a constant to get two modified functions such that (4.21) holds.

It is easy to verify that the conclusions from (4.18) to (4.20) also hold.

Since  $r \leq s$ , we would like to adopt the *cutting-off* technique [26, 27] and define a series of functions

$$\partial_{\ell} W = \begin{cases} (-\tau\beta\partial_x)^{\ell} W, & 0 \leq \ell \leq r - 1, \\ 0, & r \leq \ell \leq s. \end{cases} \tag{4.22}$$

By this treatment, the smoothness assumption can be well controlled and we can get rid of the effect of the stage number.

Since  $W(x) \in H^{r+1}(I)$ , we know every  $\partial_{\ell} W \in H^2(I)$  and hence is continuous everywhere by an application of the Sobolev embedding theorem. Using integration by parts, after some manipulations we can get the consistency property

$$\tau \mathcal{H}^{\vartheta}(\partial_{\ell} W, v) = \left( (-\tau\beta\partial_x)\partial_{\ell} W, v \right)_{I_h}, \quad \forall v \in V_h. \tag{4.23}$$

Furthermore, the approximation property (4.18) with  $\aleph = \max(k+1-\ell, 1)$  and the definition (4.22) show

$$\|\mathbb{G}_{\vartheta}^{\perp}(\partial_{\ell} W)\|_{L^2(I)} + h^{\frac{1}{2}}\|(\mathbb{G}_{\vartheta}^{\perp}(\partial_{\ell} W))^{\pm}\|_{L^2(I_h)} \leq Ch^{k+1}\|W\|_{H^2(I)}, \tag{4.24}$$

no matter whether  $k + 1 \geq r$  or not. Here and below we assume  $\lambda \leq 1$  without losing generality, since it is small enough.

For  $0 \leq \ell \leq s$ , we define a series of error function in the finite element space

$$\mathcal{E}_{\ell}^{\vartheta} = \mathbb{D}_{\ell} W_h - \mathbb{G}_{\vartheta}(\partial_{\ell} W) \in V_h, \tag{4.25}$$

which leads to the decomposition  $\mathbb{D}_{\ell} W_h - \partial_{\ell} W = \mathcal{E}_{\ell}^{\vartheta} - \mathbb{G}_{\vartheta}^{\perp}(\partial_{\ell} W)$  as usual. Due to the triangle inequality and (4.24), it is sufficient to obtain (4.8) by proving

$$\|\mathcal{E}_0^{\vartheta}\|_{L^2(I)} \leq C \left[ h^{k+1}\|W\|_{H^2(I)} + \tau^r\|W\|_{H^r(I)} \right], \tag{4.26}$$

with a special choice  $\vartheta = \Theta$ .

To do that, we have to set up two relationships among  $\|\mathcal{E}_{\ell}^{\vartheta}\|_{L^2(I)}$  for  $\ell = 0, 1, \dots, s$ , in the forward and reverse direction, respectively. For ease of reading, we would like to only state them in the following two lemmas and put aside their proofs in the next two small subsections.

**Lemma 4.2** *For any  $\vartheta \neq \frac{1}{2}$ , there exists a bounding constant  $C = C(\vartheta) > 0$  such that for  $0 \leq \ell \leq s - 1$  there holds*

$$\|\mathcal{E}_{\ell+1}^{\vartheta}\|_{L^2(I)} \leq C\lambda\|\mathcal{E}_0^{\vartheta}\|_{L^2(I)} + C \left[ h^{k+1}\|W\|_{H^2(I)} + \tau^r\|W\|_{H^r(I)} \right]. \tag{4.27}$$

**Lemma 4.3** For  $\vartheta = \Theta$ , there exists a bounding constant  $C > 0$  such that

$$\|\mathcal{E}_0^\vartheta\|_{L^2(I)} \leq C \sum_{1 \leq \ell \leq s-1} \|\mathcal{E}_\ell^\vartheta\|_{L^2(I)} + C \left[ h^{k+1} \|W\|_{H^2(I)} + \tau^r \|W\|_{H^r(I)} \right]. \tag{4.28}$$

Till now (4.26) is implied by collecting Lemmas 4.2 and 4.3 if  $\lambda$  is small enough. This completes the proof of Lemma 4.1 and ends this subsection.

**4.2.1 Proof of Lemma 4.2**

We can prove this lemma by (4.5), which is equivalent to condition (4.3a). By adding and subtracting some terms involving  $\mathbb{G}_\vartheta(\partial_i W)$  three times, we have

$$\left( \mathcal{E}_{\ell+1}^\vartheta, v \right)_{I_h} = \mathcal{I}_1(v) + \mathcal{I}_2(v) + \mathcal{I}_3(v),$$

where

$$\begin{aligned} \mathcal{I}_1(v) &= \tau \mathcal{H}^\vartheta(\mathcal{E}_\ell^\vartheta, v) - \tau \beta \sum_{0 \leq \kappa \leq \ell} q_{\ell, \kappa}(\vartheta) \left\langle \llbracket \mathcal{E}_\kappa^\vartheta \rrbracket, \llbracket v \rrbracket \right\rangle_{I_h}, \\ \mathcal{I}_2(v) &= \tau \mathcal{H}^\vartheta(\mathbb{G}_\vartheta(\partial_\ell W), v) - \left\langle \mathbb{G}_\vartheta(\partial_{\ell+1} W), v \right\rangle_{I_h}, \\ \mathcal{I}_3(v) &= -\tau \beta \sum_{0 \leq \kappa \leq \ell} q_{\ell, \kappa}(\vartheta) \left\langle \llbracket \mathbb{G}_\vartheta(\partial_\kappa W) \rrbracket, \llbracket v \rrbracket \right\rangle_{I_h}. \end{aligned}$$

In what follows we estimate the above terms one by one.

Using (2.5c) for the first term, and using the Cauchy–Schwartz inequality and the inverse inequality (2.2) for the second term, we have

$$\mathcal{I}_1(v) \leq C \lambda \sum_{0 \leq \kappa \leq \ell} \|\mathcal{E}_\kappa^\vartheta\|_{L^2(I)} \|v\|_{L^2(I)}. \tag{4.29}$$

Due to (4.20) and (4.23), it follows from definition (4.22) that

$$\begin{aligned} \mathcal{I}_2(v) &= \left( -\tau \beta \partial_x(\partial_\ell W) - \mathbb{G}_\vartheta(\partial_{\ell+1} W), v \right)_{I_h} \\ &= \begin{cases} \left( \mathbb{G}_\vartheta^\perp(\partial_{\ell+1} W), v \right)_{I_h}, & 0 \leq \ell \leq r - 2, \\ \left( -\tau \beta \partial_x(\partial_\ell W), v \right)_{I_h}, & \ell = r - 1, \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

Using (4.24) for the first case and (4.22) for the second case, respectively, an application of Cauchy–Schwartz inequality yields a unified inequality

$$\mathcal{I}_2(v) \leq C \left[ h^{k+1} \|W\|_{H^2(I)} + \tau^r \|W\|_{H^r(I)} \right] \|v\|_{L^2(I)}. \tag{4.30}$$

Since  $\llbracket \partial_\kappa W \rrbracket = 0$  and  $\lambda \leq 1$ , we can use (4.24) and (2.2) to get

$$\mathcal{I}_3(v) = \tau \beta \sum_{0 \leq \kappa \leq \ell} q_{\ell, \kappa}(\vartheta) \left\langle \llbracket \mathbb{G}_\vartheta^\perp(\partial_\kappa W) \rrbracket, \llbracket v \rrbracket \right\rangle_{I_h} \leq C h^{k+1} \|W\|_{H^2(I)} \|v\|_{L^2(I)}. \tag{4.31}$$

Summing up the above three conclusions and taking  $v = \mathcal{E}_{\ell+1}^\vartheta \in V_h$ , we finally obtain

$$\|\mathcal{E}_{\ell+1}^\vartheta\|_{L^2(I)} \leq C \lambda \sum_{0 \leq \kappa \leq \ell} \|\mathcal{E}_\kappa^\vartheta\|_{L^2(I)} + C \left[ h^{k+1} \|W\|_{H^2(I)} + \tau^r \|W\|_{H^r(I)} \right],$$

for  $0 \leq \ell \leq s - 1$ . This completes the proof of Lemma 4.2.

### 4.2.2 Proof of Lemma 4.3

We can prove this lemma by (4.6), which is mainly related to condition (4.3b). Substitute (4.4) into the left hand side (LHS) of this condition and expand each term by the relationship (4.5). By changing the summation orders for those terms involving  $q_{\kappa,\ell}(\vartheta)$ , we can easily get

$$\begin{aligned} \text{LHS} &= \tau \sum_{0 \leq \ell \leq s-1} \alpha_{\ell+1} \mathcal{H}^\vartheta(\mathbb{D}_\ell W_h, v) - \tau\beta \sum_{0 \leq \kappa \leq s-1} \psi_\kappa(\vartheta) \left\langle \llbracket \mathbb{D}_\kappa W_h \rrbracket, \llbracket v \rrbracket \right\rangle_{I_h} \\ &= \tau \mathcal{H}^\vartheta \left( \sum_{0 \leq \ell \leq s-1} \left[ \alpha_{\ell+1} \mathbb{D}_\ell W_h - \psi_\ell(\vartheta) \mathbb{L}_\vartheta \llbracket \mathbb{D}_\ell W_h \rrbracket \right], v \right), \end{aligned} \tag{4.32}$$

where the second identity in (4.20) is used at the last step, and

$$\psi_\ell(\vartheta) = \sum_{\ell \leq \kappa \leq s-1} \alpha_{\kappa+1} q_{\kappa,\ell}(\vartheta). \tag{4.33}$$

Next we consider the right hand side (RHS) of condition (4.6). An application of the Taylor expansion up to  $r$ th order derivative yields

$$W(x - \tau\beta) - W(x) = (-\tau\beta \partial_x) \left[ \sum_{0 \leq \ell \leq r-1} \frac{1}{(\ell + 1)!} \partial_\ell W(x) + \tilde{W}(x) \right], \tag{4.34}$$

with the truncation function

$$\tilde{W}(x) = \frac{1}{r!(\tau\beta)} \int_0^{\tau\beta} \partial_x^r W(x - \tilde{x})(\tilde{x} - \tau\beta)^r d\tilde{x}.$$

It is easy to see that  $(\tilde{W}, 1)_{I_h} = 0$  and

$$\|\tilde{W}\|_{L^2(I)} \leq C\tau^r \|W\|_{H^r(I)}. \tag{4.35}$$

By integration by part for the definition of  $\tilde{W}(x)$ , say,

$$\tilde{W}(x) = \frac{1}{r!(\tau\beta)} \left[ \int_0^{\tau\beta} \partial_x^{r-1} W(x - \tilde{x})(\tilde{x} - \tau\beta)^{r-1} r d\tilde{x} + \partial_x^{r-1} W(x)(-\tau\beta)^r \right],$$

the derivative order of  $W(\cdot)$  is dropped to be  $r - 1$ . With this new formula and noticing the relationship of integration and the norm in Hilbert space, we are able to get

$$\|\tilde{W}\|_{H^\sharp(I)} \leq C\tau^{r-1} \|W\|_{H^\sharp(I)}, \quad \text{with } \sharp = \max(k + 2 - r, 1). \tag{4.36}$$

Substituting (4.34) into RHS and using the consistency property (4.23) for every  $\partial_\ell W$  and  $\tilde{W}$ , we can obtain from the first identity in (4.20) that

$$\text{RHS} = \tau \mathcal{H}^\vartheta \left( \sum_{0 \leq \ell \leq s-1} \alpha_{\ell+1} \mathbb{G}_\vartheta(\partial_\ell W) + \mathbb{G}_\vartheta \tilde{W}, v \right), \tag{4.37}$$

where Lemma 3.1 has been used. Here the upper bound of summation index is raised from  $r - 1$  to  $s - 1$ , since  $\partial_\ell W = 0$  for  $\ell \geq r$ , due to (4.22).

Due to (4.32) and (4.37), it follows from condition (4.6) that

$$Q^\vartheta \stackrel{\text{def}}{=} \sum_{0 \leq \ell \leq s-1} \left[ \alpha_{\ell+1} \mathcal{E}_\ell^\vartheta - \psi_\ell(\vartheta) \mathbb{L}_\vartheta \llbracket \mathbb{D}_\ell W_h \rrbracket \right] - \mathbb{G}_\vartheta \tilde{W} \in V_h \tag{4.38}$$

satisfies the variational form  $\mathcal{H}^\vartheta(\varrho^\vartheta, v) = 0$  for any  $v \in V_h$ . By successively taking  $v = \varrho^\vartheta$  and  $v = \partial_x \varrho^\vartheta$  here, we can see that  $\varrho^\vartheta$  must be a constant. This concludes

$$\varrho^\vartheta = 0, \tag{4.39}$$

since the overall mean is equal to zero. In fact, it is trivial to verify  $(\varrho^\vartheta, 1)_{I_h} = 0$  as following:

- By (4.21), we have  $(\mathbb{L}_\vartheta \llbracket \mathbb{D}_\ell W_h \rrbracket, 1)_{I_h} = 0$  for  $\ell \geq 0$  and  $(\mathbb{G}_\vartheta \tilde{W}, 1)_{I_h} = (\tilde{W}, 1)_{I_h} = 0$ .
- Furthermore, we also have  $(\mathcal{E}_\ell^\vartheta, 1)_{I_h} = 0$  for different cases:
  - For  $\ell = 0$ , condition (4.3c) implies  $(W_h^0, 1)_{I_h} = (W, 1)_{I_h} = (\mathbb{G}_\vartheta W, 1)_{I_h}$ ;
  - Otherwise, for  $\ell \geq 1$ , the periodicity means  $(\mathbb{G}_\vartheta(\partial_\ell W), 1)_{I_h} = (\partial_\ell W, 1)_{I_h} = 0$ , and (4.7) shows  $(\mathbb{D}_\ell W_h, 1)_{I_h} = 0$ .

Lemma 3.5 with  $m = 1$  implies the main property

$$\psi_0(\Theta) = \sum_{0 \leq \kappa \leq s-1} \alpha_{\kappa+1} q_{\kappa,0}(\Theta) = 0. \tag{4.40}$$

Thanks to this property, we can get rid of the trouble term  $\mathbb{L}_\vartheta \llbracket \mathbb{D}_0 W_h \rrbracket$  in (4.38). At this moment it follows from (4.39) and  $\alpha_1 \neq 0$  (due to Lemma 3.1) that

$$\|\mathcal{E}_0^\vartheta\|_{L^2(I)} \leq C \sum_{1 \leq \ell \leq s-1} \left[ \|\mathcal{E}_\ell^\vartheta\|_{L^2(I)} + \|\mathbb{L}_\vartheta \llbracket \mathbb{D}_\ell W_h \rrbracket\|_{L^2(I)} \right] + C \|\mathbb{G}_\vartheta \tilde{W}\|_{L^2(I)}. \tag{4.41}$$

Here and below  $\vartheta$  is fixed to be  $\Theta$ .

Due to the continuity of  $\partial_\ell W$ , as mentioned after its definition, we have  $\llbracket \mathbb{D}_\ell W_h \rrbracket = \llbracket \mathbb{D}_\ell W_h - \partial_\ell W \rrbracket = \llbracket \mathcal{E}_\ell^\vartheta \rrbracket - \llbracket \mathbb{G}_\vartheta^\perp \partial_\ell W \rrbracket$ . Then it follows from (4.19) and the triangle inequality that

$$\|\mathbb{L}_\vartheta \llbracket \mathbb{D}_\ell W_h \rrbracket\|_{L^2(I)} \leq Ch^{\frac{1}{2}} \|\llbracket \mathcal{E}_\ell^\vartheta \rrbracket\|_{L^2(I_h)} + Ch^{\frac{1}{2}} \|\llbracket \mathbb{G}_\vartheta^\perp \partial_\ell W \rrbracket\|_{L^2(I_h)}.$$

Together with (2.2) and (4.24) for each term, this deduces

$$\|\mathbb{L}_\vartheta \llbracket \mathbb{D}_\ell W_h \rrbracket\|_{L^2(I)} \leq C \|\mathcal{E}_\ell^\vartheta\|_{L^2(I)} + Ch^{k+1} \|W\|_{H^r(I)}. \tag{4.42}$$

By the triangle inequality and (4.18), we have

$$\|\mathbb{G}_\vartheta \tilde{W}\|_{L^2(I)} \leq \|\tilde{W}\|_{L^2(I)} + \|\mathbb{G}_\vartheta^\perp \tilde{W}\|_{L^2(I)} \leq \|\tilde{W}\|_{L^2(I)} + Ch^\sharp \|\tilde{W}\|_{H^r(I)}.$$

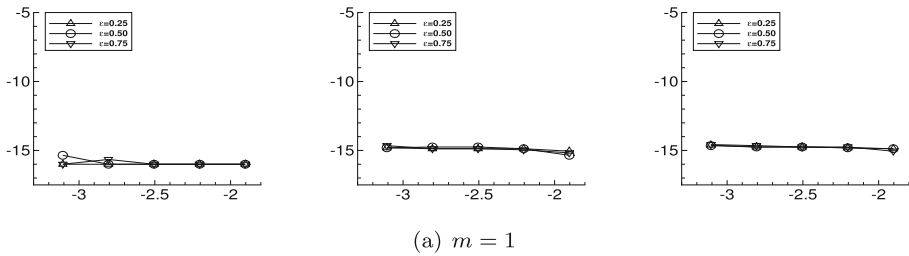
The two terms on the right hand side are bounded by (4.35) and (4.36), respectively. Since  $\lambda \leq 1$ , we can get the unified inequality

$$\|\mathbb{G}_\vartheta \tilde{W}\|_{L^2(I)} \leq C \left[ h^{k+1} \|W\|_{H^r(I)} + \tau^r \|W\|_{H^r(I)} \right]. \tag{4.43}$$

Substituting (4.42) and (4.43) into (4.41) completes the proof of Lemma 4.3.

### 5 Numerical Experiments

In this section we present some numerical experiments to verify the proposed theoretical results. Let  $\beta = 1$  in (1.1) for all tests. All schemes are taken from the two examples given in Sect. 3.



**Fig. 1** The  $L^2$ -norm amplification of the RKDG(4, 4,  $k$ ) solutions every  $m$ -step:  $k = 1, 2, 3$  from left to right. Here  $\varepsilon = 0.25, 0.50, 0.75, z = 1$  and  $y = 3$

### 5.1 Verification on Stability Results

Take the uniform meshes with  $J = 64$ , as an example. With standard orthogonal basis functions of the discontinuous finite element space, the ESTDG method can be written into  $\tilde{u}^{n+1} = \mathbb{K}\tilde{u}^n$ , where  $\tilde{u}^n$  is the solution vector made up of the expansion coefficients of  $u^n$ . The spectral norm  $\|\mathbb{K}^m\|_2$  describes the  $L^2$ -norm amplification every  $m$  step time marching [29].

#### 5.1.1 The RKDG Method

Consider the RKDG(4, 4,  $k$ ) method with the numerical flux parameters (3.38), where  $\varepsilon = 0.25, 0.50, 0.75$  and  $z = 1$ . For  $y = 3$  and  $y = 1$ , we respectively plot in Figs. 1 and 2 the quantity

$$\max(\|\mathbb{K}^m\|_2^2 - 1, 10^{-16}) \tag{5.1}$$

for different CFL number  $\lambda$  in the logarithmic coordinates, with  $k = 1, 2, 3$  from left to right.

- For  $y = 3$ , this quantity is always close to  $10^{-16}$  and thus implies the monotonicity stability.
- For  $y = 1$ , the data points increase along the line of slope 5 only for  $k \geq 2$  and  $m = 1$ . These numerical results show the strong(2) stability at least and the monotonicity stability for  $k \leq 1$ .

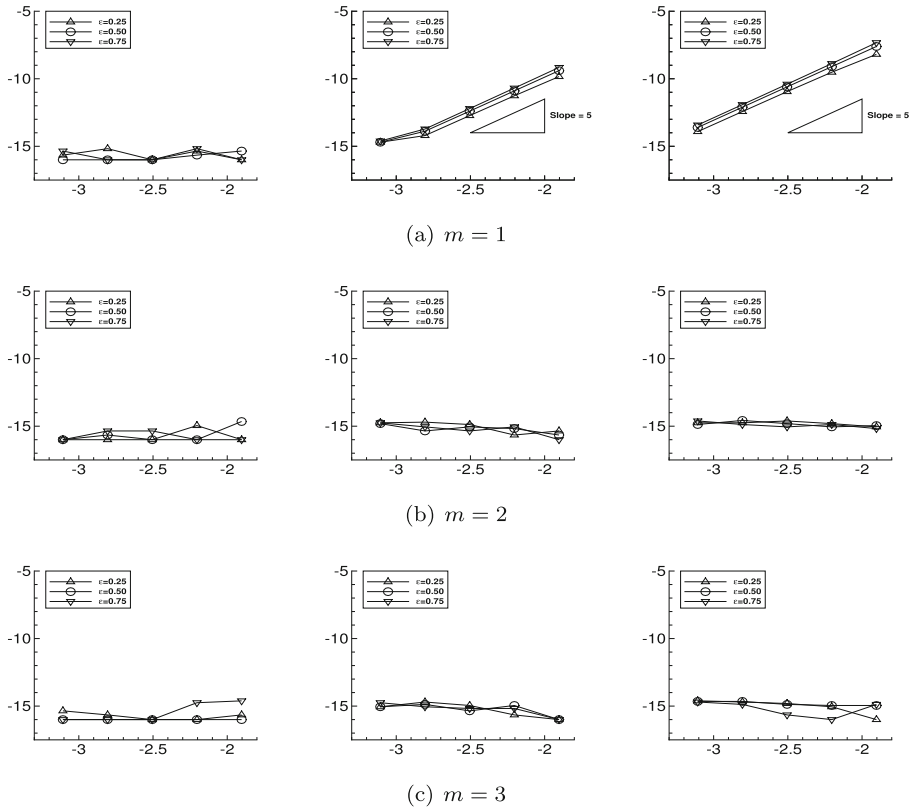
These observation verify what we have stated in Sect. 3.2.1.

To show the difference between the strong stability and the monotonicity stability, we take  $k = 3$  as an example and plot in Fig. 3 the  $L^2$ -norm evolution at the first twelve steps, where  $\lambda = 0.02$  and  $\varepsilon = 0.50$ . The initial solution is taken as the first unit singular vector of  $\mathbb{K}$ . For  $y = 1$ , we can see in the left picture that the  $L^2$ -norm overshoots at the first step and decreases every two and three steps. But for  $y = 3$ , the monotonicity stability is clearly observed in the right picture. This sharply verifies our theoretical results given in Sect. 3.2.1.

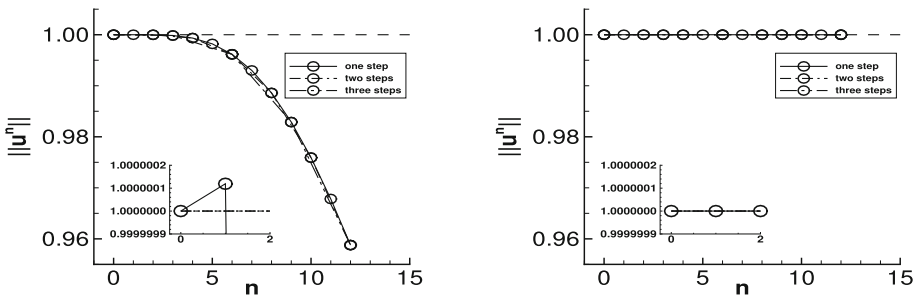
#### 5.1.2 The LWDG Method

Consider the LWDG(2,  $k$ ) method. As an example, the numerical flux parameters are defined as  $\theta_{00} = \theta_{10} = 1/2 + \varepsilon$  and  $\theta_{11} = 1/2 - \varepsilon$ , with  $\varepsilon = 0.25, 0.50, 0.75$ . We plot in Fig. 4 some pictures about the quantity (5.1) for  $k = 0, 1, 2$  and  $m = 1, 2, 3, 4, 5$ .

- If  $k = 0$ , this quantity is close to  $10^{-16}$  and shows the monotonicity stability.



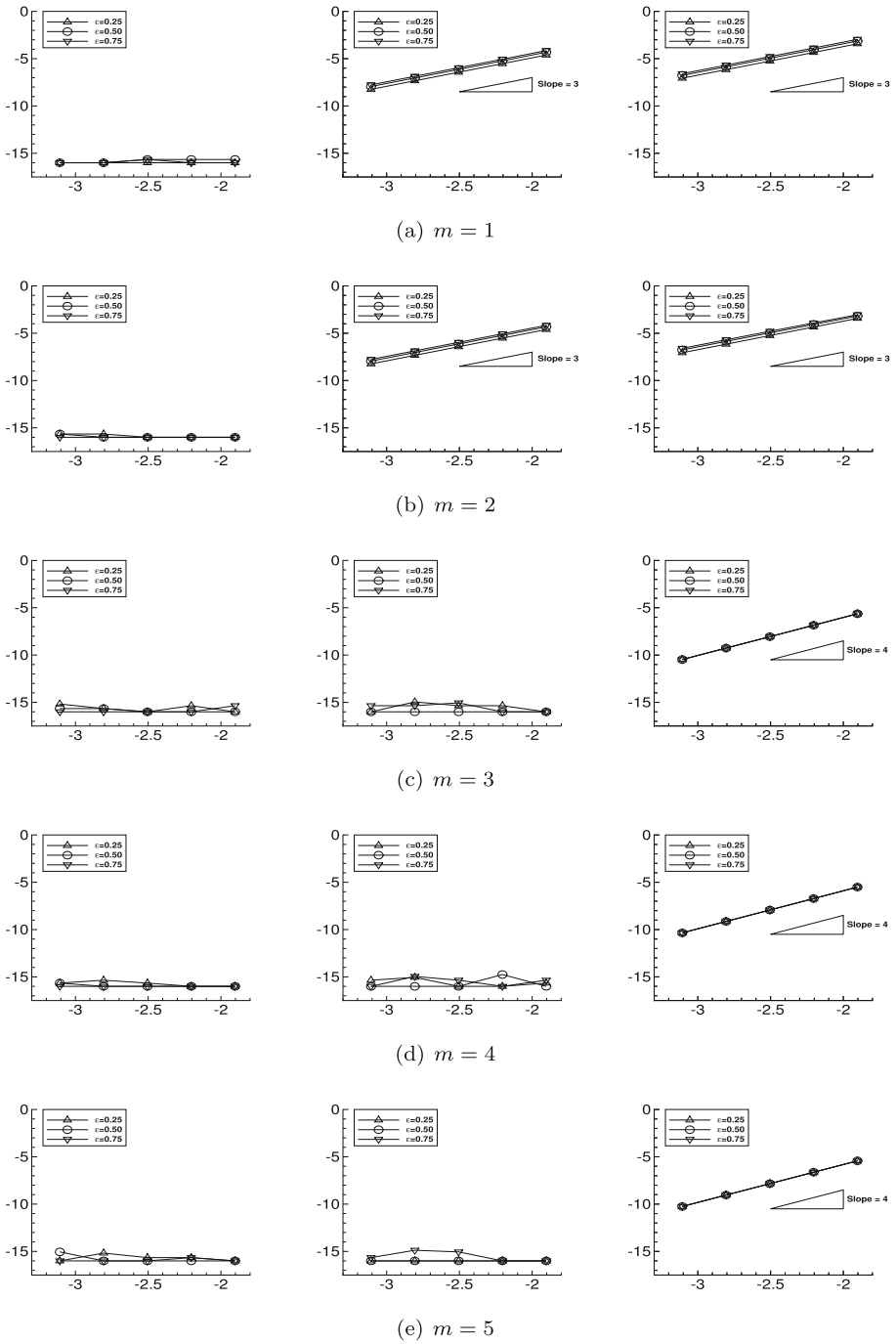
**Fig. 2** The  $L^2$ -norm amplification of the RKDG(4, 4,  $k$ ) solutions every  $m$ -step:  $k = 1, 2, 3$  from left to right. Here  $\varepsilon = 0.25, 0.50, 0.75, z = 1$  and  $y = 1$



**Fig. 3** The  $L^2$ -norm evolution for the RKDG(4, 4, 3) method. Left:  $y = 1$ ; Right:  $y = 3$ . Here  $z = 1, \lambda = 0.02$  and  $\varepsilon = 0.50$

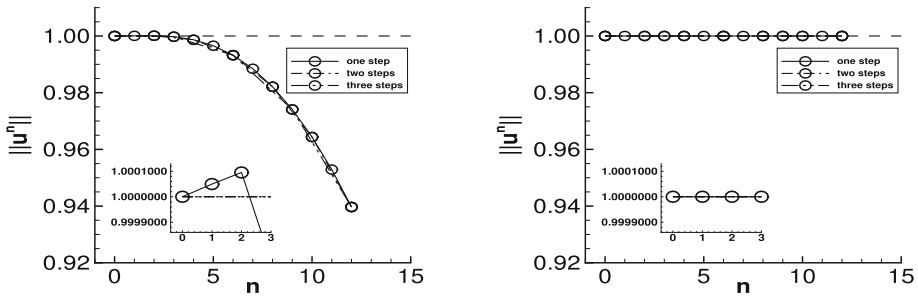
- If  $k = 1$ , the data points increase along the line of slope 3 for  $m \leq 2$  and are close to  $10^{-16}$  for  $m \geq 3$ . This verifies the strong(3) stability for  $k = 1$ .
- If  $k = 2$ , the data points increase with slope 3 (odd) for  $m \leq 2$  and with slope 4 (even) for  $m \geq 3$ . This shows the weak(4) stability.

The above observations well support the results listed in Table 1.



**Fig. 4** The  $L^2$ -norm amplification of the LWDG(2,  $k$ ) solution every  $m$ -step:  $\theta_{00} = \theta_{10} = 1/2 + \varepsilon$  and  $\theta_{11} = 1/2 - \varepsilon$ . Here  $k = 0, 1, 2$  from left to right and  $\varepsilon = 0.25, 0.50, 0.75$





**Fig. 5** The  $L^2$ -norm evolution for the LWDG(2, 1) method. Left:  $\theta_{11} = 0$ ; Right:  $\theta_{11} = 1$ . Here  $\lambda = 0.02$  and  $\theta_{00} = \theta_{10} = 1$

In Fig. 5, the left picture plots the  $L^2$ -norm evolution of the LWDG(2, 1) solution at the previous twelve steps, where  $\lambda = 0.02$  and  $\varepsilon = 0.50$ . The initial solution vector is taken as the first unit singular vector of  $\mathbb{K}^2$ . We can see that the monotonicity decreasing is lost at the first two steps and conclude that the scheme can not have the strong(2) stability. As a comparison, we also plot in the right picture for the LWDG(2,1) method with  $\theta_{11} = 1/2 + \varepsilon$  and the others are kept the same. We can see the monotonicity stability for this case, as we have predicted in theory.

**5.2 Verification on the Error Estimate**

In this subsection we investigate the numerical accuracy of the ESTDG method with two initial solutions. Since the numerical results are almost the same, we only present the experiment data for the RKDG(4,4,k) method on nonuniform mesh, which is constructed by perturbing the uniform mesh nodes randomly by at most 10%. Take the final time  $T = 1$ , and the time step  $\tau = 0.05h_{\min}$  in what follows, where  $h_{\min}$  is the minimal length.

First we take a sufficiently smooth initial solution, for example,

$$U_0(x) = \sin(2\pi x).$$

In Tables 5 and 6, we give the error and convergence order in the  $L^2$ -norm for  $y = 3$  and  $y = 1$  respectively. We can clearly observe the optimal convergence order, which supports the result in Theorem 4.1.

Next we investigate the smoothness requirement proposed in this paper. To do that, we take  $k = 3$  and the initial solution

$$U_0(x) = [\sin(2\pi x)]^{\epsilon+2/3},$$

and  $\epsilon$  is a positive integer. This function belongs to  $H^{\epsilon+1}(I)$ , but not  $H^{\epsilon+2}(I)$ . In Table 7, the optimal convergence order is clearly observed when  $\epsilon = r$ , but not  $\epsilon = r - 1$ . This indicates that the regularity requirement in Theorem 4.1 appears to be sharp.

**5.3 Discussions on  $\Theta = 1/2$**

In this subsection we give some discussions on the stability performance and the convergence order when  $\Theta = 1/2$ . As an example, we consider the standard RKDG(3, 3, k) method [32] with numerical flux parameters  $\theta_{00}$ ,  $\theta_{11}$  and  $\theta_{22}$ , for which the average numerical flux

**Table 5** The  $L^2$ -norm errors and convergence orders of the RKDG(4, 4,  $k$ ) method with the numerical flux parameter (3.38) and  $\gamma = 3$

$k$	$J$	$\varepsilon = 0.25$		$\varepsilon = 0.50$		$\varepsilon = 0.75$	
		Error	Order	Error	Order	Error	Order
1	160	7.32E-05		5.28E-05		4.90E-05	
	320	1.83E-05	2.00	1.31E-05	2.01	1.24E-05	1.98
	640	4.56E-06	2.00	3.32E-06	1.99	3.09E-06	2.00
	1280	1.14E-06	2.00	8.27E-07	2.00	7.73E-07	2.00
	2560	2.85E-07	2.00	2.07E-07	2.00	1.93E-07	2.00
2	160	2.11E-07		3.42E-07		4.91E-07	
	320	2.67E-08	2.98	4.27E-08	3.00	6.17E-08	2.99
	640	3.34E-09	3.00	5.32E-09	3.01	7.68E-09	3.01
	1280	4.18E-10	3.00	6.66E-10	3.00	9.60E-10	3.00
	2560	5.23E-11	3.00	8.32E-11	3.00	1.20E-10	3.00
3	160	6.03E-10		4.88E-10		5.21E-10	
	320	3.71E-11	4.02	2.99E-11	4.03	2.96E-11	4.14
	640	2.31E-12	4.01	1.90E-12	3.98	1.83E-12	4.01
	1280	1.44E-13	4.00	1.16E-13	4.03	1.16E-13	3.98
	2560	8.95E-15	4.01	7.26E-15	4.00	7.34E-15	3.98

Nonuniform mesh

**Table 6** The  $L^2$ -norm errors and convergence orders of the RKDG(4, 4,  $k$ ) method with the numerical flux parameter (3.38) and  $\gamma = 1$

$k$	$J$	$\varepsilon = 0.25$		$\varepsilon = 0.50$		$\varepsilon = 0.75$	
		Error	Order	Error	Order	Error	Order
1	160	8.03E-05		5.55E-05		4.99E-05	
	320	2.01E-05	2.00	1.39E-05	2.00	1.24E-05	2.01
	640	5.01E-06	2.00	3.47E-06	2.00	3.13E-06	1.98
	1280	1.25E-06	2.00	8.67E-07	2.00	7.87E-07	1.99
	2560	3.13E-07	2.00	2.17E-07	2.00	1.97E-07	2.00
2	160	2.03E-07		4.83E-07		4.22E-07	
	320	2.49E-08	3.03	3.03E-08	3.99	5.31E-08	2.99
	640	3.13E-09	2.99	1.91E-09	3.99	6.64E-09	3.00
	1280	3.91E-10	3.00	1.18E-10	4.01	8.30E-10	3.00
	2560	4.94E-11	2.99	7.44E-11	3.99	1.04E-10	3.00
3	160	6.48E-10		4.83E-10		4.78E-10	
	320	3.92E-11	4.05	3.03E-11	3.99	2.91E-11	4.04
	640	2.49E-12	3.98	1.91E-12	3.99	1.86E-12	3.97
	1280	1.58E-13	3.98	1.18E-13	4.01	1.15E-13	4.02
	2560	9.78E-15	4.01	7.44E-15	3.99	7.23E-15	3.99

Nonuniform mesh

**Table 7** The  $L^2$ -norm errors and convergence orders of the RKDG(4, 4, 3) method on nonuniform mesh

$\epsilon$	$J$	RKDG(4, 4, 3), $y = 3$				RKDG(4, 4, 3), $y = 1$			
		Error	Order	Error	Order	Error	Order	Error	Order
0.25	160	3.87E-08		2.20E-08		3.70E-08		2.43E-08	
	320	3.06E-09	3.66	1.37E-09	4.00	2.92E-09	3.66	1.52E-09	4.00
	640	2.45E-10	3.64	8.57E-11	4.00	2.35E-10	3.64	9.51E-11	4.00
	1280	1.97E-11	3.64	5.35E-12	4.00	1.91E-11	3.62	5.94E-12	4.00
	2560	1.59E-12	3.63	3.34E-13	4.00	1.55E-12	3.62	3.71E-13	4.00
0.50	160	5.24E-08		1.66E-08		4.82E-08		1.74E-08	
	320	4.05E-09	3.69	1.02E-09	4.02	3.76E-09	3.68	1.08E-09	4.01
	640	3.12E-10	3.70	6.36E-11	4.01	2.93E-10	3.68	6.72E-11	4.01
	1280	2.41E-11	3.70	3.97E-12	4.00	2.28E-11	3.68	4.19E-12	4.00
	2560	1.85E-12	3.70	2.48E-13	4.00	1.77E-12	3.68	2.62E-13	4.00
0.75	160	6.45E-08		1.56E-08		5.82E-08		1.59E-08	
	320	4.82E-09	3.74	9.55E-10	4.03	4.43E-09	3.72	9.78E-10	4.03
	640	3.62E-10	3.74	5.91E-11	4.01	3.37E-10	3.72	6.07E-11	4.01
	1280	2.73E-11	3.73	3.68E-12	4.01	2.56E-11	3.71	3.78E-12	4.00
	2560	2.07E-12	3.72	2.30E-13	4.00	1.96E-12	3.71	2.36E-13	4.00

Here  $\epsilon = r - 1$  on the left column and  $\epsilon = r$  on the right column

**Table 8** The  $L^2$ -norm errors and convergence orders of the RKDG(3, 3,  $k$ ) method with different triplets  $(\theta_{00}, \theta_{11}, \theta_{22})$  satisfying  $\Theta = 1/2$ : regular nonuniform mesh

$k$	$J$	(0.5, 0.5, 0.5)		(0.52, 0.48, 0.5)		(1, 0, 0.5)	
		Error	Order	Error	Order	Error	Order
1	160	3.99E-03		3.31E-03		6.27E-04	
	320	1.98E-03	1.01	1.42E-03	1.22	1.59E-04	1.98
	640	9.91E-04	1.00	5.90E-04	1.26	4.00E-05	1.99
	1280	4.95E-04	1.00	2.61E-04	1.18	1.01E-05	1.99
	2560	2.48E-04	1.00	1.27E-04	1.04	2.50E-06	2.01
2	160	1.96E-06		1.84E-06		3.39E-07	
	320	4.70E-07	2.06	4.98E-07	1.89	4.28E-08	2.99
	640	1.16E-07	2.01	1.28E-07	1.96	5.40E-09	2.99
	1280	2.90E-08	2.00	3.19E-08	2.00	6.81E-10	2.99
	2560	7.25E-09	2.00	7.86E-09	2.02	8.39E-11	3.02

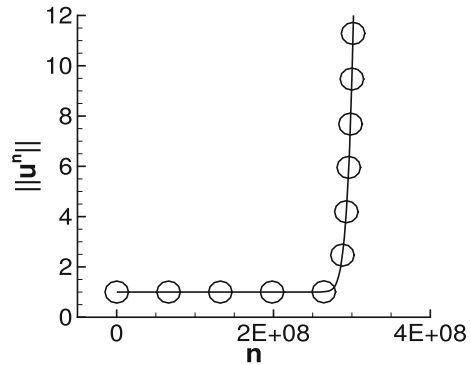
parameter is

$$\Theta = \frac{1}{6}(\theta_{00} + \theta_{11} + 4\theta_{22}).$$

In Table 8 we give three examples that the related schemes convergent with different convergence orders. In this test, we take the final time  $T = \pi$ , and the regular nonuniform mesh [14]

$$x_{j+\frac{1}{2}} = \begin{cases} j/J, & j \text{ is even,} \\ j/J + 0.1/J, & j \text{ is odd.} \end{cases} \tag{5.2}$$

**Fig. 6** The  $L^2$ -norm evolution of numerical solution of the RKDG(3, 3, 2) method with  $(\theta_{00}, \theta_{11}, \theta_{22}) = (1, 1, 0.25)$ . Here  $\lambda = 0.001$  and  $J = 64$



The time step is set as  $\tau = 0.1h_{\min}$ . Different to the same numerical flux parameter, the last parameter triplet  $(\theta_{00}, \theta_{11}, \theta_{22}) = (1, 0, 0.5)$  gives the optimal order even when the mesh is nonuniform for  $k = 1$  and  $k = 2$ .

For the above three schemes, the stability under the CFL condition is implied by the convergence. However, the stability conclusion is inconclusive when  $\Theta = 1/2$ . Below we give an example that the RKDG(3, 3, 2) scheme with  $(\theta_{00}, \theta_{11}, \theta_{22}) = (1, 1, 0.25)$  seems to be linearly unstable. To this end we take the uniform mesh with  $J = 64$ , and take the initial solution as the local  $L^2$  projection of  $\sqrt{2}\sin(8\pi x)$ . Even with a very small CFL number,  $\lambda = 0.001$ , one can clearly observe in Fig. 6 that the  $L^2$ -norm of numerical solution exponentially increases, which indicates a possible instability.

## 6 Conclusion

In this paper we have presented the  $L^2$ -norm stability analysis and optimal error estimate for the ESTDG method, which adopts the explicit single-step time-marching and the spatial DG discretization with stage-dependent numerical flux parameters. By a unified analysis framework, we successfully address many difficulties in theoretical analysis and then set up the detailed  $L^2$ -norm stability results for the RKDG method with downwind treatments and the LWDG method with different numerical flux parameters for the auxiliary variables. The main technique used in this paper is the matrix transferring process based on temporal differences of stage solutions, in order to achieve a good energy equation with some important indices to carry out the energy analysis. Motivated by the studies for the RKDG methods with the same numerical flux parameters, the averaged numerical flux parameter is proposed in this paper to measure the upwind effect in the fully discrete ESTDG method. In order to obtain the optimal error estimate for the ESTDG method with stage-dependent numerical flux parameters, we put forward a new proof framework by proposing a series of space-time approximation functions for any given spatial function. During this procedure, many techniques proposed in the matrix transferring process, together with the averaged numerical flux parameter, still play essential roles to establish the corresponding approximation property. In future work, we will extend the above works to variable-coefficient linear hyperbolic problems and nonlinear conservation laws in one and/or multidimensional cases.

## 7 Appendix

In this section we give some supplemental materials for those conclusions unproved in Sect. 3. This process involves many notations and manipulations of matrices.

To do that, we give some elemental notations here. Associated with the multistep number  $m$  and the stage number  $s$ , we introduce some column vectors and square matrices of size  $ms$ , whose component is either 0 or 1. More specifically, we denote  $\mathbf{1}(m, s) = (1, 1, \dots, 1)^\top$  and let  $\mathbf{e}_i(m, s)$ , for  $0 \leq i \leq ms - 1$ , be the unit vector which has 1 at the  $i$ -th position. Let  $\underline{\mathbf{I}}(m, s)$  be the identity matrix and  $\underline{\mathbf{E}}(m, s)$  be the shifting matrix which has 1 at the lower second diagonal line. Then we define

$$\underline{\mathbf{L}}(m, s) = \left[ \underline{\mathbf{I}}(m, s) - \underline{\mathbf{E}}(m, s) \right]^{-1} - \underline{\mathbf{I}}(m, s) = \sum_{1 \leq \kappa \leq ms-1} \underline{\mathbf{E}}(m, s)^\kappa, \quad (7.1)$$

which has 1 at the strictly lower region. For simplicity of notations, we would like to denote, for example

$$\mathbf{1}(m) = \mathbf{1}(m, s), \quad \mathbf{1} = \mathbf{1}(1, s), \quad \hat{\mathbf{1}} = \mathbf{1}(m, 1).$$

This notation rule will be used throughout the entire section.

### 7.1 Matrix Description of the Ultimate Spatial Matrix

In this subsection we devote to presenting a matrix description of how to get the ultimate spatial matrix. To do that, we define the  $ms$  order matrices

$$\underline{\mathbf{C}}(m) = \{c_{ij}(m)\}, \quad \underline{\mathbf{D}}(m) = \{d_{ij}(m)\}, \quad \underline{\mathbf{W}}(m; \vartheta) = \{d_{ij}(m)(\theta_{ij}(m) - \vartheta)\}, \quad (7.2a)$$

related to the description of the ESTDG method, and

$$\underline{\mathbf{\Sigma}}(m) = \{\sigma_{ij}(m)\}, \quad \underline{\mathbf{\Phi}}(m) = \{\phi_{ij}(m)\}, \quad \underline{\mathbf{Q}}(m; \vartheta) = \{q_{ij}(m; \vartheta)\}, \quad (7.2b)$$

related to the definition of temporal differences of stage solutions. Here all indices  $i$  and  $j$  are taken from 0 to  $ms - 1$ , and  $\vartheta$  is the parameter as mentioned in subsection 3.1.1. We would like to remark that all data in the above matrices are set to be zero, if they are not clearly stated or defined.

#### 7.1.1 Elemental Formula

The ultimate spatial matrix is obtained by running Algorithm 1 for  $\ell = 1, 2, \dots, \zeta$ , where the crucial calculation is the increment accumulation in Step 2.

To do that, we define a lower triangle matrix  $\underline{\mathbf{A}}^*(m) = \{a_{ij}^*(m)\}_{0 \leq i, j \leq ms-1}$ , whose entries are all defined to be zero except

$$a_{ij}^*(m) = (1 - \delta_{ij}/2)a_{i+1, j}^{(j)}(m), \quad \text{for } j \leq i \leq ms - 1 \quad \text{and} \quad 0 \leq j \leq \zeta - 1.$$

Since  $\{\tilde{q}_{ij}(m)\}_{0 \leq i, j \leq ms-1}$  is a lower triangle matrix, all summation ranges in Step 2 can be enlarged to  $\{0, 1, \dots, ms - 1\}$ . Gathering up the related operation till the matrix transferring process stops, we can obtain a unified description for the increment procedure at any fixed position. More specifically, the integrated calculation at every  $(i', j')$  position reads (dropping  $(m)$  here for convenience)

$$g_{i'j'} \leftarrow g_{i'j'} - a_{i'j'}^*; \quad g_{i'j'} \leftarrow g_{i'j'} + a_{k'j'}^* \tilde{q}_{k'i'}, \quad g_{i'j'} \leftarrow g_{i'j'} + a_{i'k'}^* \tilde{q}_{k'j'},$$

where  $i', j'$  and  $\kappa'$  go through  $\{0, 1, \dots, ms - 1\}$ . As a result, the total increment at Step 2 of Algorithm 1 can be expressed in the matrix form

$$\underline{\mathbf{G}}(m) = (2\vartheta - 1)\underline{\mathbf{A}}^*(m) + \underline{\mathbf{Q}}(m; \vartheta)^\top \underline{\mathbf{A}}^*(m) + \underline{\mathbf{A}}^*(m) \underline{\mathbf{Q}}(m; \vartheta),$$

where the definition (3.16), i.e.,  $\tilde{q}_{ij}(m) = q_{ij}(m; \vartheta) + \vartheta \delta_{ij}$  is used.

From Step 3 of Algorithm 1, we have the ultimate spatial matrix (below the last row and column is dropped, since they are always zero)

$$\begin{aligned} \mathbb{B}(m) &= \underline{\mathbf{G}}(m) + \underline{\mathbf{G}}(m)^\top \\ &= \left(\vartheta - \frac{1}{2}\right)\underline{\mathbf{B}}^*(m) + \frac{1}{2}\left[\underline{\mathbf{B}}^*(m)\underline{\mathbf{Q}}(m; \vartheta) + \underline{\mathbf{Q}}(m; \vartheta)^\top \underline{\mathbf{B}}^*(m)\right], \end{aligned} \tag{7.3}$$

with the symmetric matrix

$$\underline{\mathbf{B}}^*(m) = 2\underline{\mathbf{A}}^*(m) + 2\underline{\mathbf{A}}^*(m)^\top = \{b_{ij}^*(m)\}_{0 \leq i, j \leq ms-1}. \tag{7.4}$$

The entry at the lower triangular zone is defined as

$$b_{ij}^*(m) = \begin{cases} 2a_{i+1, j}^{(j)}(m), & 0 \leq j \leq \zeta - 1 \text{ and } j \leq i \leq ms - 1, \\ 0, & \text{otherwise,} \end{cases} \tag{7.5}$$

as the same as in the ultimate spatial matrix in [26] for the RKDG methods with a fixed numerical flux parameter.

To investigate the property of the second term in (7.3), we just need to study the perturbation matrix

$$\underline{\mathbf{Z}}(m; \vartheta) = \underline{\mathbf{B}}^*(m)\underline{\mathbf{Q}}(m; \vartheta) = \{z_{ij}(m; \vartheta)\}_{0 \leq i, j \leq ms-1}. \tag{7.6}$$

Taking into account on definition of the contribution index, we only pay attention on those left-top entries in (7.6). In what follows we try to deduce a convenient and unified formula for

$$z_{ij}(m; \vartheta) = \sum_{0 \leq \ell \leq ms-1} b_{i\ell}^*(m)q_{\ell, j}(m; \vartheta), \quad 0 \leq i, j \leq \zeta - 1. \tag{7.7}$$

The formula of every  $b_{i\ell}^*(m)$  has been given in [26], but is variant according to the size relationship of  $i$  and  $\ell$ . In this paper we have to rebuild an equivalent and unified formula, as stated in the next lemma.

**Lemma 7.1** For  $0 \leq i \leq \zeta - 1$ , there holds

$$b_{i\ell}^*(m) = 2 \sum_{0 \leq \kappa \leq i} (-1)^\kappa \alpha_{i-\kappa}(m) \alpha_{\ell+1+\kappa}(m), \quad 0 \leq \ell \leq ms - 1. \tag{7.8}$$

Here and below we define  $\alpha_{i'}(m) = 0$  if  $i' > ms$  for simplicity.

We put aside the proof of this lemma in Sect. 7.1.3. Substituting (7.8) into (7.7) deduces for any  $0 \leq i, j \leq \zeta - 1$  that

$$z_{ij}(m; \vartheta) = \sum_{0 \leq \kappa \leq i} 2(-1)^\kappa \alpha_{i-\kappa}(m) \underbrace{\sum_{0 \leq \ell \leq ms-1} \alpha_{\ell+1+\kappa}(m)q_{\ell, j}(m; \vartheta)}_{\pi_{\kappa, j}(m; \vartheta)}. \tag{7.9}$$

In what follows we would like to set up a useful formula of  $\pi_{\kappa, j}(m; \vartheta)$  by those data to define the ESTDG method.

### 7.1.2 Formula of $\pi_{\kappa,j}(m; \vartheta)$

Due to (3.13) and (3.8), we can respectively obtain

$$\underline{Q}(m; \vartheta) \underline{\Sigma}(m) = \underline{\Phi}(m) \underline{W}(m; \vartheta), \quad \underline{\Phi}(m) \underline{D}(m) = \underline{\Sigma}(m). \tag{7.10}$$

This implies  $\underline{Q}(m; \vartheta) = \underline{\Sigma}(m) \underline{D}(m)^{-1} \underline{W}(m; \vartheta) \underline{\Sigma}(m)^{-1}$ . With the short notation

$$\mathbf{y}^\top(m) = \sum_{0 \leq \ell \leq ms-1} \alpha_{\ell+1+\kappa}(m) \mathbf{e}_\ell^\top(m) \underline{\Sigma}(m), \tag{7.11}$$

it follows from (7.9) and  $q_{\ell,j}(m; \vartheta) = \mathbf{e}_\ell^\top(m) \underline{Q}(m; \vartheta) \mathbf{e}_j(m)$  that

$$\pi_{\kappa,j}(m; \vartheta) = \mathbf{y}^\top(m) \cdot \left[ \underline{D}^{-1}(m) \underline{W}(m; \vartheta) \right] \cdot \left[ \underline{\Sigma}(m)^{-1} \mathbf{e}_j(m) \right]. \tag{7.12}$$

Below we are going to express three terms in (7.12). To that end, we start this work from the calculation of  $\underline{\Sigma}(m)^{-1}$ .

By denoting (here and below we omit  $(m)$  for the matrix entry)

$$\underline{S}(m) = \underline{I}(m) - \underline{C}(m) \underline{E}(m) = \begin{pmatrix} 1 & & & & & \\ -c_{11} & 1 & & & & \\ -c_{21} & -c_{22} & 1 & & & \\ \vdots & \vdots & & \ddots & & \\ -c_{ms-1,1} & -c_{ms-1,2} & \cdots & -c_{ms-1,ms-1} & 1 & \end{pmatrix},$$

the definition procedure of the temporal differences of stage solutions can be written into the matrix form

$$\left( \begin{array}{c|c} \underline{\Sigma}(m) & \\ \hline \sigma_{ms,0} & \cdots \sigma_{ms,ms-1} \sigma_{ms,ms} \end{array} \right) = \left( \begin{array}{c|c} 1 & \\ \hline \underline{\Phi}(m) & \end{array} \right) \left( \begin{array}{c|c} 1 & \\ \hline -\underline{C}(m) \mathbf{e}_0(m) & \underline{S}(m) \end{array} \right).$$

Recalling the definition of the evolution identity, the matrix inversion on both sides of the above identity yields

$$\left( \begin{array}{c|c} \underline{\Sigma}(m)^{-1} & \\ \hline \alpha_0 & \cdots \alpha_{ms-1} \alpha_{ms} \end{array} \right) = \left( \begin{array}{c|c} 1 & \\ \hline \underline{S}(m)^{-1} \underline{C}(m) \mathbf{e}_0(m) & \underline{S}(m)^{-1} \underline{D}(m) \underline{\Sigma}(m)^{-1} \end{array} \right),$$

where we have used (7.10) to get  $\underline{\Phi}(m)^{-1} = \underline{D}(m) \underline{\Sigma}(m)^{-1}$ . Comparing with the matrices entries on both sides, we can achieve the following equalities for every column in the matrix  $\underline{\Sigma}(m)^{-1}$ ,

$$\underline{\Sigma}(m)^{-1} \mathbf{e}_0(m) = [\underline{I}(m) + \underline{E}(m) \underline{S}(m)^{-1} \underline{C}(m)] \mathbf{e}_0(m) \stackrel{\text{def}}{=} \mathbf{q}(m), \tag{7.13a}$$

$$\underline{\Sigma}(m)^{-1} \mathbf{e}_j(m) = \underbrace{\underline{E}(m) \underline{S}(m)^{-1} \underline{D}(m)}_{\underline{K}(m)} \underline{\Sigma}(m)^{-1} \mathbf{e}_{j-1}(m), \quad j \geq 1, \tag{7.13b}$$

and for every evolution coefficient in (3.11),

$$\alpha_0(m) = \mathbf{e}_{ms-1}(m)^\top \underline{S}(m)^{-1} \underline{C}(m) \mathbf{e}_0(m), \tag{7.14a}$$

$$\alpha_j(m) = \underbrace{\mathbf{e}_{ms-1}(m)^\top \underline{\mathbf{S}}(m)^{-1} \underline{\mathbf{D}}(m)}_{\mathbf{p}^\top(m)} \underline{\mathbf{S}}(m)^{-1} \mathbf{e}_{j-1}(m), \quad j \geq 1. \tag{7.14b}$$

Then, an induction process for (7.13) yields that

$$\underline{\mathbf{S}}(m)^{-1} \mathbf{e}_j(m) = \underline{\mathbf{K}}(m)^j \mathbf{q}(m), \quad j \geq 0, \tag{7.15}$$

and the matrix identity

$$\underline{\mathbf{S}}(m)^{-1} \underline{\mathbf{E}}(m) = \underline{\mathbf{K}}(m) \underline{\mathbf{S}}(m)^{-1}. \tag{7.16}$$

For any  $\kappa \geq 0$ , substituting (7.14b) into (7.11) yields

$$\begin{aligned} \mathbf{y}^\top(m) &= \mathbf{p}(m)^\top \underline{\mathbf{S}}(m)^{-1} \left[ \sum_{0 \leq \ell \leq ms-1} \mathbf{e}_{\ell+\kappa}(m) \mathbf{e}_\ell(m)^\top \right] \underline{\mathbf{S}}(m) \\ &= \mathbf{p}(m)^\top \underline{\mathbf{S}}(m)^{-1} \underline{\mathbf{E}}(m)^\kappa \underline{\mathbf{S}}(m) \\ &= \mathbf{p}(m)^\top [\underline{\mathbf{S}}(m)^{-1} \underline{\mathbf{E}}(m) \underline{\mathbf{S}}(m)]^\kappa = \mathbf{p}(m)^\top \underline{\mathbf{K}}(m)^\kappa, \end{aligned} \tag{7.17}$$

where (7.16) is used at the last step. Substituting (7.17) and (7.15) into (7.12), we finally have

$$\pi_{\kappa,j}(m; \vartheta) = \mathbf{p}(m)^\top \underline{\mathbf{K}}(m)^\kappa \underline{\mathbf{D}}^{-1}(m) \underline{\mathbf{W}}(m; \vartheta) \underline{\mathbf{K}}(m)^j \mathbf{q}(m). \tag{7.18}$$

In order to investigate the relationship between this quantity and the multistep number, we would like to make some (right) Kronecker product of matrices [25] to simplify each term in (7.18). For example, we will use

$$\mathbf{e}_0(m) = \hat{\mathbf{e}}_0 \otimes \mathbf{e}_0, \quad \mathbf{e}_{ms-1}(m)^\top = \hat{\mathbf{e}}_{m-1}^\top \otimes \mathbf{e}_{s-1}^\top, \quad \underline{\mathbf{I}}(m) = \hat{\underline{\mathbf{I}}} \otimes \underline{\mathbf{I}}, \tag{7.19}$$

which implies

$$\underline{\mathbf{E}}(m) = \hat{\underline{\mathbf{I}}} \otimes \underline{\mathbf{E}} + \hat{\underline{\mathbf{E}}} \otimes \mathbf{e}_0 \mathbf{e}_{s-1}^\top. \tag{7.20}$$

Due to the definition (3.3), we derive

$$\underline{\mathbf{C}}(m) = \hat{\underline{\mathbf{I}}} \otimes \underline{\mathbf{C}}, \quad \underline{\mathbf{D}}(m) = \frac{1}{m} \hat{\underline{\mathbf{I}}} \otimes \underline{\mathbf{D}}, \quad \underline{\mathbf{W}}(m; \vartheta) = \frac{1}{m} \hat{\underline{\mathbf{I}}} \otimes \underline{\mathbf{W}}(\vartheta), \tag{7.21}$$

where  $\underline{\mathbf{W}}(\vartheta) = \underline{\mathbf{W}}(1; \vartheta)$ . Based on these identities, by some lengthy and tedious matrices manipulations, we can get the following important conclusions

$$\underline{\mathbf{S}}(m)^{-1} = \hat{\underline{\mathbf{L}}} \otimes \underline{\mathbf{S}}^{-1} \underline{\mathbf{C}} \mathbf{e}_0 \mathbf{e}_{s-1}^\top \underline{\mathbf{S}}^{-1} + \hat{\underline{\mathbf{I}}} \otimes \underline{\mathbf{S}}^{-1}, \tag{7.22a}$$

$$\underline{\mathbf{K}}(m) = \frac{1}{m} \left[ \hat{\underline{\mathbf{L}}} \otimes \mathbf{q} \mathbf{p}^\top + \hat{\underline{\mathbf{I}}} \otimes \underline{\mathbf{E}} \underline{\mathbf{S}}^{-1} \underline{\mathbf{D}} \right], \tag{7.22b}$$

$$\mathbf{p}(m)^\top = \frac{1}{m} \hat{\underline{\mathbf{I}}}^\top \otimes \mathbf{p}^\top, \tag{7.22c}$$

$$\mathbf{q}(m) = \hat{\underline{\mathbf{I}}} \otimes \mathbf{q}. \tag{7.22d}$$

In this process, we have used the following simple conclusions

$$\hat{\underline{\mathbf{E}}} + \hat{\underline{\mathbf{E}}} \hat{\underline{\mathbf{L}}} = \hat{\underline{\mathbf{L}}}, \quad \hat{\mathbf{e}}_{m-1}^\top + \hat{\mathbf{e}}_{m-1}^\top \hat{\underline{\mathbf{L}}} = \hat{\underline{\mathbf{I}}}^\top, \quad \hat{\mathbf{e}}_0 + \hat{\underline{\mathbf{L}}} \hat{\mathbf{e}}_0 = \hat{\underline{\mathbf{I}}}, \tag{7.23}$$

and an important identity as a corollary of (7.14a) and  $\alpha_0(m) = 1$ ,

$$\mathbf{e}_{ms-1}^\top(m) \underline{\mathbf{S}}(m)^{-1} \underline{\mathbf{C}}(m) \mathbf{e}_0(m) = 1. \tag{7.24}$$

For ease of reading, we present the verifications of (7.22) in Sect. 7.1.4.



With the help of (7.21), substituting (7.22c) and (7.22d) into (7.18) yields the final simplification expression

$$\pi_{\kappa,j}(m; \vartheta) = \frac{1}{m} \left( \hat{\mathbf{1}}^\top \otimes \mathbf{p}^\top \right) \underline{\mathbf{K}}(m)^\kappa \left( \hat{\mathbf{I}} \otimes \underline{\mathbf{D}}^{-1} \underline{\mathbf{W}}(\vartheta) \right) \underline{\mathbf{K}}(m)^j \left( \mathbf{1} \otimes \mathbf{q} \right). \tag{7.25}$$

If needed, we can use (7.22b) to further deal with  $\underline{\mathbf{K}}(m)$ .

### 7.1.3 Proof of Lemma 7.1

To end this subsection, we need to prove the skipped Lemma 7.1. Since all related manipulation does not depend on the spatial discretization, the results given in [26, Lemma 3.1] still hold. Hence, for  $0 \leq j' \leq \zeta$  and  $j' < i' \leq ms$  we have

$$a_{i'j'}^{(j')} (m) = \sum_{0 \leq \kappa \leq j'} (-1)^\kappa \alpha_{i'+\kappa}(m) \alpha_{j'-\kappa}(m), \tag{7.26a}$$

and for  $1 \leq i' \leq \zeta$  we have

$$a_{i'i'}^{(i')} (m) = \sum_{-i' \leq \kappa \leq i'} (-1)^\kappa \alpha_{i'+\kappa}(m) \alpha_{i'-\kappa}(m). \tag{7.26b}$$

Based on the formulas in (7.26), we can prove this lemma by simple discussions for different cases of  $\ell$ .

If  $\ell > i$ , since  $\mathbb{B}^*(m)$  is symmetric, it follows from (7.5) that

$$b_{i\ell}^*(m) = b_{\ell i}^*(m) = 2a_{\ell+1,i}^{(i)}(m).$$

This proves (7.8) by using (7.26a) with  $i' = \ell + 1$  and  $j' = i$ .

Otherwise, if  $\ell \leq i$ , we similarly have from (7.26a) that

$$b_{i\ell}^*(m) = 2a_{i+1,\ell}^{(\ell)}(m) = 2 \sum_{0 \leq \kappa \leq \ell} (-1)^\kappa \alpha_{i+1+\kappa}(m) \alpha_{\ell-\kappa}(m).$$

To show it can be written in (7.8), we just need to show  $\mathcal{Y} = 0$ , with

$$\begin{aligned} \mathcal{Y} &\stackrel{\text{def}}{=} \sum_{0 \leq \kappa \leq \ell} (-1)^\kappa \alpha_{i+1+\kappa}(m) \alpha_{\ell-\kappa}(m) - \sum_{0 \leq \kappa \leq i} (-1)^\kappa \alpha_{i-\kappa}(m) \alpha_{\ell+1+\kappa}(m) \\ &= \sum_{0 \leq \kappa \leq \ell+i+1} (-1)^{\ell-\kappa} \alpha_\kappa(m) \alpha_{\ell+i+1-\kappa}(m). \end{aligned}$$

Here we have respectively used the replacements of index  $\kappa' = \ell - \kappa$  and  $\kappa' = \ell + 1 + \kappa$  in two summations of the first equality. The verification is easy as follows.

- If  $\ell + i + 1$  is odd, the replacement  $\kappa' = i + \ell + 1 - \kappa$  implies  $\mathcal{Y} = (-1)^{i+\ell+1} \mathcal{Y}$  and hence  $\mathcal{Y} = 0$ .
- Otherwise, if  $\ell + i + 1$  is even, denoted by  $2L$ , a simple replacement of summation index again reduces

$$(-1)^{\ell-L} \mathcal{Y} = \sum_{-L \leq \kappa \leq L} (-1)^\kappa \alpha_{L+\kappa}(m) \alpha_{L-\kappa}(m) = a_{L,L}^{(L)}(m),$$

where the last step uses (7.26b). Since  $L < \zeta$ , it follows  $a_{L,L}^{(L)}(m) = 0$  from the definition of  $\zeta$ . This implies  $\mathcal{Y} = 0$  also.

Till now we sum up the above conclusions and complete the proof of this lemma.

### 7.1.4 Verifications of (7.22)

To verify the first identity (7.22a), we start from the definition of  $\underline{S}(m)$ . Substituting the identities (7.19), (7.21) and (7.20), we have

$$\underline{S}(m) = \underline{I}(m) - \underline{C}(m)\underline{E}(m) = \hat{\underline{I}} \otimes \underline{I} - (\hat{\underline{I}} \otimes \underline{C})(\hat{\underline{I}} \otimes \underline{E} + \hat{\underline{E}} \otimes e_0 e_{s-1}^\top).$$

Expanding the right-hand side and using the definition of  $\underline{S}$ , after some manipulations we have

$$\begin{aligned} \underline{S}(m) &= \hat{\underline{I}} \otimes (\underline{I} - \underline{C}\underline{E}) - \hat{\underline{E}} \otimes \underline{C}e_0 e_{s-1}^\top \\ &= \hat{\underline{I}} \otimes \underline{S} - \hat{\underline{E}} \otimes \underline{C}e_0 e_{s-1}^\top = (\hat{\underline{I}} \otimes \underline{I} - \hat{\underline{E}} \otimes \underline{C}e_0 e_{s-1}^\top \underline{S}^{-1})(\hat{\underline{I}} \otimes \underline{S}). \end{aligned}$$

Since  $(\hat{\underline{E}})^m$  is a zero matrix, the inverse of the first matrix is expressed by

$$\hat{\underline{I}} \otimes \underline{I} + \sum_{1 \leq i \leq m-1} (\hat{\underline{E}})^i \otimes (\underline{C}e_0 e_{s-1}^\top \underline{S}^{-1})^i.$$

Using (7.24), we have for any  $i \geq 1$  that

$$(\underline{C}e_0 e_{s-1}^\top \underline{S}^{-1})^i = \underline{C}e_0 (e_{s-1}^\top \underline{S}^{-1} \underline{C}e_0)^{i-1} e_{s-1}^\top \underline{S}^{-1} = \underline{C}e_0 e_{s-1}^\top \underline{S}^{-1}.$$

Summing up the above identities, we have

$$\begin{aligned} \underline{S}(m)^{-1} &= (\hat{\underline{I}} \otimes \underline{S}^{-1}) \left[ \hat{\underline{I}} \otimes \underline{I} + \sum_{1 \leq i \leq m-1} (\hat{\underline{E}})^i \otimes \underline{C}e_0 e_{s-1}^\top \underline{S}^{-1} \right] \\ &= (\hat{\underline{I}} \otimes \underline{S}^{-1}) (\hat{\underline{I}} \otimes \underline{I} + \hat{\underline{L}} \otimes \underline{C}e_0 e_{s-1}^\top \underline{S}^{-1}) \\ &= \hat{\underline{I}} \otimes \underline{S}^{-1} + \hat{\underline{L}} \otimes \underline{S}^{-1} \underline{C}e_0 e_{s-1}^\top \underline{S}^{-1}, \end{aligned}$$

where we have used the definition (7.1) of  $\hat{\underline{L}}$  at the second step. This completes the verification of (7.22a).

We start the verification of (7.22b) from the definition (7.13b) of  $\underline{K}(m)$ . Substituting the identities (7.20), (7.22a) and (7.21), we have

$$\begin{aligned} m\underline{K}(m) &= m\underline{E}(m)\underline{S}(m)^{-1}\underline{D}(m) \\ &= (\hat{\underline{I}} \otimes \underline{E} + \hat{\underline{E}} \otimes e_0 e_{s-1}^\top)(\hat{\underline{I}} \otimes \underline{S}^{-1} + \hat{\underline{L}} \otimes \underline{S}^{-1} \underline{C}e_0 e_{s-1}^\top \underline{S}^{-1})(\hat{\underline{I}} \otimes \underline{D}). \end{aligned}$$

Expanding the right-hand side, using (7.24) and the first identity in (7.23), we achieve

$$\begin{aligned} m\underline{K}(m) &= \hat{\underline{I}} \otimes \underline{E}\underline{S}^{-1}\underline{D} + \hat{\underline{L}} \otimes \underline{E}\underline{S}^{-1}\underline{C}e_0 e_{s-1}^\top \underline{S}^{-1}\underline{D} + \hat{\underline{L}} \otimes e_0 e_{s-1}^\top \underline{S}^{-1}\underline{D} \\ &= \hat{\underline{I}} \otimes \underline{E}\underline{S}^{-1}\underline{D} + \hat{\underline{L}} \otimes (\underline{I} + \underline{E}\underline{S}^{-1}\underline{C})e_0 e_{s-1}^\top \underline{S}^{-1}\underline{D} \\ &= \hat{\underline{I}} \otimes \underline{E}\underline{S}^{-1}\underline{D} + \hat{\underline{L}} \otimes \underline{q}\underline{p}^\top, \end{aligned}$$

where at the last step we have used the definitions of  $\underline{q}$  and  $\underline{p}^\top$  in (7.13a) and (7.14b). This completes the verification of (7.22b).

The third identity (7.22c) is verified along the same line. Starting from the definition of  $\underline{p}(m)^\top$  in (7.14b), and substituting the identities (7.19), (7.22a) and (7.21), we have

$$\begin{aligned} m\underline{p}(m)^\top &= m e_{ms-1}(m)^\top \underline{S}(m)^{-1} \underline{D}(m) \\ &= (\hat{e}_{m-1}^\top \otimes e_{s-1}^\top)(\hat{\underline{I}} \otimes \underline{S}^{-1} + \hat{\underline{L}} \otimes \underline{S}^{-1} \underline{C}e_0 e_{s-1}^\top \underline{S}^{-1})(\hat{\underline{I}} \otimes \underline{D}). \end{aligned}$$

Expanding the above expression and using (7.24), we have

$$\begin{aligned} m \mathbf{p}(m)^\top &= \hat{\mathbf{e}}_{m-1}^\top \otimes \mathbf{e}_{s-1}^\top \underline{\mathbf{S}}^{-1} \underline{\mathbf{D}} + \hat{\mathbf{e}}_{m-1}^\top \hat{\underline{\mathbf{L}}} \otimes \mathbf{e}_{s-1}^\top \underline{\mathbf{S}}^{-1} \underline{\mathbf{C}} \mathbf{e}_0 \mathbf{e}_{s-1}^\top \underline{\mathbf{S}}^{-1} \underline{\mathbf{D}} \\ &= \hat{\mathbf{e}}_{m-1}^\top \otimes \mathbf{e}_{s-1}^\top \underline{\mathbf{S}}^{-1} \underline{\mathbf{D}} + \hat{\mathbf{e}}_{m-1}^\top \hat{\underline{\mathbf{L}}} \otimes \mathbf{e}_{s-1}^\top \underline{\mathbf{S}}^{-1} \underline{\mathbf{D}} \\ &= (\hat{\mathbf{e}}_{m-1}^\top + \hat{\mathbf{e}}_{m-1}^\top \hat{\underline{\mathbf{L}}}) \otimes \mathbf{e}_{s-1}^\top \underline{\mathbf{S}}^{-1} \underline{\mathbf{D}} = \hat{\mathbf{1}}^\top \otimes \mathbf{p}^\top, \end{aligned}$$

where at the last step we have used the second identity in (7.23) and the definition of  $\mathbf{p}^\top$  in (7.14b). This proves (7.22c).

The fourth identity (7.22d) is verified similarly. To save the length of this paper, we omit the detailed procedure.

## 7.2 Some Proofs

In this subsection we would like to prove Lemmas 3.6 and 3.7, as well as Propositions 3.1 and 3.2.

### 7.2.1 Proof of Lemma 3.6

Recalling the definition of  $\pi_{\kappa,j}(m; \vartheta)$ , given in (7.9), it follows from (3.16) and (3.27) that  $\Theta(m) = \vartheta + \pi_{00}(m; \vartheta)$ . Substituting (7.25) implies that

$$\Theta(m) = \vartheta + \frac{1}{m} (\hat{\mathbf{1}}^\top \otimes \mathbf{p}^\top) (\hat{\underline{\mathbf{I}}} \otimes \underline{\mathbf{D}}^{-1} \underline{\mathbf{W}}(\vartheta)) (\hat{\mathbf{1}} \otimes \mathbf{q}) = \vartheta + \mathbf{p}^\top \underline{\mathbf{D}}^{-1} \underline{\mathbf{W}}(\vartheta) \mathbf{q}, \tag{7.27}$$

where the simple fact  $\hat{\mathbf{1}}^\top \hat{\underline{\mathbf{I}}} \hat{\mathbf{1}} = m$  is used. This completes the proof of Lemma 3.6.

**Remark 7.1** Taking  $m = 1$  and  $\vartheta = \Theta$  in (7.27), we use Lemma 3.6 to get

$$\mathbf{p}^\top \underline{\mathbf{D}}^{-1} \underline{\mathbf{W}}(\Theta) \mathbf{q} = 0. \tag{7.28}$$

This is just the conclusion in Lemma 3.5 with  $m = 1$ . As an essence property of the averaged numerical flux parameter, it plays an important role in the proof of Lemma 3.7.

### 7.2.2 Proof of Lemma 3.7

For convenience of notations, in what follows we use a generic notation  $C$  to denote a positive constant independent of  $m$ . Recalling the proof of [26, Proposition 3.3], we have for  $0 \leq i, j \leq \zeta - 1$  that

$$\left| b_{ij}^*(m) - \frac{2}{i!j!(i+j+1)} \right| \leq \frac{C}{m}, \tag{7.29}$$

where  $\{\frac{2}{i!j!(i+j+1)}\}_{0 \leq i, j \leq \zeta-1}$  forms a symmetric positive definite matrix congruent to an Hilbert matrix. Since  $\Theta > 1/2$ , it follows from (7.3) and (7.6) with  $\vartheta = \Theta$  that we can prove this lemma by showing that  $z_{ij}(m; \Theta)$  for  $0 \leq i, j \leq \zeta - 1$  all tends to zero as  $m$  goes to infinity. By (7.9), it is sufficient to prove

$$|\pi_{\kappa,j}(m; \Theta)| \leq \frac{C}{m}, \quad 0 \leq \kappa, j \leq \zeta - 1, \tag{7.30}$$

since [26, inequality (3.16)] shows that  $\alpha_{i-\kappa}(m)$  is bounded independent of  $m$ .

Denote  $\pi_{\kappa,j} = \pi_{\kappa,j}(m; \Theta)$  and  $\underline{W} = \underline{W}(\Theta)$  for simplicity. Below we are going to prove (7.30) for different cases of  $\kappa$  and  $j$ , where (7.28) plays an important role to well control the accumulation and growth as  $m$  goes to infinity.

- If  $\kappa = j = 0$ , we have  $\pi_{0,0} = (\hat{\mathbf{1}}^\top \hat{\mathbf{1}}) \otimes (\mathbf{p}^\top \underline{D}^{-1} \underline{W} \mathbf{q}) = 0$ , due to (7.28).
- If  $\kappa > 0$  and  $j > 0$ , we have

$$\pi_{\kappa,j} = \frac{1}{m} (\hat{\mathbf{1}}^\top \otimes \mathbf{p}^\top) [\underline{K}(m)]^{\kappa-1} \underline{\Pi}_{\kappa,j}(m) [\underline{K}(m)]^{j-1} (\hat{\mathbf{1}} \otimes \mathbf{q}), \tag{7.31}$$

where  $\underline{\Pi}_{\kappa,j}(m) = \underline{K}(m) (\hat{\mathbf{1}} \otimes \underline{D}^{-1} \underline{W}) \underline{K}(m)$ . Substituting (7.22b) into this formula and then using (7.28) to eliminate the term involving  $\hat{\underline{L}}^2$ . After some manipulations we yield

$$\begin{aligned} \underline{\Pi}_{\kappa,j}(m) &= \frac{1}{m^2} \hat{\underline{L}} \otimes [\mathbf{q} \mathbf{p}^\top \underline{D}^{-1} \underline{W} \underline{E} \underline{S}^{-1} \underline{D} + \underline{E} \underline{S}^{-1} \underline{W} \mathbf{q} \mathbf{p}^\top] \\ &\quad + \frac{1}{m^2} \hat{\underline{L}} \otimes \underline{E} \underline{S}^{-1} \underline{W} \underline{E} \underline{S}^{-1} \underline{D}. \end{aligned}$$

The row norms for all matrices (including the row vectors and column vectors) do not depend on  $m$ , except that  $\|\hat{\underline{L}}\|_\infty = m - 1$ . Hence we have

$$\|\underline{\Pi}_{\kappa,j}(m)\|_\infty \leq \frac{C}{m}.$$

Noticing  $\|\frac{1}{m} (\hat{\mathbf{1}}^\top \otimes \mathbf{p}^\top)\|_\infty \leq C$  and  $\|\underline{K}(m)\|_\infty \leq C$ , we get from (7.31) what we want to prove.

- If  $\kappa = 0$  and  $j > 0$ , we have  $\pi_{0,j} = \frac{1}{m} \underline{\Pi}_{0,j}(m) [\underline{K}(m)]^{j-1} (\hat{\mathbf{1}} \otimes \mathbf{q})$  with

$$\underline{\Pi}_{0,j}(m) = (\hat{\mathbf{1}}^\top \otimes \mathbf{p}^\top) (\hat{\mathbf{1}} \otimes \underline{D}^{-1} \underline{W}) \underline{K}(m) = \frac{1}{m} \hat{\mathbf{1}}^\top \otimes \mathbf{p}^\top \underline{D}^{-1} \underline{W} \underline{E} \underline{S}^{-1} \underline{D},$$

by some manipulations with the help of (7.22b) and (7.28). The remaining proof follows the same line as above, hence is omitted.

- If  $\kappa > 0$  and  $j = 0$ , we have  $\pi_{\kappa,0} = \frac{1}{m} (\hat{\mathbf{1}}^\top \otimes \mathbf{p}^\top) [\underline{K}(m)]^{\kappa-1} \underline{\Pi}_{\kappa,0}(m)$ , where

$$\underline{\Pi}_{\kappa,0}(m) = \underline{K}(m) (\hat{\mathbf{1}} \otimes \underline{D}^{-1} \underline{W}) (\hat{\mathbf{1}} \otimes \mathbf{q}) = \hat{\mathbf{1}} \otimes \underline{E} \underline{S}^{-1} \underline{W} \mathbf{q},$$

with the help of (7.22b) and (7.28). Then we can prove (7.31) as above.

Summing up the above conclusions, we verify (7.30) and then prove this lemma.

### 7.2.3 Proof of Propositions 3.1 and 3.2

Taking  $\vartheta = 0$  in (7.27) and substituting the definition of  $\mathbf{p}^\top$  and  $\mathbf{q}$ , we have

$$\Theta = \mathbf{e}_{s-1}^\top \underline{S}^{-1} \underline{W}(0) (\underline{I} + \underline{E} \underline{S}^{-1} \underline{C}) \mathbf{e}_0. \tag{7.32}$$

This identity will be used to prove these propositions.

Since we have assumed  $c_{\ell\kappa} \geq 0$  for any  $\ell$  and  $\kappa$  in this paper, all entries of  $\underline{S}^{-1}$  are non-negative due to the simple fact

$$\underline{S}^{-1} = (\underline{I} - \underline{E} \underline{C})^{-1} = \underline{I} + \sum_{1 \leq i \leq s-1} (\underline{E} \underline{C})^i.$$

Hence we can conclude from (7.32) that  $\Theta$  is a non-negative linear combination of the entries of  $\underline{W}(0) = \{d_{\ell\kappa}\theta_{\ell\kappa}\}_{0 \leq \ell, \kappa \leq s-1}$ . As a trivial conclusion for special case that all numerical flux parameters are the same, it is easy to conclude that  $\Theta$  is a weighted average of  $\theta_{\ell\kappa}$ . This proves Proposition 3.1.

**Remark 7.2** This is the only place that the condition  $c_{\ell\kappa} \geq 0$  is used in this paper.

For the LWDG method with the time marching coefficients (2.10), we have  $\underline{S} = \underline{I}$  and then get from (7.32) that

$$\Theta = \mathbf{e}_{s-1}^\top \underline{W}(0) \mathbf{e}_0 = d_{s-1,0} \theta_{s-1,0} = \theta_{s-1,0},$$

since  $\underline{I} + \underline{E}\underline{S}^{-1}\underline{C} = \underline{I} + \underline{E}\underline{C} = \underline{I}$ . This completes the proof of Proposition 3.2.

**Funding** Yuan Xu is supported by NSFC Grant 12301513, Natural Science Foundation of Jiangsu Province Grant BK20230374 and Natural Science Foundation of Jiangsu Higher Education Institutions of China Grant 23KJB110019. Chi-Wang Shu is supported by NSF grant DMS-2309249. Qiang Zhang is supported by NSFC Grant 12071214.

**Data Availability** The datasets generated during the current study are available from the corresponding author upon reasonable request.

## Declarations

**Conflict of interest** The authors have no relevant financial or non-financial interests to disclose.

## References

1. Ai, J., Xu, Y., Shu, C.W., Zhang, Q.:  $L^2$  error estimate to smooth solutions of high order Runge–Kutta discontinuous Galerkin method for scalar nonlinear conservation laws with and without sonic points. *SIAM J. Numer. Anal.* **60**(4), 1741–1773 (2022). <https://doi.org/10.1137/21M1435495>
2. Chavent, G., Cockburn, B.: The local projection  $P^0 P^1$ -discontinuous-Galerkin finite element method for scalar conservation laws. *RAIRO Modél. Math. Anal. Numér.* **23**(4), 565–592 (1989). <https://doi.org/10.1051/m2an/1989230405651>
3. Cheng, Y., Meng, X., Zhang, Q.: Application of generalized Gauss–Radau projections for the local discontinuous Galerkin method for linear convection–diffusion equations. *Math. Comp.* **86**(305), 1233–1267 (2017). <https://doi.org/10.1090/mcom/3141>
4. Ciarlet, P.G.: The finite element method for elliptic problems. In: *Studies in Mathematics and its Applications*, vol. 4. North-Holland Publishing Co., New York (1978)
5. Cockburn, B., Hou, S., Shu, C.W.: The Runge–Kutta local projection discontinuous Galerkin finite element method for conservation laws. IV. The multidimensional case. *Math. Comput.* **54**(190), 545–581 (1990). <https://doi.org/10.2307/2008501>
6. Cockburn, B., Lin, S.Y., Shu, C.W.: TVB Runge–Kutta local projection discontinuous Galerkin finite element method for conservation laws. III. One-dimensional systems. *J. Comput. Phys.* **84**(1), 90–113 (1989). [https://doi.org/10.1016/0021-9991\(89\)90183-6](https://doi.org/10.1016/0021-9991(89)90183-6)
7. Cockburn, B., Shu, C.W.: TVB Runge–Kutta local projection discontinuous Galerkin finite element method for conservation laws. II. General framework. *Math. Comput.* **52**(186), 411–435 (1989). <https://doi.org/10.2307/2008474>
8. Cockburn, B., Shu, C.W.: The Runge–Kutta local projection  $P^1$ -discontinuous-Galerkin finite element method for scalar conservation laws. *RAIRO Modél. Math. Anal. Numér.* **25**(3), 337–361 (1991). <https://doi.org/10.1051/m2an/1991250303371>
9. Cockburn, B., Shu, C.W.: The Runge–Kutta discontinuous Galerkin method for conservation laws. V. Multidimensional systems. *J. Comput. Phys.* **141**(2), 199–224 (1998). <https://doi.org/10.1006/jcph.1998.5892>
10. Cockburn, B., Shu, C.W.: Runge–Kutta discontinuous Galerkin methods for convection-dominated problems. *J. Sci. Comput.* **16**(3), 173–261 (2001). <https://doi.org/10.1023/A:1012873910884>

11. Gottlieb, S., Ruuth, S.J.: Optimal strong-stability-preserving time-stepping schemes with fast downwind spatial discretizations. *J. Sci. Comput.* **27**(1–3), 289–303 (2006). <https://doi.org/10.1007/s10915-005-9054-8>
12. Gottlieb, S., Shu, C.W.: Total variation diminishing Runge–Kutta schemes. *Math. Comput.* **67**(221), 73–85 (1998). <https://doi.org/10.1090/S0025-5718-98-00913-2>
13. Guo, W., Qiu, J., Qiu, J.: A new Lax–Wendroff discontinuous Galerkin method with superconvergence. *J. Sci. Comput.* **65**(1), 299–326 (2015). <https://doi.org/10.1007/s10915-014-9968-0>
14. Liu, Y., Shu, C.W., Zhang, M.: Sub-optimal convergence of discontinuous Galerkin methods with central fluxes for linear hyperbolic equations with even degree polynomial approximations. *J. Comput. Math.* **39**(4), 518–537 (2021). <https://doi.org/10.4208/jcm.2002-m2019-0305>
15. Qiu, J., Zhang, Q.: Stability, error estimate and limiters of discontinuous Galerkin methods. In: *Handbook of Numerical Methods for Hyperbolic Problems, Handbook of Numerical Analysis*, vol. 17, pp. 147–171. Elsevier, Amsterdam (2016). <https://doi.org/10.1016/bs.hna.2016.06.001>
16. Ruuth, S.J.: Global optimization of explicit strong-stability-preserving Runge–Kutta methods. *Math. Comput.* **75**(253), 183–207 (2006). <https://doi.org/10.1090/S0025-5718-05-01772-2>
17. Ruuth, S.J., Spiteri, R.J.: Two barriers on strong-stability-preserving time discretization methods. *J. Sci. Comput.* **17**(1–4), 211–220 (2002). <https://doi.org/10.1023/A:1015156832269>
18. Ruuth, S.J., Spiteri, R.J.: High-order strong-stability-preserving Runge–Kutta methods with downwind-biased spatial discretizations. *SIAM J. Numer. Anal.* **42**(3), 974–996 (2004). <https://doi.org/10.1137/S0036142902419284>
19. Shu, C.W.: Total-variation-diminishing time discretizations. *SIAM J. Sci. Stat. Comput.* **9**(6), 1073–1084 (1988). <https://doi.org/10.1137/0909073>
20. Shu, C.W.: Discontinuous Galerkin methods: general approach and stability. In: *Numerical Solutions of Partial Differential Equations, Adv. Courses Math. CRM Barcelona*, pp. 149–201. Birkhäuser, Basel (2009)
21. Shu, C.W.: Discontinuous Galerkin methods for time-dependent convection dominated problems: basics, recent developments and comparison with other methods. In: *Building Bridges: Connections and Challenges in Modern Approaches to Numerical Partial Differential Equations, Lecture Notes in Computer Science Engineering*, vol. 114, pp. 369–397. Springer, New York (2016)
22. Shu, C.W., Osher, S.: Efficient implementation of essentially nonoscillatory shock-capturing schemes. *J. Comput. Phys.* **77**(2), 439–471 (1988). [https://doi.org/10.1016/0021-9991\(88\)90177-5](https://doi.org/10.1016/0021-9991(88)90177-5)
23. Sun, Z., Shu, C.W.: Stability analysis and error estimates of Lax–Wendroff discontinuous Galerkin methods for linear conservation laws. *ESAIM Math. Model. Numer. Anal.* **51**(3), 1063–1087 (2017). <https://doi.org/10.1051/m2an/2016049>
24. Sun, Z., Shu, C.W.: Strong stability of explicit Runge–Kutta time discretizations. *SIAM J. Numer. Anal.* **57**(3), 1158–1182 (2019). <https://doi.org/10.1137/18M122892X>
25. Van Loan, C.F.: The ubiquitous Kronecker product. *J. Comput. Appl. Math.* **123**(1–2), 85–100 (2000). [https://doi.org/10.1016/S0377-0427\(00\)00393-9](https://doi.org/10.1016/S0377-0427(00)00393-9)
26. Xu, Y., Meng, X., Shu, C.W., Zhang, Q.: Superconvergence analysis of the Runge–Kutta discontinuous Galerkin methods for a linear hyperbolic equation. *J. Sci. Comput.* **84**, 23 (2020). <https://doi.org/10.1007/s10915-020-01274-1>
27. Xu, Y., Shu, C.W., Zhang, Q.: Error estimate of the fourth-order Runge–Kutta discontinuous Galerkin methods for linear hyperbolic equations. *SIAM J. Numer. Anal.* **58**(5), 2885–2914 (2020). <https://doi.org/10.1137/19M1280077>
28. Xu, Y., Zhang, Q.: Superconvergence analysis of the Runge–Kutta discontinuous Galerkin method with upwind-biased numerical flux for two dimensional linear hyperbolic equation. *Commun. Appl. Math. Comput.* **4**, 319–352 (2022). <https://doi.org/10.1007/s42967-020-00116-z>
29. Xu, Y., Zhang, Q., Shu, C.W., Wang, H.: The  $L^2$ -norm stability analysis of Runge–Kutta discontinuous Galerkin methods for linear hyperbolic equations. *SIAM J. Numer. Anal.* **57**(4), 1574–1601 (2019). <https://doi.org/10.1137/18M1230700>
30. Xu, Y., Zhao, D., Zhang, Q.: Local error estimates for Runge–Kutta discontinuous Galerkin methods with upwind-biased numerical fluxes for a linear hyperbolic equation in one-dimension with discontinuous initial data. *J. Sci. Comput.* **91**, 11 (2022). <https://doi.org/10.1007/s10915-022-01793-z>
31. Zhang, Q., Shu, C.W.: Error estimates to smooth solutions of Runge–Kutta discontinuous Galerkin methods for scalar conservation laws. *SIAM J. Numer. Anal.* **42**(2), 641–666 (2004). <https://doi.org/10.1137/S0036142902404182>
32. Zhang, Q., Shu, C.W.: Stability analysis and a priori error estimates of the third order explicit Runge–Kutta discontinuous Galerkin method for scalar conservation laws. *SIAM J. Numer. Anal.* **48**(3), 1038–1063 (2010). <https://doi.org/10.1137/090771363>

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

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.