

A *PRIORI* ERROR ESTIMATES TO SMOOTH SOLUTIONS OF THE THIRD ORDER RUNGE–KUTTA DISCONTINUOUS GALERKIN METHOD FOR SYMMETRIZABLE SYSTEMS OF CONSERVATION LAWS ^{*,**}

JUAN LUO¹, CHI-WANG SHU² AND QIANG ZHANG¹

Abstract. In this paper we present an *a priori* error estimate of the Runge–Kutta discontinuous Galerkin method for solving symmetrizable conservation laws, where the time is discretized with the third order explicit total variation diminishing Runge–Kutta method and the finite element space is made up of piecewise polynomials of degree $k \geq 2$. Quasi-optimal error estimate is obtained by energy techniques, for the so-called generalized E-fluxes under the standard temporal-spatial CFL condition $\tau \leq \gamma h$, where h is the element length and τ is time step, and γ is a positive constant independent of h and τ . Optimal estimates are also considered when the upwind numerical flux is used.

Mathematics Subject Classification. 65M60, 65M12.

Received August 22, 2013. Revised November 2, 2014.

Published online June 19, 2015.

1. INTRODUCTION

In this paper we would like to continue the works in [29–31] and present error estimates of the Runge–Kutta discontinuous Galerkin (RKDG) method for smooth solutions of symmetrizable systems of conservation laws. For simplicity of presentation, we consider the model equation in the spatial domain $I = (0, 1)$ and the time interval $[0, T]$,

$$\mathbf{u}_t + \mathbf{f}(\mathbf{u})_x = 0, \quad (x, t) \in I \times (0, T], \quad (1.1)$$

with the initial solution $\mathbf{u}_0(x)$. Here $\mathbf{u}(x, t) : \mathbb{R} \times \mathbb{R}^+ \rightarrow \mathbb{R}^m$ is the unknown vector-valued solution, and $\mathbf{f}(\mathbf{u}) : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is the given smooth flux function. We will, however, point out similarities and differences when the analysis is generalized to multiple space dimensions. Furthermore, we do not pay much attention to boundary conditions in this paper; hence the solution is considered to be periodic or compactly-supported. Generic solutions to (1.1) will contain discontinuities, which are much more difficult to analyze, see [4, 32] for

Keywords and phrases. Discontinuous Galerkin method, Runge–Kutta method, error estimates, symmetrizable system of conservation laws, energy analysis.

* *Research of Luo and Zhang is supported by NSFC Grant 11271187, 11071116 and 10931004.*

** *Research of Shu is supported by NSF Grants DMS-1112700 and DMS-1418750, and DOE Grant DE-FG02-08ER25863.*

¹ Department of Mathematics, Nanjing University, Nanjing, 210093, Jiangsu Province, P.R. China. luojuan.nikki@163.com. qzh@nju.edu.cn.

² Division of Applied Mathematics, Brown University, Providence, RI 02912, USA. shu@dam.brown.edu.

some preliminary results in this direction. In this paper we consider only smooth solutions, therefore, we assume the initial solution $\mathbf{u}_0(x)$ is smooth and we consider only a short time interval before discontinuities develop.

The symmetrizable system of conservation laws has a wide background. Well-known examples are the shallow water wave equations and the Euler equations of compressible gas dynamics. As is well-known, a conservation law system (1.1) is symmetrizable if and only if it has a convex entropy function [14]. Due to the symmetrization theory [14, 25], one can seek a mapping $\mathbf{u}(\mathbf{v}): \mathbb{R}^m \rightarrow \mathbb{R}^m$ applied to (1.1), so that when transformed, $\mathbf{u}'_{\mathbf{v}} \mathbf{v}_t + \mathbf{f}'_{\mathbf{v}} \mathbf{v}_x = 0$, the Jacobian matrix $\mathbf{u}'_{\mathbf{v}} = (\partial \mathbf{u}_i / \partial v_j)_{i=1, \dots, m}^{j=1, \dots, m}$ is symmetric positive definite and the Jacobian matrix $\mathbf{f}'_{\mathbf{v}} = \mathbf{f}'_{\mathbf{u}} \mathbf{u}'_{\mathbf{v}}$ is also symmetric.

In this paper we consider the numerical method to solve directly (1.1) by using the RKDG method. The first version of discontinuous Galerkin (DG) method was introduced in 1973 by Reed and Hill [23] for the steady linear neutron transport. It was later developed into RKDG methods by Cockburn *et al.* [5–7, 10, 12] for nonlinear hyperbolic conservation laws, which use a DG discretization in space and combine it with an explicit total variation diminishing Runge–Kutta (TVDRK) time-marching algorithm [26]. Later, this method was developed to solve equations with higher order derivatives. It is well-known that the DG method has strong stability and optimal accuracy to capture discontinuous jumps and/or sharp transient layers, and it combines the advantages of finite element and finite volume methods. For a fairly complete set of references on this methods, we refer to the review papers [9, 11].

However, up to now there has been relatively few work on stability analysis and error estimates for the fully discrete RKDG methods with the explicit TVDRK time marching to solve (1.1). The method of line version (continuous in time) of the DG scheme for linear equations has been considered in [8, 18, 20], and has been proved to maintain good L^2 -norm stability and optimal error estimates. For nonlinear equations, there exists the well-known local entropy inequality [17] for the semi-discrete DG scheme, as well as for the fully discrete DG scheme with some special time-discretizations such as the backward Euler and Crank–Nicolson algorithms. Recently, RKDG method with explicit time-discretization for nonlinear conservation laws has been analyzed in [2, 29, 31], where the (quasi)-optimal *a priori* error estimates are obtained for the second order and the third order explicit TVDRK time discretization, respectively, for scalar equations.

As to symmetrizable systems, error estimates are more difficult to obtain. In [30], error estimates are obtained for RKDG methods with a second order TVDRK time discretization, under a restrictive time step constraint in which $\tau = o(h)$. This restriction is not surprising, as RKDG methods with a second order TVDRK time discretization is linearly unconditionally unstable under the regular time step restriction $\tau \leq ch$ for any fixed constant $c > 0$, when the polynomial degree $k \geq 2$ [9]. However, the time step restriction $\tau = o(h)$ is not realistic in actual computation, and we would like to use third order TVDRK time discretization under regular CFL time step restriction. As far as the authors know, there are still no error estimates for the symmetrizable system, when explicit Runge–Kutta time discretization is used under regular CFL conditions. The RKDG3 method uses the third order explicit TVDRK time-marching combined with piecewise polynomials of arbitrary degree, and is more popular because in practice it provides better linear stability and higher order accuracy in time. However, the techniques used in [30] to obtain error estimates for RKDG2 under the more restrictive time step constraints are not applicable to RKDG3 methods under regular CFL conditions. In a recent work [31], a new technique which explores the specific dissipation natures of third order TVDRK time discretization is developed to analyze the scalar equations. The extension of this technique to symmetrizable systems is the purpose of this paper. This extension is non-trivial, as we must carefully study the numerical fluxes (referred to as the generalized E-flux) in the system case and their influence towards the errors. We can no longer use monotonicity of the fluxes in the scalar case. Our result is a quasi-optimal or optimal error estimate depending on the numerical flux, similar to that in [31] for the scalar case, under the standard temporal-spatial CFL condition $\tau \leq \gamma h$. Here h and τ are the element length and time step, respectively, and the positive constant γ is independent of h and τ .

The main line of analysis in this paper follows that in [30, 31], by using energy analysis as the main tool. It is rather straightforward to generalize the error estimates in [31] from the scalar case to *symmetric* systems. However, as we will see later in this paper, it is significantly more difficult to carry out the above error estimates

for the symmetrizable systems. In this development, we need to pay more attention to several issues owing to the symmetrizable theory, for instance, the suitable norm with local rotational matrix in each element, and Lipschitz continuity of the rotational matrix and the Jacobian matrix of the flux. In this process, many issues about vector-valued functions need to be addressed.

The essential difficulty in this paper is how to describe the numerical viscosity resulting from the DG spatial discretization. We have considered the RKDG2 method in [30], where the explicit second order TVDRK time-marching is used, for the symmetrizable system, and a description on this issue has been given there. However, the stability mechanism in the RKDG3 method is completely different from that of the RKDG2 method, where the additional stability provided by the explicit third order TVDRK time-marching plays an important role. As we have done in Section 3.3 of [31], we need to use the Lipschitz continuity of the numerical viscosity quantity to maintain the structure of this stability term in the time direction. This property is not needed in the analysis for the RKDG method with the second order time-marching, as in [29,30]. This new property demands that we have to modify the analysis in a suitable way for the symmetrizable system.

To this end, we firstly try to propose a careful classification of the necessary properties for the numerical fluxes used in the symmetrizable system case. This kind of numerical flux is defined as the *generalized E-flux* in this paper, where we relax the demand than that in [30]. For example, we only need two inequalities related to the two states under consideration, and we add the freedom on the rotation position and the adjusting matrix. Thus this new definition is more easily verified than that in [30]. Then we establish an important matrix termed *generalized numerical viscosity matrix*, to describe the total numerical stability coming from the jumps at every element interface. In order to ensure the Lipschitz continuity of this matrix, the definition in this paper is different from that in [30]. We carry out this idea through the generalized Newton difference quotients. Furthermore, we present three typical assumptions on the generalized numerical viscosity matrix, which are enough to obtain good error estimates. The details are given in Section 2.2.

An outline of this paper is as follows. In Section 2 we present the RKDG3 scheme for the symmetrizable system of conservation laws. The so-called *generalized E-flux* and the *generalized numerical viscosity matrix* are presented. Some assumptions on the numerical viscosity matrix and the smoothness of the exact solution are also given here, which yield the main conclusion about the quasi-optimal error estimate. In the remaining part of this paper, we would like to present the detailed proof to this main conclusion. In Section 3 we obtain the error representation and the corresponding error equations, and in Section 4 we provide some elementary discussions on the error functional in detail. In Section 5 we carry out the energy analysis and complete the error estimate, with some technical proofs left to the appendix. Finally, a concluding remark is given in Section 6.

2. RKDG3 SCHEME AND THE MAIN CONCLUSION

In this section we present the detailed implementation of the RKDG3 method, following the notations in [30,31]. The generalized numerical viscosity matrix is defined for the so-called generalized E-flux. Finally, we present the main conclusion on the quasi-optimal error estimates in general, and on the optimal error estimate for the upwind numerical flux.

2.1. RKDG3 scheme

Let $\mathcal{J}_h = \{I_j = (x_{j-1/2}, x_{j+1/2})\}_{j=1}^N$ be a partition of $I = (0, 1)$, with each element length being $h_j = x_{j+1/2} - x_{j-1/2}$. The maximum length of this mesh is denoted by $h = \max_j h_j$. In this paper we assume the partition is quasi-uniform, namely, there exists a positive constant ν such that $h \leq \nu h_j$ for all $j = 1, 2, \dots, N$, as h goes to zero. The discontinuous finite element space is defined as

$$\mathbb{V}_h = \{ \mathbf{v} \in [L^2(0, 1)]^m : \mathbf{v}|_{I_j} \in [P^k(I_j)]^m, j = 1, \dots, N \}, \tag{2.1}$$

where $P^k(I_j)$ denotes the space of polynomials in I_j of degree at most k . Note that the function $\mathbf{p} \in \mathbb{V}_h$ is allowed to have discontinuities across element interfaces. Two limits from the left- and the right-directions are

denoted by \mathbf{p}^- and \mathbf{p}^+ , respectively. Furthermore, the jump and the mean, respectively, are denoted by

$$[[\mathbf{p}]] = \mathbf{p}^+ - \mathbf{p}^-, \quad \text{and} \quad \{\{\mathbf{p}\}\} = \frac{1}{2}(\mathbf{p}^+ + \mathbf{p}^-). \tag{2.2}$$

We discretize the time interval $[0, T]$ with the time step τ , which could actually change from step to step but is taken as a constant in this paper for simplicity. In the RKDG3 method, we would like to seek successively the numerical solution, denoted by $\mathbf{u}_h^n(x) = \mathbf{u}_h(x, n\tau)$, in the discontinuous finite element space.

The initial solution \mathbf{u}_h^0 is taken as the approximation of $\mathbf{u}_0(x)$, for instance, the standard L^2 -projection $\mathbb{P}_h \mathbf{u}_0(x)$. It is defined as the unique function in \mathbb{V}_h such that

$$(\mathbb{P}_h \mathbf{u}_0(x) - \mathbf{u}_0(x), \mathbf{v}_h) = 0 \quad \forall \mathbf{v}_h \in \mathbb{V}_h. \tag{2.3}$$

Here $(\mathbf{p}, \mathbf{q}) = \int_I \mathbf{p}^\top \mathbf{q} \, dx$ is the inner product as usual in the space $[L^2(0, 1)]^m$. Note that this projection is locally defined on each element and hence we will also refer to it as the local L^2 -projection later.

For each $n \geq 0$, the numerical solution of the RKDG3 method from the time $n\tau$ to the next time $(n + 1)\tau$ is defined as follows: find $\mathbf{u}_h^{n,1}, \mathbf{u}_h^{n,2}$ and \mathbf{u}_h^{n+1} in the finite element space \mathbb{V}_h , such that for any $\mathbf{v}_h \equiv \mathbf{v}_h(x) \in \mathbb{V}_h$ there hold

$$(\mathbf{u}_h^{n,1}, \mathbf{v}_h) = (\mathbf{u}_h^n, \mathbf{v}_h) + \tau \mathcal{H}(\mathbf{u}_h^n, \mathbf{v}_h), \tag{2.4a}$$

$$(\mathbf{u}_h^{n,2}, \mathbf{v}_h) = \frac{3}{4}(\mathbf{u}_h^n, \mathbf{v}_h) + \frac{1}{4}(\mathbf{u}_h^{n,1}, \mathbf{v}_h) + \frac{\tau}{4} \mathcal{H}(\mathbf{u}_h^{n,1}, \mathbf{v}_h), \tag{2.4b}$$

$$(\mathbf{u}_h^{n+1}, \mathbf{v}_h) = \frac{1}{3}(\mathbf{u}_h^n, \mathbf{v}_h) + \frac{2}{3}(\mathbf{u}_h^{n,2}, \mathbf{v}_h) + \frac{2\tau}{3} \mathcal{H}(\mathbf{u}_h^{n,2}, \mathbf{v}_h). \tag{2.4c}$$

Here $\mathcal{H}(\mathbf{p}, \mathbf{q})$ is the DG spatial discretization, expressed compactly in the form

$$\mathcal{H}(\mathbf{p}, \mathbf{q}) = \sum_{1 \leq j \leq N} \int_{I_j} \mathbf{q}_x^\top \mathbf{f}(\mathbf{p}) \, dx + \sum_{1 \leq j \leq N} [[\mathbf{q}]]_{j+\frac{1}{2}}^\top \hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+)_{j+\frac{1}{2}}, \quad \forall \mathbf{p}, \mathbf{q}, \tag{2.5}$$

since the numerical solution is periodic or compactly-supported.

In (2.5), $\hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+)$ is the numerical flux to ensure the good stability and high-order accuracy. In general, it depends on the two vectors along the left and right directions at the element boundary point. The well-known examples include the global (local) Lax–Friedrichs flux, the Roe linearization flux function [24] with or without Harten’s entropy fix [15], and those fluxes constructed by the flux vector splitting [27]. They can be looked upon as exact or approximate Riemann solvers. For more details, we refer to [27, 30].

To ensure numerical stability, the time step should satisfy a temporal-spatial restriction. The detailed CFL condition will be given later; see (5.1) and (5.25). We have now completed the definition of the RKDG3 method.

2.2. Numerical flux

In this subsection we would like to recall some issues for the numerical fluxes used in practice. Then an important matrix will be defined to describe the numerical viscosity for symmetrizable systems.

2.2.1. Some important matrices in the symmetrizable theory

It follows from the symmetrizable theory that $\mathbf{f}'_{\mathbf{u}}$ has a strong relationship with an important symmetric matrix, since

$$\mathbf{Q}^{1/2} \mathbf{f}'_{\mathbf{u}} \mathbf{Q}^{-1/2} = \mathbf{Q}^{1/2} \mathbf{f}'_{\mathbf{v}} \mathbf{Q}^{1/2} \equiv \mathbf{K}, \tag{2.6}$$

where $\mathbf{Q} \equiv \mathbf{Q}(\mathbf{u}) = \mathbf{v}'_{\mathbf{u}}$ is the Jacobian matrix of the transformation $\mathbf{v} = \mathbf{v}(\mathbf{u})$. Note that \mathbf{K} is a symmetric matrix with the same spectrum as $\mathbf{f}'_{\mathbf{u}}$. Thus there exists the eigenvector decomposition $\mathbf{K} = \mathbf{X}^{-1} \tilde{\mathbf{K}} \mathbf{X}$, where $\tilde{\mathbf{K}} = \text{diag}\{\lambda_i\}_{i=1}^m$. Let $\tilde{\mathbf{K}}^+ = \text{diag}\{\max(\lambda_i, 0)\}_{i=1}^m$ and $\tilde{\mathbf{K}}^- = \text{diag}\{\min(\lambda_i, 0)\}_{i=1}^m$. The positive and negative

parts of \mathbf{K} are then defined as $\mathbf{K}^\pm = \mathbf{X}^{-1}\tilde{\mathbf{K}}^\pm\mathbf{X}$. The standard absolute value matrix of \mathbf{K} is defined as $|\mathbf{K}| = \mathbf{K}^+ - \mathbf{K}^-$.

In this paper, we will use another important matrix and its generalized absolute value matrix, respectively,

$$\mathbf{H} = \mathbf{v}'_u \mathbf{f}'_u = \mathbf{Q}^{1/2} \mathbf{K} \mathbf{Q}^{1/2}, \quad \langle \mathbf{H} \rangle = \mathbf{Q}^{1/2} |\mathbf{K}| \mathbf{Q}^{1/2}. \tag{2.7}$$

Note that $\langle \mathbf{H} \rangle$ is a symmetric positive semidefinite matrix.

Furthermore, in this paper we will use $\|\cdot\|_M$ to represent the length of a vector, or the spectrum norm of a matrix, respectively. Namely, $\|\mathbf{p}\|_M = (\sum_{i=1}^m p_i^2)^{1/2}$ for any vector $\mathbf{p} = (p_1, \dots, p_m)^\top$, and $\|\mathbf{C}\|_M = \max_{\|\mathbf{p}\|_M=1} \|\mathbf{C}\mathbf{p}\|_M$ for any matrix \mathbf{C} . If the matrix \mathbf{C} is symmetric, then $\|\mathbf{C}\|_M$ is equal to the spectral radius of this matrix, denoted by $\rho(\mathbf{C})$. We will also use the following inequalities

$$|\mathbf{p}^\top \mathbf{C} \mathbf{q}| \leq \|\mathbf{C}\|_M \|\mathbf{p}\|_M \|\mathbf{q}\|_M, \quad \left| \mathbf{p}^\top \mathbf{C} \mathbf{q} \right| \leq \left(\mathbf{p}^\top |\mathbf{C}| \mathbf{p} \right)^{1/2} \left(\mathbf{q}^\top |\mathbf{C}| \mathbf{q} \right)^{1/2}, \tag{2.8}$$

for any matrix \mathbf{C} and any vectors \mathbf{p} and \mathbf{q} . Note that $\mathbf{C} = |\mathbf{C}|$ if \mathbf{C} is positive semidefinite. Both inequalities in (2.8) are named Cauchy–Schwarz inequalities in this paper.

2.2.2. Generalized numerical fluxes

Now we would like to recall some elementary properties of numerical fluxes for symmetrizable systems.

In general, the numerical flux $\hat{\mathbf{f}}(\mathbf{a}, \mathbf{b})$ is locally Lipschitz continuous with respect to each argument, and consistent with the true flux $\mathbf{f}(\mathbf{p})$, namely, $\hat{\mathbf{f}}(\mathbf{p}, \mathbf{p}) = \mathbf{f}(\mathbf{p})$. Furthermore, successful fluxes used in practice should satisfy certain conditions to ensure stability and convergence of RKDG schemes.

In this paper we would like to present an abstract framework to describe these conditions, suitable for further analysis. This abstract framework is naturally defined for the scalar case (when $m = 1$) in terms of monotonicity, where the numerical flux $\hat{f}(p^-, p^+)$ is assumed to be a nondecreasing function of its first argument and a nonincreasing function of its second argument. Such fluxes are called monotone fluxes. Following [22], a more general class of fluxes are termed as entropy fluxes (E-fluxes), namely, there always holds for any q between p^- and p^+ , that

$$\llbracket \mathbf{p} \rrbracket (f(q) - \hat{f}(p^-, p^+)) \geq 0. \tag{2.9}$$

However, this description is not trivial to extend to the system case.

For symmetric systems of conservation laws, an important concept along this line is the so-called generalized E-flux, proposed in [16]. That is to say that

$$\llbracket \mathbf{p} \rrbracket^\top \{ \mathbf{f}(\mathbf{r}) - \hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+) \} \geq 0, \tag{2.10}$$

for any \mathbf{r} standing on the straight-line segment with two endpoints \mathbf{p}^\pm . However, this condition can not describe desirable numerical fluxes when solving symmetrizable (but not symmetric) systems of conservation laws.

Thanks to (2.6), we would like in this paper to propose an extension of the generalized E-fluxes to symmetrizable systems, by the help of the local rotation matrix \mathbf{Q} .

Definition 2.1. The numerical flux $\hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+)$ is called a generalized E-flux, if we have a rotation position \mathbf{s}_κ and a positive semidefinite matrix \mathbf{Y}_κ , such that

$$\llbracket \mathbf{p} \rrbracket^\top \mathbf{Q}(\mathbf{s}_\kappa) \{ \mathbf{f}(\mathbf{r}_\kappa) - \hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+) \} + \llbracket \mathbf{p} \rrbracket^\top \mathbf{Y}_\kappa \llbracket \mathbf{p} \rrbracket \geq 0, \quad \kappa = 1, 2, \tag{2.11}$$

for $\mathbf{r}_1 = \mathbf{p}^-$ and $\mathbf{r}_2 = \mathbf{p}^+$. Here $\mathbf{s}_\kappa = \mathbf{s}_\kappa(\mathbf{p}^-, \mathbf{p}^+)$ lies in the standard super-rectangle with two vertices \mathbf{p}^\pm . Furthermore, each element in $\mathbf{Y}_\kappa = \mathbf{Y}_\kappa(\mathbf{p}^-, \mathbf{p}^+)$ has the order $\mathcal{O}(\llbracket \mathbf{p} \rrbracket \|_M)$, with the common bound depending only on the local Lipschitz constant of \mathbf{Q} in the above super-rectangle.

Note that the adjusting matrix \mathbf{Y}_κ is introduced to overcome the trouble in [30] on how to seek a reasonable rotation position \mathbf{s}_κ . It is easy to seek a good rotation position for the scalar case and for symmetric systems; in such cases we can take $\mathbf{Y}_\kappa = 0$. In general, with the help of the new freedom \mathbf{Y}_κ , we can take the rotation position \mathbf{s}_κ as an arbitrary point along the straight line between \mathbf{p}^- and \mathbf{p}^+ . Hence, this definition enlarges the class of generalized E-fluxes, given in [16, 30].

Many numerical fluxes can be verified easily to be generalized E-fluxes. A detailed example will be given in the appendix, for the (global/local) Lax–Friedrichs flux. More examples can be found in [21].

2.2.3. Generalized numerical viscosity matrix

It is well-known that the numerical viscosity in the semi-discrete DG method comes from the square of the jumps at element boundary points. Below we would like to define an important matrix for the generalized E-flux, to describe the strength of the numerical stability in space. Similar quantity has been defined and analyzed for the scalar case; we refer to [29].

To extend to the system case, we would like to use the following notation. Let $\mathbf{g} = (g_1, g_2, \dots, g_m)^\top$ be an m -dimensional vector-valued function with respect to the m -dimensional variable. For any given two vectors $\mathbf{a} = (a_1, a_2, \dots, a_m)^\top$ and $\mathbf{b} = (b_1, b_2, \dots, b_m)^\top$, define $\mathbf{a}^{(0)} = \mathbf{a}$, $\mathbf{a}^{(m)} = \mathbf{b}$, and

$$\mathbf{a}^{(j)} = (b_1, \dots, b_{j-1}, b_j, a_{j+1}, \dots, a_m)^\top, \quad j = 1, 2, \dots, m - 1. \tag{2.12}$$

Note that $\mathbf{g}(\mathbf{a}) - \mathbf{g}(\mathbf{b}) = \sum_{j=1}^m [\mathbf{g}(\mathbf{a}^{(j-1)}) - \mathbf{g}(\mathbf{a}^{(j)})]$. The generalized Newton difference quotient is an m -dimensional matrix, denoted by $\mathcal{D}\mathbf{g}[\mathbf{a}, \mathbf{b}]$, with the element in the (i, j) -entry being

$$\left(\mathcal{D}\mathbf{g}[\mathbf{a}, \mathbf{b}]\right)_{ij} = \frac{\mathbf{g}_i(\mathbf{a}^{(j-1)}) - \mathbf{g}_i(\mathbf{a}^{(j)})}{a_j - b_j}. \tag{2.13}$$

If the denominator is equal to zero, the term should be understood as the limit when the denominator goes to zero.

Definition 2.2. Let $\hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+)$ be a generalized E-flux, locally Lipschitz continuous and consistent with $\mathbf{f}(\mathbf{u})$. The generalized numerical viscosity matrix $\mathcal{A}(\hat{\mathbf{f}}; \mathbf{p})$ is defined at each element boundary point, in the form

$$\begin{aligned} \mathcal{A}(\hat{\mathbf{f}}; \mathbf{p}) &= \frac{1}{2}\mathcal{A}_1(\hat{\mathbf{f}}; \mathbf{p}) + \frac{1}{2}\mathcal{A}_2(\hat{\mathbf{f}}; \mathbf{p}) \\ &\equiv \frac{1}{2}\mathbf{Q}(\mathbf{s}_1)\mathcal{D}\hat{\mathbf{f}}^{(1)}[\mathbf{p}^-, \mathbf{p}^+] - \frac{1}{2}\mathbf{Q}(\mathbf{s}_2)\mathcal{D}\hat{\mathbf{f}}^{(2)}[\mathbf{p}^-, \mathbf{p}^+], \end{aligned} \tag{2.14}$$

where $\hat{\mathbf{f}}^{(1)}(\mathbf{r}) = \hat{\mathbf{f}}(\mathbf{r}, \mathbf{p}^+)$ and $\hat{\mathbf{f}}^{(2)}(\mathbf{r}) = \hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{r})$. Here \mathbf{s}_1 and \mathbf{s}_2 are the two rotation positions, stated in Definition 2.1.

We would like to mention that the above definition confirms the result for the linear constant system, with $\mathbf{f}(\mathbf{u}) = \mathbf{A}\mathbf{u}$ and \mathbf{A} being a constant matrix. In this case, all upwind numerical fluxes used in practice are equal to $\hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+) = \mathbf{A}^+\mathbf{p}^- + \mathbf{A}^-\mathbf{p}^+$. Due to the symmetrizable theory, there exists a symmetric positive definite matrix \mathbf{A}_0 such that $\mathbf{A}\mathbf{A}_0$ is symmetric. A tedious stability analysis [21] for the semi-discrete DG method will give the total numerical viscosity,

$$\frac{1}{2} \sum_{1 \leq j \leq N} \left[\llbracket \mathbf{u}_h \rrbracket^\top \langle \mathbf{A}_0^{-1} \mathbf{A} \rangle \llbracket \mathbf{u}_h \rrbracket \right]_{j+\frac{1}{2}}. \tag{2.15}$$

The generalized numerical viscosity matrix given by (2.14) is the same as the above involved matrix, namely, $\frac{1}{2}(\mathbf{A}_0^{-1} \mathbf{A})$. From this viewpoint, Definition 2.2 is nicer than that in [30].

Through observation on many numerical fluxes, we would like in this paper to make some elementary assumptions on the generalized numerical viscosity matrix, for the convenience of analysis. They read

A1 The matrix $\mathcal{A}(\hat{\mathbf{f}}; \mathbf{p})$ is nearly positive definite, in the sense

$$[\mathbf{p}]^\top \mathcal{A}(\hat{\mathbf{f}}; \mathbf{p}) [\mathbf{p}] \geq -L_{\text{np}} \|\mathbf{p}\|_{\text{M}}^3, \quad \forall \mathbf{p}. \tag{2.16a}$$

A2 The matrix $\mathcal{A}(\hat{\mathbf{f}}; \mathbf{p})$ is Lipschitz continuous, in the sense,

$$\|\mathcal{A}(\hat{\mathbf{f}}; \mathbf{p}) - \mathcal{A}(\hat{\mathbf{f}}; \{\mathbf{p}\})\|_{\text{M}} \leq L_{\text{Lip}} \|\mathbf{p}\|_{\text{M}}, \quad \forall \mathbf{p}. \tag{2.16b}$$

A3 The matrix $\mathcal{A}(\hat{\mathbf{f}}; \mathbf{p})$ is upper consistent with $\frac{1}{2}\langle \mathbf{H}(\{\mathbf{p}\}) \rangle$, in the sense: there exists a positive semidefinite matrix $\mathcal{U} = \mathcal{U}(\mathbf{f}, \hat{\mathbf{f}}; \mathbf{p})$, such that

$$\left| [\mathbf{p}]^\top \left(\frac{1}{2} \langle \mathbf{H}(\{\mathbf{p}\}) \rangle + \mathcal{U} - \mathcal{A}(\hat{\mathbf{f}}; \mathbf{p}) \right) [\mathbf{p}] \right| \leq L_{\text{uc}} \|\mathbf{p}\|_{\text{M}}^3, \quad \forall \mathbf{p}. \tag{2.16c}$$

Note that the above three bounding constants are all independent of $\|\mathbf{p}\|_{\text{M}}$, but may depend on the local Lipschitz constants of the true flux $\mathbf{f}(\mathbf{u})$ and/or the numerical flux $\hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+)$.

The assumption **A1** is provided directly from Definition 2.1. Actually, it follows from (2.11) and (2.14) that

$$\begin{aligned} [\mathbf{p}]^\top \mathcal{A}(\hat{\mathbf{f}}; \mathbf{p}) [\mathbf{p}] &= \frac{1}{2} [\mathbf{p}]^\top \mathbf{Q}(\mathbf{s}_1) \{ \mathbf{f}(\mathbf{p}^+) - \hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+) \} + \frac{1}{2} [\mathbf{p}]^\top \mathbf{Q}(\mathbf{s}_2) \{ \mathbf{f}(\mathbf{p}^-) - \hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+) \} \\ &\geq - [\mathbf{p}]^\top \mathbf{Y}_1 [\mathbf{p}] - [\mathbf{p}]^\top \mathbf{Y}_2 [\mathbf{p}], \end{aligned}$$

where $\|\mathbf{Y}_\kappa\|_{\text{M}} \leq L_\kappa \|\mathbf{p}\|_{\text{M}}$ for $\kappa = 1, 2$. The remaining two assumptions can be verified easily for many numerical fluxes. An example is given in the appendix, for the local/global Lax–Friedrich flux (A.1).

It is worthy to point out that the generalized numerical viscosity matrix satisfies the kernel property

$$\mathcal{A}_\kappa(\hat{\mathbf{f}}; \mathbf{p}) [\mathbf{p}] = \mathbf{Q}(\mathbf{s}_\kappa) \left[\mathbf{f}(\mathbf{r}_\kappa) - \hat{\mathbf{f}}(\mathbf{p}^-, \mathbf{p}^+) \right], \tag{2.17}$$

for $\kappa = 1, 2$, where $\mathbf{r}_1 = \mathbf{p}^-$ and $\mathbf{r}_2 = \mathbf{p}^+$. This property will play an important role in our analysis, as in [30], in order to express explicitly the numerical viscosity provided by the numerical flux. Detailed discussion can be found in section 4.4 and in section A.4 on the proof of Lemma 5.4.

Remark 2.3. In the above discussion, we do not pay attention to whether the generalized numerical viscosity matrix $\mathcal{A}(\hat{\mathbf{f}}; \mathbf{p})$ is symmetric or not.

Remark 2.4. Assumption **A3** can be verified easily for the upwind numerical fluxes (namely, one-sided numerical fluxes for each characteristic variable after suitable local characteristic decomposition, see Sect. 5.6 for more details), for example, the upwind numerical flux for a linear flux and the Steger–Warming flux [27] for Euler equations. In such cases, we can take $\mathcal{U} = 0$, and consequently there holds $\mathcal{A}(\hat{\mathbf{f}}; \mathbf{p})$ tends to $\frac{1}{2}\langle \mathbf{H}(\{\mathbf{p}\}) \rangle$ as $\|\mathbf{p}\|_{\text{M}}$ goes to zero.

2.3. The main conclusion

In this paper, some standard norms will be used. Let \mathbf{p} be a vector-valued function and/or matrix-valued function. The $[L^2(I)]^m$ norm and the infinity norm, respectively, are denoted by

$$\|\mathbf{p}\| = \left(\int_I \|\mathbf{p}(x)\|_{\text{M}}^2 dx \right)^{1/2}, \quad \|\mathbf{p}\|_\infty = \text{ess sup}_{x \in (0,1)} \|\mathbf{p}(x)\|_{\text{M}}. \tag{2.18}$$

We also use $\|\cdot\|_s$ to denote the standard norm in the Sobolev space $[H^s(I)]^m$, in which the function and its derivatives up to sth order are all in $[L^2(I)]^m$.

In order to obtain the *a priori* error estimates to the RKDG3 method, we would like in this paper to assume that the exact solution $\mathbf{u}(x, t)$ to problem (1.1) is sufficiently smooth. Namely,

- H1 Both $\|\mathbf{u}\|_{k+1}(t)$ and $\|\mathbf{u}_{tttt}\|(t)$ are bounded uniformly by a constant for any time $t \in [0, T]$. Furthermore, the exact solution is continuous and bounded. For the purpose of optimal error estimate, we assume that \mathbf{u} has a higher order smoothness, *i.e.*, $\|\mathbf{u}\|_{k+2}(t)$ is also bounded uniformly for any time $t \in [0, T]$.
- H2 Each component of $\mathbf{f}(\mathbf{p})$ and $\mathbf{f}'_{\mathbf{u}}(\mathbf{p})$ is bounded for all $\mathbf{p} \in \mathbb{R}^m$. Moreover, $\mathbf{f}'_{\mathbf{u}}(\mathbf{p})$ is Lipschitz continuous with the bounding constant C_{\star}^{Flux} .
- H3 The rotation matrix $\mathbf{Q}(\mathbf{p})$ is symmetric positive definite uniformly, namely, there exist two constants γ_{\star} and γ^{\star} , such that $0 < \gamma_{\star} \leq \|\mathbf{Q}(\mathbf{p})\|_{\text{M}} \leq \gamma^{\star}$ holds for all $\mathbf{p} \in \mathbb{R}^m$. Furthermore, $\mathbf{Q}(\mathbf{p})$ is Lipschitz continuous with the bounding constant C_{\star}^{Rot} .

As a consequence, the bounding constants in the assumptions A1–A3 can be simplified by three common constants, respectively.

A remark is given here. Due to the boundedness of the exact solution, the last two assumptions H2 and H3 are reasonable with a cut-off modification on $\mathbf{f}(\mathbf{u})$ and $\mathbf{u}(\mathbf{v})$. We refer to [29] for more details.

Now we present the main conclusion in the following theorem.

Theorem 2.5. *Let \mathbf{u}_h be the numerical solution of the RKDG3 scheme (2.4), using piecewise polynomials with arbitrary degree $k \geq 2$, defined on any quasi-uniform triangulations of $I = (0, 1)$. The numerical flux is a generalized E-flux associated with the generalized numerical viscosity matrix satisfying assumptions A1–A3. Let \mathbf{u} be the exact solution of problem (1.1), which satisfies the above smoothness assumptions H1–H3. Then we have the following error estimate*

$$\max_{n\tau \leq T} \|\mathbf{u}(x, t^n) - \mathbf{u}_h^n\| \leq C(h^{k+\sigma} + \tau^3), \tag{2.19}$$

under a standard CFL condition $\tau \leq \gamma h$ with a suitably fixed constant $\gamma > 0$, where the bounding constant $C > 0$ is independent of h and τ . Here $\sigma = \frac{1}{2}$ for generalized E-fluxes; and $\sigma = 1$ for upwind numerical fluxes.

The proof is technical and long, which will be given in several steps in the subsequent sections. To show the main ideas clearly, we focus our attention on the quasi-optimal error estimate ($\sigma = \frac{1}{2}$) for generalized E-fluxes. The optimal error estimate for upwind numerical fluxes can be proved in a similar line; a sketch will be given in Section 5.6.

3. THE ERROR EQUATIONS

To obtain the error estimate, we have to first establish the error equations. This process is similar to that in [30, 31].

3.1. Reference stage solutions

Following [31], three reference functions are defined as the local time discretization of the third order explicit TVDRK algorithm for the exact solution of the conservation law (1.1). Let $\mathbf{u}^{(0)} = \mathbf{u}$, and

$$\mathbf{u}^{(1)} = \mathbf{u}^{(0)} - \tau \mathbf{f}(\mathbf{u}^{(0)})_x, \quad \mathbf{u}^{(2)} = \frac{3}{4} \mathbf{u}^{(0)} + \frac{1}{4} \mathbf{u}^{(1)} - \frac{1}{4} \tau \mathbf{f}(\mathbf{u}^{(1)})_x. \tag{3.1}$$

Owing to the Sobolev embedding theory [1], it follows from the smoothness assumptions H1–H3 that the above reference values are all continuous and bounded in the whole domain $I \times [0, T]$.

A simple manipulation yields the local truncation error in time, denoted by $\mathcal{E}(x, t)$, such that

$$\mathbf{u}^{(0)}(x, t + \tau) = \frac{1}{3} \mathbf{u}^{(0)}(x, t) + \frac{2}{3} \mathbf{u}^{(2)}(x, t) - \frac{2}{3} \tau [\mathbf{f}(\mathbf{u}^{(2)}(x, t))]_x + \mathcal{E}(x, t). \tag{3.2}$$

There holds uniformly $\|\mathcal{E}(x, t)\| = \mathcal{O}(\tau^4)$ for any time, due to assumption H1. Similar discussion can be found in [31] for the scalar case.

Define $\mathbf{u}^{n,\ell} = \mathbf{u}^{(\ell)}(x, t^n)$ to be the reference stage solution. Here and below we may drop the index ℓ if $\ell = 0$. For convenience of notations, we denote three differences

$$\mathbb{E}_1^n \mathbf{u}^n = \mathbf{u}^{n,1} - \mathbf{u}^n, \quad \mathbb{E}_2^n \mathbf{u}^n = 4\mathbf{u}^{n,2} - \mathbf{u}^{n,1} - 3\mathbf{u}^n, \quad \mathbb{E}_3^n \mathbf{u}^n = \frac{1}{2}(3\mathbf{u}^{n+1} - 2\mathbf{u}^{n,2} - \mathbf{u}^n), \quad (3.3)$$

to describe the evolution of the solution at each time stage.

Multiply the test function $\mathbf{v}_h \in \mathbb{V}_h$ on both sides of equations (3.1) and (3.2), respectively. Let $t = t^n$ and then integrate them in each element. Due to the consistency of the numerical flux, this process yields a set of equalities for $\ell = 0, 1, 2$,

$$(\mathbb{E}_{\ell+1}^n \mathbf{u}^n, \mathbf{v}_h) = \tau \mathcal{H}^{n,\ell}(\mathbf{u}^{n,\ell}, \mathbf{v}_h) + (\mathcal{E}^{n,\ell}, \mathbf{v}_h), \quad \forall \mathbf{v}_h \in \mathbb{V}_h, \quad (3.4)$$

which are almost the same as the RKDG3 scheme (2.4). In this paper we would like to denote $\mathcal{E}^{n,2} = \mathcal{E}(x, t^n)$, and denote $\mathcal{E}^n = \mathcal{E}^{n,1} = 0$ just for simplicity.

3.2. Projection properties

Define the stage error $\mathbf{e}^{n,\ell} = \mathbf{u}^{n,\ell} - \mathbf{u}_h^{n,\ell}$. As the usual treatment in a finite element analysis, we divide the stage error in the form

$$\mathbf{e}^{n,\ell} = [\pi_h \mathbf{u}^{n,\ell} - \mathbf{u}_h^{n,\ell}] - [\pi_h \mathbf{u}^{n,\ell} - \mathbf{u}^{n,\ell}] \equiv \boldsymbol{\xi}^{n,\ell} - \boldsymbol{\eta}^{n,\ell}, \quad (3.5)$$

where π_h is a suitable projection. The projection error $\boldsymbol{\eta}^{n,\ell}$ can be estimated easily.

It is enough to take $\pi_h = \mathbb{P}_h$ as the standard L^2 -projection (refer to (2.3)), in order to obtain the quasi-optimal error estimate. Since $\mathbf{u}(x, t)$ is smooth enough, a standard scaling argument [3, 20] yields that the projection error $\boldsymbol{\eta}^{n,\ell}$ satisfies

$$\|\boldsymbol{\eta}^{n,\ell}\| + h\|\boldsymbol{\eta}_x^{n,\ell}\| + h^{1/2}\|\boldsymbol{\eta}^{n,\ell}\|_{\Gamma_h} \leq C_1 h^{k+1}, \quad \forall n: n\tau \leq T, \ell = 0, 1, 2. \quad (3.6a)$$

Here Γ_h is the union of all element interfaces, and, for any function \mathbf{p} , we denote

$$\|\mathbf{p}\|_{\Gamma_h} = \left(\sum_{1 \leq j \leq N} \frac{1}{2} \|\mathbf{p}_{j+\frac{1}{2}}^+\|_{\mathbb{M}}^2 + \frac{1}{2} \|\mathbf{p}_{j+\frac{1}{2}}^-\|_{\mathbb{M}}^2 \right)^{1/2}.$$

It follows from the interpolation theory [3] that

$$\|\boldsymbol{\eta}^{n,\ell}\|_{\infty} \leq C_2 h^{k+\frac{1}{2}}, \quad \forall n: n\tau \leq T, \ell = 0, 1, 2. \quad (3.6b)$$

Since the projection \mathbb{P}_h is linear and independent of the time, $\mathbb{E}_\ell^n \boldsymbol{\eta}^n$ is also the projection error of $\mathbb{E}_\ell^n \mathbf{u}^n$ under the same projection. Thus we also have the estimates

$$\|\mathbb{E}_\ell^n \boldsymbol{\eta}^n\| + h^{1/2}\|\mathbb{E}_\ell^n \boldsymbol{\eta}^n\|_{\Gamma_h} \leq C_3 h^{k+1} \tau, \quad \forall n: n\tau \leq T, \ell = 1, 2, 3. \quad (3.6c)$$

Note that the above bounding constants, C_1, C_2 and C_3 , depend solely on the smoothness of the exact solution.

3.3. The error equations

The remaining work in this paper is to estimate the errors in the finite element space, namely $\boldsymbol{\xi}^{n,\ell} = \pi_h \mathbf{e}^{n,\ell} \in \mathbb{V}_h$. To this end, we need to set up the error equations as follows. Subtracting (3.4) from (2.4) gives the error equations for $\ell = 0, 1, 2$,

$$\begin{aligned} & (\mathbb{E}_{\ell+1}^n \boldsymbol{\xi}^n, \mathbf{v}_h) = \tau \mathcal{K}^{n,\ell}(\mathbf{v}_h) \\ & \equiv (\mathbb{E}_{\ell+1}^n \boldsymbol{\eta}^n + \mathcal{E}^{n,\ell}, \mathbf{v}_h) + \tau \mathcal{H}(\mathbf{u}^{n,\ell}, \mathbf{v}_h) - \tau \mathcal{H}(\mathbf{u}_h^{n,\ell}, \mathbf{v}_h), \quad \forall \mathbf{v}_h \in \mathbb{V}_h. \end{aligned} \quad (3.7)$$

Here $\mathcal{K}^{n,\ell}(v_h)$ is called the error functional at each time stage. To state it more clearly, we would like to separate it into four parts with different meanings, through a simple re-arranging process and using consistency of the numerical flux. Namely,

$$\mathcal{K}^{n,\ell}(v_h) \equiv \mathcal{L}^{n,\ell}(\mathbf{e}^{n,\ell}, \mathbf{v}_h) + \mathcal{N}^{n,\ell}(v_h) + \mathcal{V}^{n,\ell}(v_h) + \mathcal{T}^{n,\ell}(v_h). \quad (3.8)$$

The terms on the right-hand side are named the linear part, the nonlinear part, the viscosity part, and the time-marching part, respectively, which are defined below.

The linear part is defined as a bilinear functional with respect to \mathbf{w} and \mathbf{v}_h , with the given Jacobian matrix $\mathbf{f}'_{\mathbf{u}}(\mathbf{u}^{n,\ell})$. It reads in the form

$$\begin{aligned} \mathcal{L}^{n,\ell}(\mathbf{w}, \mathbf{v}_h) &= \mathcal{L}_{\text{main}}(\mathbf{f}'_{\mathbf{u}}(\mathbf{u}^{n,\ell}); \mathbf{w}, \mathbf{v}_h) \\ &= \sum_{1 \leq j \leq N} \left[\llbracket \mathbf{v}_h \rrbracket_{j+\frac{1}{2}}^\top \mathbf{f}'_{\mathbf{u}}(\mathbf{u}^{n,\ell}) \{\!\! \{ \mathbf{w} \}\!\!\} \right]_{j+\frac{1}{2}} + \int_I (\mathbf{v}_h)_x^\top \mathbf{f}'_{\mathbf{u}}(\mathbf{u}^{n,\ell}) \mathbf{w} \, dx. \end{aligned} \quad (3.9a)$$

Here and below the integration on I means the sum of those integrations on every elements I_j . The remaining three parts are linear functionals with respect to \mathbf{v}_h . The nonlinear part is given in the form

$$\begin{aligned} \mathcal{N}^{n,\ell}(v_h) &= \sum_{1 \leq j \leq N} \left[\llbracket \mathbf{v}_h \rrbracket_{j+\frac{1}{2}}^\top \left[\mathbf{f}(\mathbf{u}^{n,\ell}) - \{\!\! \{ \mathbf{f}(\mathbf{u}_h^{n,\ell}) \}\!\!\} - \mathbf{f}'_{\mathbf{u}}(\mathbf{u}^{n,\ell}) \{\!\! \{ \mathbf{e}^{n,\ell} \}\!\!\} \right] \right]_{j+\frac{1}{2}} \\ &\quad + \int_I (\mathbf{v}_h)_x^\top \left[\mathbf{f}(\mathbf{u}^{n,\ell}) - \mathbf{f}(\mathbf{u}_h^{n,\ell}) - \mathbf{f}'_{\mathbf{u}}(\mathbf{u}^{n,\ell}) \mathbf{e}^{n,\ell} \right] \, dx. \end{aligned} \quad (3.9b)$$

If $\mathbf{f}(\mathbf{u})$ is a linear flux, the nonlinear part disappears. Further, the viscosity part solely depends on the numerical solutions, defined as

$$\mathcal{V}^{n,\ell}(v_h) = \sum_{1 \leq j \leq N} \left[\llbracket \mathbf{v}_h \rrbracket_{j+\frac{1}{2}}^\top \left(\{\!\! \{ \mathbf{f}(\mathbf{u}_h^{n,\ell}) \}\!\!\} - \hat{\mathbf{f}}(\mathbf{u}_h^{n,\ell}) \right) \right]_{j+\frac{1}{2}}. \quad (3.9c)$$

The time-marching part is defined as

$$\mathcal{T}^{n,\ell}(v_h) = \tau^{-1}(\zeta^{n,\ell}, v_h), \quad (3.9d)$$

where $\zeta^{n,\ell} = \mathbb{E}_{\ell+1}^n \boldsymbol{\eta}^n + \mathcal{E}^{n,\ell}$, representing the evolution of the projection error and the local truncation error in time.

Remark 3.1. In the splitting of the error (3.8) above, we have introduced the reference vector at each element boundary point, which depends on the specific projection being used. For the local L^2 -projection \mathbb{P}_h , the reference vector is taken as the simple average $\{\!\! \{ \mathbf{u}_h^{n,\ell} \}\!\!\}$. This corresponds to three terms, namely, $\{\!\! \{ \mathbf{w} \}\!\!\}$ in (3.9a), and $\{\!\! \{ \mathbf{f}(\mathbf{u}_h^{n,\ell}) \}\!\!\}$ in both (3.9b) and (3.9c).

To obtain the optimal error estimate for the upwind numerical flux, we will use the local Gauss–Radau projection. In this case, we would need to introduce a different reference vector at each element boundary point, and make corresponding modifications; see Section 5.6 for more details.

4. ELEMENTARY ESTIMATES

In this section we would like to set up some basic discussions on the error functional. For notational convenience, in this section we will drop the super-index, n and ℓ , for the reference stage solutions, the numerical solutions and the operators.

4.1. Preliminaries

4.1.1. Notations

In this paper we will use notations C, K, ε to denote generic positive constants independent of h, τ and n . Here ε is a small positive constant, and K depends solely on the inverse constants, to be specified in Section 4.1.3 below. To emphasize the nonlinearity of the flux $\mathbf{f}(\mathbf{u})$ and the transformation $\mathbf{u}(\mathbf{v})$, we would like to use C_\star to denote a nonnegative constant which vanishes, *i.e.* $C_\star = 0$ for a linear flux $\mathbf{f}(\mathbf{u}) = \mathbf{A}\mathbf{u}$. Furthermore, for any function \mathbf{p} , we would like to use a macro notation

$$C(\mathbf{p}) = C + C_\star h^{-1} \|\mathbf{p}\|_\infty + C_\star h^{-2} \|\mathbf{p}\|_\infty^2. \quad (4.1)$$

If the bounding constant depends on ε , we will denote it by $C(\varepsilon; \mathbf{p})$. Note that the above notations may have a different value in each occurrence.

In this paper we use the notations \mathbf{u}_c and \mathbf{u}_b to denote the evaluation of the reference solutions at the element center points and the element boundary points, respectively. Recalling $\mathbf{Q} = \mathbf{v}'_{\mathbf{u}}$, we define the piecewise constant rotation matrix as

$$\mathbf{Q}_c = \mathbf{Q}(\mathbf{u}_c), \quad (4.2)$$

which is equal to the evaluation of the matrix at the center point in each element. Also we denote $\mathbf{Q}_b = \mathbf{Q}(\mathbf{u}_b)$ at every element boundary point.

4.1.2. Local focus shifting

Let \mathbf{C} be one of the matrix-valued functions $\mathbf{f}'_{\mathbf{u}}, \mathbf{Q}, \mathbf{Q}^{1/2}, \mathbf{H}$ and $\langle \mathbf{H} \rangle$. Their focus shifting (*i.e.* change of the vector at which this matrix-valued function is evaluated) causes many complexity in the following analysis. Due to the Lipschitz continuity and the famous Wielandt–Hoffman Theorem [13], we have

$$\|\mathbf{C}(\mathbf{a}) - \mathbf{C}(\mathbf{b})\|_{\mathbf{M}} \leq C_\star \|\mathbf{a} - \mathbf{b}\|_{\mathbf{M}}, \quad (4.3)$$

where \mathbf{a} and \mathbf{b} are two considered focuses. Furthermore, we will use the inequalities $\|\mathbf{f}(\mathbf{a}) - \mathbf{f}(\mathbf{b})\|_{\mathbf{M}} \leq C \|\mathbf{a} - \mathbf{b}\|_{\mathbf{M}}$ and

$$\|\mathbf{f}(\mathbf{a}) - \mathbf{f}(\mathbf{b}) - \mathbf{f}'_{\mathbf{u}}(\mathbf{b})(\mathbf{a} - \mathbf{b})\|_{\mathbf{M}} \leq \frac{1}{2} C_\star^{\text{Flux}} \|\mathbf{a} - \mathbf{b}\|_{\mathbf{M}}^2, \quad (4.4)$$

in the local linearization process. We refer to [19] for more details.

The focuses considered in this paper will be taken from the reference solutions $\mathbf{u}_c, \mathbf{u}_b$ and $\mathbf{u}(x)$, or the numerical solutions $\{\{\mathbf{u}_h\}\}$ and $\mathbf{s}_i = \mathbf{s}_i(\mathbf{u}_h)$ in the definition of generalized E-flux, in the same cell or in adjacent cells, at the same time stage level, or at different time stage levels in the same time-marching step. In all these cases, we have

$$\|\mathbf{a} - \mathbf{b}\|_{\mathbf{M}} \leq C(h + \tau) + \|\llbracket \mathbf{e} \rrbracket\|_{\mathbf{M}} \leq C(h + \tau) + \|\mathbf{e}\|_\infty, \quad (4.5)$$

since the reference solution \mathbf{u} is Lipschitz continuous and thus

$$\llbracket \mathbf{u}_h \rrbracket = \llbracket \boldsymbol{\eta} \rrbracket - \llbracket \boldsymbol{\xi} \rrbracket = -\llbracket \mathbf{e} \rrbracket, \quad (4.6)$$

at each element boundary point. Note that (4.5) simply gives a very crude bound. In many cases, we do not have all the terms on the right-hand side of this inequality.

4.1.3. The inverse properties

Some inverse properties will be used in this paper. For any function $\mathbf{v}_h \in \mathbb{V}_h$, there exists a positive constant μ , independent of \mathbf{v}_h and h , such that

$$(i) \ \|(\mathbf{v}_h)_x\| \leq \mu h^{-1} \|\mathbf{v}_h\|; \quad (ii) \ \|\mathbf{v}_h\|_{\Gamma_h} \leq \mu^{1/2} h^{-1/2} \|\mathbf{v}_h\|; \quad (iii) \ \|\mathbf{v}_h\|_\infty \leq \mu h^{-1/2} \|\mathbf{v}_h\|.$$

For more details of these inverse properties, we refer the reader to [3].

To cope with the local moving of focus near the element boundary, we will use the following inequalities

$$\|\mathbf{Q}_b^{1/2}\mathbf{v}_h\|_{\Gamma_h} \leq K\mu^{1/2}\left[h^{-1/2}\|\mathbf{Q}_c^{1/2}\mathbf{v}_h\| + C_\star h^{1/2}\|\mathbf{v}_h\|\right], \quad (4.7a)$$

$$\|\mathbf{Q}^{1/2}(\{\mathbf{u}_h\})\mathbf{v}_h\|_{\Gamma_h}^2 \leq K\|\mathbf{Q}_b^{1/2}\mathbf{v}_h\|_{\Gamma_h}^2 + C_\star h^{-1}\|e\|_\infty\|\mathbf{v}_h\|^2, \quad (4.7b)$$

for any $\mathbf{v}_h \in \mathbb{V}_h$. These can be obtained from the application of focus shifting and the elementary inverse properties.

4.2. The linear part

Starting from this subsection we will discuss separately the four parts in the error functional. Let $S_{\max} = \max \varrho(\mathbf{f}'_u(\mathbf{u}))$, representing the maximum of the flow speed; here $\varrho(\mathbf{f}'_u)$ is the spectral radius, and the maximum is taken over all vectors in \mathbb{R}^m .

Lemma 4.1. *The linear part is a bounded bilinear functional in $\mathbb{V}_h \times \mathbb{V}_h$. Namely, there exist bounding constants K and C_\star independent of h and τ , such that*

$$|\mathcal{L}(\mathbf{w}_h, \mathbf{Q}_c\mathbf{v}_h)| \leq KS_{\max}\mu h^{-1}\|\mathbf{Q}_c^{1/2}\mathbf{w}_h\|\|\mathbf{Q}_c^{1/2}\mathbf{v}_h\| + C_\star\|\mathbf{w}_h\|\|\mathbf{v}_h\| \quad (4.8)$$

holds for any $\mathbf{w}_h, \mathbf{v}_h \in \mathbb{V}_h$.

Proof. For any point x in the interior of each element, we move the focus of the rotation matrix \mathbf{Q} from \mathbf{u}_c to $\mathbf{u}(x)$. Then from (3.9a) we get a splitting $\mathcal{L}(\mathbf{w}_h, \mathbf{Q}_c\mathbf{v}_h) = \mathcal{L}^{(1)}(\mathbf{w}_h, \mathbf{v}_h) + \mathcal{L}^{(2)}(\mathbf{w}_h, \mathbf{v}_h)$, where

$$\begin{aligned} \mathcal{L}^{(1)}(\mathbf{w}_h, \mathbf{v}_h) &= \int_I (\mathbf{v}_h)_x^\top \mathbf{H}(\mathbf{u})\mathbf{w}_h \, dx + \sum_{1 \leq j \leq N} \left[[(\mathbf{v}_h)_x]^\top \mathbf{H}(\mathbf{u}_b) \{\mathbf{w}_h\} \right]_{j+\frac{1}{2}}, \\ \mathcal{L}^{(2)}(\mathbf{w}_h, \mathbf{v}_h) &= \int_I (\mathbf{v}_h)_x^\top \mathbf{M}(\mathbf{u})\mathbf{f}'_u(\mathbf{u})\mathbf{w}_h \, dx + \sum_{1 \leq j \leq N} \left[[(\mathbf{M}(\mathbf{u}_b)\mathbf{v}_h)^\top \mathbf{f}'_u(\mathbf{u}_b) \{\mathbf{w}_h\}] \right]_{j+\frac{1}{2}}, \end{aligned}$$

where $\mathbf{M}(\mathbf{u}) = \mathbf{Q}(\mathbf{u}_c) - \mathbf{Q}(\mathbf{u})$. Note that $\mathbf{M}(\mathbf{u})$ may be discontinuous at the element interface, and hence $\mathbf{M}(\mathbf{u}_b)$ has two limits from different directions.

Now we consider the first term $\mathcal{L}^{(1)}(\mathbf{w}_h, \mathbf{v}_h)$. The two terms involved here are denoted by $\mathcal{L}_{\text{int}}^{(1)}$ and $\mathcal{L}_{\text{bry}}^{(1)}$, respectively. Recall that $\mathbf{H}(\mathbf{u}) = \mathbf{Q}^{1/2}(\mathbf{u})\mathbf{K}(\mathbf{u})\mathbf{Q}^{1/2}(\mathbf{u})$, and $\mathbf{K}(\mathbf{u})$ is a symmetric matrix satisfying $\|\mathbf{K}(\mathbf{u})\|_{\text{M}} = \varrho(\mathbf{f}'_u(\mathbf{u}))$. Noting (2.8), a simple application of the inverse property (i) yields that

$$\begin{aligned} |\mathcal{L}_{\text{int}}^{(1)}| &\leq \left| \int_I (\mathbf{v}_h)_x \mathbf{H}(\mathbf{u}_c)\mathbf{w}_h \, dx \right| + \left| \int_I (\mathbf{v}_h)_x [\mathbf{H}(\mathbf{u}) - \mathbf{H}(\mathbf{u}_c)]\mathbf{w}_h \, dx \right| \\ &\leq \|\mathbf{K}(\mathbf{u}_c)\|_\infty \|(\mathbf{Q}_c^{1/2}\mathbf{v}_h)_x\| \|\mathbf{Q}_c^{1/2}\mathbf{w}_h\| + \|\mathbf{H}(\mathbf{u}) - \mathbf{H}(\mathbf{u}_c)\|_\infty \|(\mathbf{v}_h)_x\| \|\mathbf{w}_h\| \\ &\leq S_{\max}\mu h^{-1} \|\mathbf{Q}_c^{1/2}\mathbf{w}_h\| \|\mathbf{Q}_c^{1/2}\mathbf{v}_h\| + C_\star \|\mathbf{v}_h\| \|\mathbf{w}_h\|, \end{aligned} \quad (4.10)$$

since $\|\mathbf{H}(\mathbf{u}) - \mathbf{H}(\mathbf{u}_c)\|_{\text{M}} \leq C_\star h$ holds for all x , due to (4.3) and the smoothness assumption H1 on \mathbf{u} . Also we can use (4.7a) and the boundedness of \mathbf{Q}_c (see assumption H3), to get

$$\begin{aligned} |\mathcal{L}_{\text{bry}}^{(1)}| &\leq S_{\max} \|\mathbf{Q}_b^{1/2} [(\mathbf{v}_h)_x]_{\Gamma_h}\| \|\mathbf{Q}_b^{1/2} \{\mathbf{w}_h\}\|_{\Gamma_h} \\ &\leq KS_{\max}\mu h^{-1} \|\mathbf{Q}_c^{1/2}\mathbf{w}_h\| \|\mathbf{Q}_c^{1/2}\mathbf{v}_h\| + C_\star \|\mathbf{v}_h\| \|\mathbf{w}_h\|. \end{aligned} \quad (4.11)$$

Since $\|\mathbf{M}(\mathbf{u})\|_{\text{M}} = \mathcal{O}(h)$ also holds for all x (as we have discussed on the focus shifting), using Cauchy-Schwarz inequality, as well as the inverse properties (i) and (ii), we achieve

$$|\mathcal{L}^{(2)}(\mathbf{w}_h, \mathbf{v}_h)| \leq C_\star \|\mathbf{w}_h\| \|\mathbf{v}_h\|. \quad (4.12)$$

Hence, collecting up the above inequalities completes the proof of this lemma. \square

In this paper the following conclusion plays an important role. The two inequalities reflect the approximate skew-symmetric and negative semidefinite properties, respectively.

Lemma 4.2. *For any functions \mathbf{w}_h and \mathbf{v}_h in \mathbb{V}_h , there holds*

$$|\mathcal{L}(\mathbf{w}_h, \mathbf{Q}_c \mathbf{v}_h) + \mathcal{L}(\mathbf{v}_h, \mathbf{Q}_c \mathbf{w}_h)| \leq C_* \|\mathbf{w}_h\| \|\mathbf{v}_h\|. \quad (4.13)$$

As a corollary, we have $|\mathcal{L}(\mathbf{v}_h, \mathbf{Q}_c \mathbf{v}_h)| \leq C_* \|\mathbf{v}_h\|^2$.

Proof. We just need to estimate $\mathcal{L}^{(1)}(\mathbf{w}_h, \mathbf{v}_h)$ again, which is defined in the proof of Lemma 4.1. By $\mathbf{H}_{\text{prj}}(\mathbf{u})$ we denote the piecewise linear interpolation of $\mathbf{H}(\mathbf{u})$. Since this matrix is symmetric, we can have

$$(\mathbf{w}_h)_x^\top \mathbf{H}_{\text{prj}} \mathbf{v}_h + (\mathbf{v}_h)_x^\top \mathbf{H}_{\text{prj}} \mathbf{w}_h = \left[\mathbf{w}_h^\top \mathbf{H}_{\text{prj}} \mathbf{v}_h \right]_x - \mathbf{w}_h^\top (\mathbf{H}_{\text{prj}})_x \mathbf{v}_h, \quad (4.14)$$

where $(\mathbf{H}_{\text{prj}})_x|_{I_j} = [\mathbf{H}(\mathbf{u}_{j+\frac{1}{2}}) - \mathbf{H}(\mathbf{u}_{j-\frac{1}{2}})]/[x_{j+\frac{1}{2}} - x_{j-\frac{1}{2}}]$. As we have discussed in Section 4.1.2, we know that $\|(\mathbf{H}_{\text{prj}})_x\|_{\text{M}}$ is bounded for all x .

Since the solution is periodic or compactly-supported, we can solve the integration in each element and obtain

$$\mathcal{L}^{(1)}(\mathbf{w}_h, \mathbf{v}_h) + \mathcal{L}^{(1)}(\mathbf{v}_h, \mathbf{w}_h) = \Xi_{\text{bry}} + \Xi_{\text{int}}, \quad (4.15)$$

where

$$\begin{aligned} \Xi_{\text{bry}} &= \sum_{1 \leq j \leq N} \left[\left[\mathbf{w}_h \right]^\top \mathbf{H}(\mathbf{u}_b) \{\mathbf{v}_h\} + \left[\mathbf{v}_h \right]^\top \mathbf{H}(\mathbf{u}_b) \{\mathbf{w}_h\} - \left[(\mathbf{w}_h)^\top \mathbf{H}(\mathbf{u}_b) \mathbf{v}_h \right] \right]_{j+\frac{1}{2}}, \\ \Xi_{\text{int}} &= \int_I \left\{ (\mathbf{w}_h)_x^\top \left[\mathbf{H} - \mathbf{H}_{\text{prj}} \right] \mathbf{v}_h + (\mathbf{v}_h)_x^\top \left[\mathbf{H} - \mathbf{H}_{\text{prj}} \right] \mathbf{w}_h - \mathbf{w}_h^\top (\mathbf{H}_{\text{prj}})_x \mathbf{v}_h \right\} dx. \end{aligned}$$

It is easy to see that $\Xi_{\text{bry}} = 0$, since the term in the bracket is equal to zero. Further, we have $|\Xi_{\text{int}}| \leq C_* \|\mathbf{w}_h\| \|\mathbf{v}_h\|$ by the inverse property (i), since the former two matrices involved here are of the size $\mathcal{O}(h)$. Finally, by collecting up the above analysis and (4.12) we complete the proof of this lemma. \square

In this paper we would like to use the following compact notations to denote the collection of jumps on the element boundary, like

$$\llbracket \mathbf{p} \rrbracket^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle \llbracket \mathbf{q} \rrbracket = \sum_{1 \leq j \leq N} \left[\llbracket \mathbf{p} \rrbracket^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle \llbracket \mathbf{q} \rrbracket \right]_{j+\frac{1}{2}}, \quad \forall \mathbf{p}, \forall \mathbf{q}. \quad (4.17)$$

This kind of compact notations will also be used to the matrices \mathcal{A}, \mathcal{S} , etc.

Lemma 4.3. *There exists a constant C independent of h and τ , such that*

$$|\mathcal{L}(\boldsymbol{\eta}, \mathbf{Q}_c \mathbf{v}_h)| \leq \varepsilon \llbracket \mathbf{v}_h \rrbracket^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle \llbracket \mathbf{v}_h \rrbracket + \varepsilon \|\mathbf{v}_h\|^2 + C \varepsilon^{-1} h^{2k+1}, \quad \forall \mathbf{v}_h \in \mathbb{V}_h. \quad (4.18)$$

Here ε is any given positive constant.

Proof. Similar to what we have done in Lemma 4.1, we move the focus of the rotation matrix in definition (3.9a). This yields $\mathcal{L}(\boldsymbol{\eta}, \mathbf{Q}_c \mathbf{v}_h) = \mathcal{L}^{(3)}(\mathbf{v}_h) + \mathcal{L}^{(4)}(\mathbf{v}_h)$, where

$$\begin{aligned} \mathcal{L}^{(3)}(\mathbf{v}_h) &= \sum_{1 \leq j \leq N} \left[\left[\mathbf{v}_h \right]^\top \mathbf{H}(\mathbf{u}_b) \{\boldsymbol{\eta}\} \right]_{j+\frac{1}{2}}, \\ \mathcal{L}^{(4)}(\mathbf{v}_h) &= \int_I (\mathbf{v}_h)_x^\top \mathbf{Q}_c \left[\mathbf{f}'_u(\mathbf{u}) - \mathbf{f}'_u(\mathbf{u}_c) \right] \boldsymbol{\eta} dx + \sum_{1 \leq j \leq N} \left[\llbracket \mathbf{M}_b \mathbf{v}_h \rrbracket^\top \mathbf{f}'_u(\mathbf{u}_b) \{\boldsymbol{\eta}\} \right]_{j+\frac{1}{2}}. \end{aligned}$$

The first term is bounded by Cauchy–Schwarz inequality (2.8) and Young’s inequality. The projection property (3.6a) yields that

$$|\mathcal{L}^{(3)}(\mathbf{v}_h)| \leq \varepsilon \llbracket \mathbf{v}_h \rrbracket^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle \llbracket \mathbf{v}_h \rrbracket + C\varepsilon^{-1}h^{2k+1}. \tag{4.20}$$

Since $\|\mathbf{f}'_{\mathbf{u}}(\mathbf{u}) - \mathbf{f}'_{\mathbf{u}}(\mathbf{u}_c)\|_M = \mathcal{O}(h)$ holds for all x , as we have discussed in Section 4.1.2, we can bound the second term in the form

$$|\mathcal{L}^{(4)}(\mathbf{v}_h)| \leq C \left[\|\mathbf{v}_h\|^2 + \|\boldsymbol{\eta}\|^2 + h\|\boldsymbol{\eta}\|_{\Gamma_h}^2 \right] \leq \varepsilon \|\mathbf{v}_h\|^2 + C\varepsilon^{-1}h^{2k+1}, \tag{4.21}$$

using the inverse properties (i) and (ii), as well as the projection property (3.6a). This completes the proof of this lemma. \square

Remark 4.4. In the above discussion, the average $\{\{\boldsymbol{\eta}\}\}$ comes from the setting of the reference vector and the local L^2 -projection. When the local Gauss–Radau projection is used to obtain optimal error estimate for upwind numerical fluxes, this term will be changed into the new setting along the upwind direction. See Section 5.6 for details.

4.3. The nonlinear part

Now we turn to show that the nonlinear part does not cause serious trouble in the error estimate. The conclusion is stated in the following lemma.

Lemma 4.5. *There exists a constant $C_\star \geq 0$ independent of h and τ , such that*

$$|\mathcal{N}(\mathbf{Q}_c \mathbf{v}_h)| \leq C_\star \|\mathbf{v}_h\|^2 + C_\star h^{-2} \|e\|_\infty^2 \left[\|\boldsymbol{\xi}\|^2 + h^{2k+2} \right], \quad \forall \mathbf{v}_h \in \mathbb{V}_h. \tag{4.22}$$

Proof. Noticing (4.4), we can bound the nonlinear term $\mathcal{N}(\mathbf{Q}_c \mathbf{v}_h)$ in the form

$$\begin{aligned} |\mathcal{N}(\mathbf{Q}_c \mathbf{v}_h)| &\leq C_\star \sum_{1 \leq j \leq N} \left[\left(\|\llbracket \mathbf{v}_h \rrbracket\|_M \|\{\{e\}\}\|_M^2 \right)_{j+\frac{1}{2}} + \int_{I_j} \|(\mathbf{v}_h)_x\|_M \|e\|_M^2 dx \right] \\ &\leq C_\star \|e\|_\infty \left[\|\llbracket \mathbf{v}_h \rrbracket\|_{\Gamma_h} \|\llbracket e \rrbracket\|_{\Gamma_h} + \|(\mathbf{v}_h)_x\| \|e\| \right] \\ &\leq C_\star \|e\|_\infty \left[\mu^{1/2} h^{-1/2} \|\mathbf{v}_h\| \left[\mu^{1/2} h^{-1/2} \|\boldsymbol{\xi}\| + \|\boldsymbol{\eta}\|_{\Gamma_h} \right] + \mu h^{-1} \|\mathbf{v}_h\| (\|\boldsymbol{\xi}\| + \|\boldsymbol{\eta}\|) \right] \\ &\leq C_\star h^{-1} \|e\|_\infty \|\mathbf{v}_h\| \left[\|\boldsymbol{\xi}\| + h^{k+1} \right]. \end{aligned}$$

Here we have used the inverse properties (i) and (ii), together with the projection property (3.6a). This and Young’s inequality complete the proof of this lemma. \square

4.4. The numerical viscosity part

In this subsection we consider the numerical viscosity part, which is an important feature of the DG method.

4.4.1. A few propositions

The generalized numerical viscosity matrix $\mathcal{A}(\hat{\mathbf{f}}; \mathbf{u}_h)$ is not guaranteed to be positive semidefinite, causing some inconvenience to the analysis. However, it follows from assumptions A1 and A2 that there exists a positive semidefinite matrix $\mathcal{M}(\hat{\mathbf{f}}; \mathbf{p})$ such that

$$\mathcal{S}(\hat{\mathbf{f}}; \mathbf{u}_h) = \mathcal{A}(\hat{\mathbf{f}}; \mathbf{u}_h) + \mathcal{M}(\hat{\mathbf{f}}; \mathbf{u}_h) \tag{4.23}$$

is a symmetric positive semidefinite matrix, and

$$\|\mathcal{M}(\hat{\mathbf{f}}; \mathbf{u}_h)\|_M \leq C_\star \|\llbracket \mathbf{u}_h \rrbracket\|_M. \tag{4.24}$$

The total numerical viscosity in the DG spatial discretization can be shown explicitly in the form

$$[[\mathbf{u}_h]^\top \mathcal{S}[\mathbf{u}_h]] = \sum_{1 \leq j \leq N} \left[[[\mathbf{u}_h]^\top \mathcal{S}(\hat{\mathbf{f}}; \mathbf{u}_h)[\mathbf{u}_h]] \right]_{j+\frac{1}{2}}. \tag{4.25}$$

Note that this treatment is not essential, but only for the analysis. For example, the second Cauchy–Schwarz inequality in (2.8) holds for $\mathcal{S}(\hat{\mathbf{f}}; \mathbf{u}_h)$, but not for $\mathcal{A}(\hat{\mathbf{f}}; \mathbf{u}_h)$.

Associated with this issue, we will encounter another expression $[[\mathbf{u}_h]^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle [\mathbf{u}_h]]$. Below we would like to present some useful propositions between them.

Proposition 4.6. *There exists a bounding constant C_\star such that*

$$\frac{1}{2} [[\mathbf{u}_h]^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle [\mathbf{u}_h]] \leq [[\mathbf{u}_h]^\top \mathcal{S}[\mathbf{u}_h]] + C(e)(\|\boldsymbol{\xi}\|^2 + h^{2k+2}).$$

Note that the notation C_\star is included in $C(e)$; see (4.1).

Proof. Noticing (4.23), a simple manipulation yields

$$\frac{1}{2} [[\mathbf{u}_h]^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle [\mathbf{u}_h]] + [[\mathbf{u}_h]^\top \mathcal{U}[\mathbf{u}_h]] = [[\mathbf{u}_h]^\top \mathcal{S}[\mathbf{u}_h]] + \Lambda_3 + \Lambda_4, \tag{4.26}$$

where \mathcal{U} is the symmetric positive semidefinite matrix given in assumption A3, and

$$\begin{aligned} \Lambda_3 &= \frac{1}{2} [[\mathbf{u}_h]^\top \left(\langle \mathbf{H}(\mathbf{u}_b) \rangle - \langle \mathbf{H}(\{\{\mathbf{u}_h\}\}) \rangle - \mathcal{M}(\hat{\mathbf{f}}; \mathbf{u}_h) \right) [\mathbf{u}_h]], \\ \Lambda_4 &= [[\mathbf{u}_h]^\top \left(\frac{1}{2} \langle \mathbf{H}(\{\{\mathbf{u}_h\}\}) \rangle + \mathcal{U} - \mathcal{A}(\hat{\mathbf{f}}, \mathbf{u}_h) \right) [\mathbf{u}_h]]. \end{aligned}$$

We can bound easily $|\Lambda_3|$ by using (4.3), (4.5) and (4.24). The bound of $|\Lambda_4|$ is given by assumption A3, (4.6) and the inverse property (ii). Together with (4.6) again, we have completed the proof of this proposition. \square

Proposition 4.7. *Let ε be any given positive constant. There holds*

$$\frac{1}{2} [[\boldsymbol{\xi}]^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle [\boldsymbol{\xi}]] \leq (1 + \varepsilon) [[\mathbf{u}_h]^\top \mathcal{S}[\mathbf{u}_h]] + C(\varepsilon; e)(\|\boldsymbol{\xi}\|^2 + h^{2k+1}). \tag{4.27}$$

Proof. Since $\langle \mathbf{H} \rangle$ is a symmetric positive semidefinite matrix, we can use Cauchy–Schwarz inequality and Young’s inequality to get

$$\begin{aligned} [[\boldsymbol{\xi}]^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle [\boldsymbol{\xi}]] &= [[\boldsymbol{\eta} - \mathbf{u}_h]^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle [\boldsymbol{\eta} - \mathbf{u}_h]] \\ &\leq (1 + \varepsilon) [[\mathbf{u}_h]^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle [\mathbf{u}_h]] + (1 + \varepsilon^{-1}) [[\boldsymbol{\eta}]^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle [\boldsymbol{\eta}]]. \end{aligned} \tag{4.28}$$

Then we can finish the proof of this proposition by using the projection property of the finite element space, together with Proposition 4.6. \square

4.4.2. Boundedness

Let $F_{\max} = \max \|\mathbf{Q}^{-1/2}(\mathbf{s})\mathcal{A}(\hat{\mathbf{f}}; \mathbf{s}, \mathbf{s})\mathbf{Q}^{-1/2}(\mathbf{s})\|_{\mathbb{M}}$ represent the maximum strength of the numerical viscosity. Here the maximum is taken over all vectors in \mathbb{R}^m . Due to the smoothness assumptions H1–H3, we know that F_{\max} is a finite number.

The next lemma will show a rough amplification in each time-marching, for any functions in the finite element space.

Lemma 4.8. *Let $\varepsilon \leq 1$ be any positive constant. For any $\mathbf{v}_h \in \mathbb{V}_h$, there holds*

$$\begin{aligned} |\mathcal{V}(\mathbf{Q}_c \mathbf{v}_h)| &\leq \varepsilon F_{\max} \mu h^{-1} \|\mathbf{Q}_c^{1/2} \mathbf{v}_h\|^2 + K \varepsilon^{-1} \llbracket \mathbf{u}_h \rrbracket^\top \mathcal{S} \llbracket \mathbf{u}_h \rrbracket \\ &\quad + C_* h^{-1} \|\mathbf{e}\|_\infty \left[\|\boldsymbol{\xi}\|^2 + \|\mathbf{v}_h\|^2 + h^{2k+2} \right]. \end{aligned} \quad (4.29)$$

Proof. Recall that $\mathbf{r}_1 = \mathbf{u}_h^+$ and $\mathbf{r}_2 = \mathbf{u}_h^-$, with the local rotation positions \mathbf{s}_1 and \mathbf{s}_2 . It follows from the definition (2.14) that $\mathcal{V}(\mathbf{Q}_c \mathbf{v}_h) = \mathcal{V}^{(1)}(\mathbf{v}_h) + \mathcal{V}^{(2)}(\mathbf{v}_h)$, where

$$\mathcal{V}^{(1)}(\mathbf{v}_h) = \llbracket \mathbf{v}_h \rrbracket^\top \mathcal{A}(\hat{\mathbf{f}}; \mathbf{u}_h) \llbracket \mathbf{u}_h \rrbracket, \quad (4.30a)$$

$$\mathcal{V}^{(2)}(\mathbf{v}_h) = -\frac{1}{2} \sum_{\kappa=1,2} \sum_{1 \leq j \leq N} \left[\llbracket \mathbf{M}_\kappa \mathbf{v}_h \rrbracket^\top (\mathbf{f}(\mathbf{r}_\kappa) - \hat{\mathbf{f}}(\mathbf{u}_h)) \right]_{j+\frac{1}{2}}. \quad (4.30b)$$

Here $\mathbf{M}_\kappa = \mathbf{Q}(\mathbf{u}_c) - \mathbf{Q}(\mathbf{s}_\kappa)$ describes the focus shifting of the local rotation matrix. Below we will estimate these terms one by one.

Noticing (4.23), an application of Cauchy–Schwarz inequality to the first term yields the estimate

$$\begin{aligned} |\mathcal{V}^{(1)}(\mathbf{v}_h)| &\leq \|\llbracket \mathbf{v}_h \rrbracket^\top \mathcal{S} \llbracket \mathbf{u}_h \rrbracket\| + \|\llbracket \mathbf{v}_h \rrbracket^\top \mathcal{M} \llbracket \mathbf{u}_h \rrbracket\| \\ &\leq K \varepsilon^{-1} \llbracket \mathbf{u}_h \rrbracket^\top \mathcal{S} \llbracket \mathbf{u}_h \rrbracket + \varepsilon \llbracket \mathbf{v}_h \rrbracket^\top \mathcal{S} \llbracket \mathbf{v}_h \rrbracket + \|\llbracket \mathbf{v}_h \rrbracket^\top \mathcal{M} \llbracket \mathbf{u}_h \rrbracket\|. \end{aligned} \quad (4.31)$$

The second term is equal to the sum of $\mathcal{V}^{(1a)} = \varepsilon \llbracket \mathbf{v}_h \rrbracket^\top \mathcal{A}(\hat{\mathbf{f}}; \{\mathbf{u}_h\}) \llbracket \mathbf{v}_h \rrbracket$ and

$$\mathcal{V}^{(1b)} = \varepsilon \llbracket \mathbf{v}_h \rrbracket^\top \left[\mathcal{A}(\hat{\mathbf{f}}; \mathbf{u}_h) - \mathcal{A}(\hat{\mathbf{f}}; \{\mathbf{u}_h\}) + \mathcal{M}(\hat{\mathbf{f}}; \mathbf{u}_h) \right] \llbracket \mathbf{v}_h \rrbracket, \quad (4.32)$$

which are bounded by the inverse property (4.7b), and by assumption A2 and (4.24), respectively. This gives

$$\varepsilon \llbracket \mathbf{v}_h \rrbracket^\top \mathcal{S} \llbracket \mathbf{v}_h \rrbracket \leq \varepsilon K F_{\max} \mu h^{-1} \|\mathbf{Q}_c^{1/2} \mathbf{v}_h\|^2 + C_* h^{-1} \|\mathbf{e}\|_\infty \|\mathbf{v}_h\|^2, \quad (4.33)$$

since $\varepsilon \leq 1$. The last term can be bounded by using (4.24) and (4.6), which reads

$$\|\llbracket \mathbf{v}_h \rrbracket^\top \mathcal{M} \llbracket \mathbf{u}_h \rrbracket\| \leq C_* h^{-1} \|\mathbf{e}\|_\infty \left[\|\mathbf{v}_h\|^2 + \|\boldsymbol{\xi}\|^2 + h^{2k+2} \right]. \quad (4.34)$$

As we have shown in Section 4.1.2, we know that both $\|\mathbf{M}_\kappa\|_M \leq C_*(h + \|\mathbf{e}\|_\infty)$ and $\|\mathbf{f}(\mathbf{u}_h^+) - \hat{\mathbf{f}}(\mathbf{u}_h^-, \mathbf{u}_h^+)\|_M \leq C \|\llbracket \mathbf{u}_h \rrbracket\|_M$ hold at every element boundary point. The inverse property (ii) and the projection property (3.6a) yield

$$|\mathcal{V}^{(2)}(\mathbf{v}_h)| \leq C_* h^{-1} \|\mathbf{e}\|_\infty (\|\boldsymbol{\xi}\| + h^{k+1}) \|\mathbf{v}_h\|. \quad (4.35)$$

Finally, summing up the above estimates completes the proof of this lemma. For convenience we use a new notation ε instead of εK in the final conclusion. \square

Furthermore, for a special test function $\mathbf{v}_h = \boldsymbol{\xi}$, we can show explicitly the total numerical viscosity in the DG discretization. It is stated in the following lemma.

Lemma 4.9. *For any given positive constant $\varepsilon < 1$, there holds*

$$\mathcal{V}(\mathbf{Q}_c \boldsymbol{\xi}) \leq -(1 - \varepsilon) \llbracket \mathbf{u}_h \rrbracket^\top \mathcal{S} \llbracket \mathbf{u}_h \rrbracket + C(\varepsilon; \mathbf{e}) \left[\|\boldsymbol{\xi}\|^2 + h^{2k+1} \right]. \quad (4.36)$$

Proof. At this time, we still use the splitting (4.30), and need to pay more attention to the first term. It reads

$$\mathcal{V}^{(1)}(\boldsymbol{\xi}) = \llbracket \boldsymbol{\xi} \rrbracket^\top \mathcal{A} \llbracket \mathbf{u}_h \rrbracket = \llbracket \boldsymbol{\eta} \rrbracket^\top \mathcal{S} \llbracket \mathbf{u}_h \rrbracket - \llbracket \boldsymbol{\xi} \rrbracket^\top \mathcal{M} \llbracket \mathbf{u}_h \rrbracket - \llbracket \mathbf{u}_h \rrbracket^\top \mathcal{S} \llbracket \mathbf{u}_h \rrbracket, \quad (4.37)$$

following (4.23) and (4.6). The first two terms on the right-hand side are denoted by Λ_1 and Λ_2 , respectively.

The definition (2.14) implies that every component in \mathcal{A} are bounded uniformly, owing to assumptions H2 and H3. Hence we have $\|\mathcal{S}\|_{\mathcal{M}} \leq C + C_* \|e\|_{\infty}$, due to (4.24) and (4.6). Employing Young’s inequality, we get

$$|A_1| \leq \varepsilon [\mathbf{u}_h]^\top \mathcal{S} [\mathbf{u}_h] + \frac{1}{4\varepsilon} [\boldsymbol{\eta}]^\top \mathcal{S} [\boldsymbol{\eta}] \leq \varepsilon [\mathbf{u}_h]^\top \mathcal{S} [\mathbf{u}_h] + \varepsilon^{-1} [C + C_* \|e\|_{\infty}] \|[\boldsymbol{\eta}]\|_{T_h}^2.$$

Using (4.24) and (4.6), as well as the inverse property (ii), we have

$$|A_2| \leq \|\mathcal{M}\|_{\infty} \|[\mathbf{u}_h]\|_{T_h} \|[\boldsymbol{\eta}]\|_{T_h} \leq C_* h^{-1} \|e\|_{\infty} \left[\|\boldsymbol{\xi}\|^2 + h \|[\boldsymbol{\eta}]\|_{T_h}^2 \right].$$

Substituting the above inequalities into (4.37) yields

$$\mathcal{V}^{(1)}(\boldsymbol{\xi}) \leq -(1 - \varepsilon) [\mathbf{u}_h]^\top \mathcal{S} [\mathbf{u}_h] + C(\varepsilon; e) \left[\|\boldsymbol{\xi}\|^2 + h^{2k+2} \right], \tag{4.38}$$

where the projection property (3.6a) has been used.

The next term $\mathcal{V}^{(2)}(\boldsymbol{\xi})$ can be bounded by (4.35). This completes the proof of this lemma. \square

Remark 4.10. Actually the numerical viscosity of the DG method can be shown by another process. For example, if we use the reference vector along the upwind direction at each element boundary point, the numerical viscosity can be shown explicitly by the quantity $[\mathbf{u}_h]^\top \langle \mathbf{H}(\mathbf{u}_b) \rangle [\mathbf{u}_h]$. When the local Gauss–Radau projection is used, the numerical viscosity is shown to come mainly from the linear part. This complex line of analysis is only necessary to obtain the optimal error estimate, when the upwind numerical flux is used. In this case, the new viscosity part will basically disappear. See section 5.6.

4.5. The time-marching part

By using the approximation property (3.6c), and the estimate of the local truncation error in time, we can easily get the following conclusion.

Lemma 4.11. *There exists a constant $C > 0$ independent of h and τ , such that*

$$|\mathcal{T}(\mathbf{Q}_c \mathbf{v}_h)| \leq C \left[\|\mathbf{v}_h\|^2 + h^{2k+2} + \tau^6 \right], \quad \forall \mathbf{v}_h \in \mathbb{V}_h. \tag{4.39}$$

Proof. The proof is a simple application of Young’s inequality, so it is omitted. \square

5. PROOF OF THE THEOREM

In this section we come back to prove Theorem 2.5, under a standard CFL condition

$$\max\{S_{\max}, F_{\max}\} \mu \tau h^{-1} \leq \lambda_{\max}, \tag{5.1}$$

where S_{\max} and F_{\max} have been defined in the previous section. Note that $\lambda_{\max} \leq 1$ is a suitable CFL number, independent of h and τ .

5.1. Three useful differences

The following energy analysis depends strongly on three important differences, which have been presented in [31] for the scalar case. They read

$$\mathbb{G}_1^n \boldsymbol{\xi}^n = \boldsymbol{\xi}^{n,1} - \boldsymbol{\xi}^n, \quad \mathbb{G}_2^n \boldsymbol{\xi}^n = 2\boldsymbol{\xi}^{n,2} - \boldsymbol{\xi}^{n,1} - \boldsymbol{\xi}^n, \quad \mathbb{G}_3^n \boldsymbol{\xi}^n = \boldsymbol{\xi}^{n+1} - 2\boldsymbol{\xi}^{n,2} + \boldsymbol{\xi}^n. \tag{5.2}$$

It is worthy to point out that the group $\{\mathbb{E}_\ell^n \boldsymbol{\xi}^n\}_{\ell=1,2,3}$ and the group $\{\mathbb{G}_\ell^n \boldsymbol{\xi}^n\}_{\ell=1,2,3}$ can be linearly expressed by each other.

It is convenient to understand \mathbb{G}_ℓ^n as an operator acting on different objects defined on time stage levels, for example, on matrix-valued functions, or on linear functionals. If the object is the error in the finite element space, we would omit the object and use a simplified notation $\mathbb{G}_\ell^n = \mathbb{G}_\ell^n \boldsymbol{\xi}^n$, in this section and in the appendix.

The following conclusion shows that the above differences \mathbb{G}_ℓ^n have a strong relationship to approximations of time derivatives.

Lemma 5.1. *For the RKDG3 method (2.4) and $\ell = 0, 1, 2$, we have*

$$(\mathbb{G}_{\ell+1}^n, \mathbf{v}_h) = \frac{\tau}{\ell+1} \mathbb{G}_\ell^n \mathcal{K}^n(\mathbf{v}_h), \quad \mathbf{v}_h \in \mathbb{V}_h. \tag{5.3}$$

Here \mathbb{G}_0^n is just the identity operator, namely $\mathbb{G}_0^n \mathcal{K}^n = \mathcal{K}^n$.

The proof of this lemma is trivial by the suitable linear combinations of those error equations in (3.7). The detailed process is omitted; we refer to [31]. Based on this lemma, we can achieve easily a crude estimate for each stage in L^2 -norm within one time step. The proof of the following lemma is given in the appendix.

Lemma 5.2. *Let n be any time level and $\ell = 0, 1, 2$. Assume $\|\mathbf{e}^{n,\kappa}\|_\infty \leq h$ holds for any $\kappa \leq \ell$. Then there holds*

$$\|\boldsymbol{\xi}^{n,\ell+1}\|^2 \leq C \sum_{0 \leq \kappa \leq \ell} \|\boldsymbol{\xi}^{n,\kappa}\|^2 + Ch^{2k+2} + C\tau^8, \tag{5.4}$$

under the CFL condition (5.1), where the bounding constant C is independent of n, h and τ . Here $\boldsymbol{\xi}^{n,3} = \boldsymbol{\xi}^{n+1}$ for notational convenience.

5.2. The energy equation

To carry out the energy analysis, we need to establish the energy equation. To this end, we take three test functions $\mathbf{Q}_c^n \boldsymbol{\xi}^n$, $\mathbf{Q}_c^{n,1} \boldsymbol{\xi}^{n,1}$ and $\mathbf{Q}_c^{n,2} \boldsymbol{\xi}^{n,2}$ in (3.7), for $\ell = 0, 1, 2$, respectively. Recall that the piecewise constant rotation matrix is given as

$$\mathbf{Q}_c^{n,\ell} = \mathbf{Q}(\mathbf{u}_c^{n,\ell}), \quad \ell = 0, 1, 2. \tag{5.5}$$

A simple manipulation yields the energy equation in the form

$$3\|(\mathbf{Q}_c^{n+1})^{1/2} \boldsymbol{\xi}^{n+1}\|^2 - 3\|(\mathbf{Q}_c^n)^{1/2} \boldsymbol{\xi}^n\|^2 = \Pi_1 + \Pi_2 + \Pi_3, \tag{5.6}$$

which are used to reflect the DG discretization, the TVDRK time-marching, and the variation of the rotation matrix in time, respectively. In details, they read

$$\Pi_1 = \tau[\mathcal{K}^n(\mathbf{Q}_c^n \boldsymbol{\xi}^n) + \mathcal{K}^{n,1}(\mathbf{Q}_c^{n,1} \boldsymbol{\xi}^{n,1}) + 4\mathcal{K}^{n,2}(\mathbf{Q}_c^{n,2} \boldsymbol{\xi}^{n,2})], \tag{5.7a}$$

$$\Pi_2 = (\mathbb{G}_2^n, \mathbf{Q}_c^n \mathbb{G}_2^n) + 3(\mathbb{G}_1^n, \mathbf{Q}_c^n \mathbb{G}_3^n) + 3(\mathbb{G}_2^n, \mathbf{Q}_c^n \mathbb{G}_3^n) + 3(\mathbb{G}_3^n, \mathbf{Q}_c^n \mathbb{G}_3^n), \tag{5.7b}$$

$$\Pi_3 = 3(\boldsymbol{\xi}^{n+1}, \tilde{\mathbf{Q}}_c^{n,3} \boldsymbol{\xi}^{n+1}) - (\mathbb{E}^{n,2} \boldsymbol{\xi}^n, \tilde{\mathbf{Q}}_c^{n,1} \boldsymbol{\xi}^{n,1}) - 4(\mathbb{E}^{n,3} \boldsymbol{\xi}^n, \tilde{\mathbf{Q}}_c^{n,2} \boldsymbol{\xi}^{n,2}). \tag{5.7c}$$

Here $\tilde{\mathbf{Q}}_c^{n,\ell} = \mathbf{Q}_c^{n,\ell} - \mathbf{Q}_c^n$ for $\ell = 1, 2, 3$, reflect the focus shifting, with the notations $\mathbf{Q}_c^{n,3} = \mathbf{Q}_c^{n+1}$ and $\mathbf{u}_c^{n,3} = \mathbf{u}_c^{n+1}$ for convenience.

Below we will estimate the above three terms one by one. For notational convenience, we would like to use a macro notation

$$\Psi_{\kappa_1, \kappa_2}^n = \sum_{0 \leq \ell \leq \kappa_1} C(e^{n,\ell}) \left[\sum_{0 \leq \ell \leq \kappa_2} \|\boldsymbol{\xi}^{n,\ell}\|^2 + h^{2k+1} + \tau^6 \right], \tag{5.8}$$

where κ_1 and κ_2 are either 2 or 3, and the notation $C(e^{n,\ell})$ has been given in (4.1).

5.3. Estimates to the first term and the third term

The estimate to Π_1 is straightforward, since the test function and the error functional are taken at the same time stage level. Namely, we use Lemmas 4.2 and 4.3 for the linear part, Lemma 4.5 for the nonlinear part, Lemma 4.9 for the viscosity part, and Lemma 4.11 for the time-marching part. All the test functions are taken as $\mathbf{v}_h = \boldsymbol{\xi}^{n,\ell}$. Then we use Proposition 4.7 and obtain

$$\tau \mathcal{K}^{n,\ell}(\mathbf{Q}_c^{n,\ell} \boldsymbol{\xi}^{n,\ell}) \leq -(1 - \varepsilon) [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}] \tau + C(\varepsilon; \mathbf{e}^{n,\ell}) \left[\|\boldsymbol{\xi}^{n,\ell}\|^2 + h^{2k+1} + \tau^6 \right] \tau,$$

for $\ell = 0, 1, 2$. By taking ε small enough, we achieve the estimate

$$\Pi_1 \leq -\frac{3}{4} \sum_{\ell=0,1,2} [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}] \tau + \Psi_{2,2}^n \tau. \tag{5.9}$$

Note that the last term gives a very crude bound.

The estimate to Π_2 is more complex, so we present the detailed discussion in the next subsection. As for the third term, a simple application of Cauchy–Schwarz inequality yields

$$\Pi_3 \leq C \left[\|\boldsymbol{\xi}^n\|^2 + \|\boldsymbol{\xi}^{n,1}\|^2 + \|\boldsymbol{\xi}^{n,2}\|^2 + \|\boldsymbol{\xi}^{n+1}\|^2 \right] \tau. \tag{5.10}$$

since $\|\tilde{\mathbf{Q}}_c\|_M = \|\mathbf{Q}_c^{n,\ell} - \mathbf{Q}_c^n\|_M = \mathcal{O}(\tau)$ holds for all x , as we have discussed before.

5.4. Estimate to the second term

The four terms involved in Π_2 are denoted by $\Pi_2^{(\kappa)}$ for $\kappa = 1, 2, 3, 4$. The estimates to them are a little long and technical, depending strongly on the three differences (5.2) and the relationships among them.

One main trick used here is that we consider the first two terms $\Pi_2^{(1)}$ and $\Pi_2^{(2)}$ together. By taking the test functions $\mathbf{Q}_c^n \mathbb{G}_2^n$ and $\mathbf{Q}_c^n \mathbb{G}_1^n$, respectively, in Lemma 5.1 for $\ell = 1, 2$, we have

$$\begin{aligned} \Pi_2^{(1)} + \Pi_2^{(2)} &= -(\mathbb{G}_2^n, \mathbf{Q}_c^n \mathbb{G}_2^n) + 2(\mathbb{G}_2^n, \mathbf{Q}_c^n \mathbb{G}_2^n) + 3(\mathbb{G}_3^n, \mathbf{Q}_c^n \mathbb{G}_1^n) \\ &= -\|\mathbb{G}_2^n\|_n^2 + \tau \mathbb{G}_1^n \mathcal{K}^n(\mathbf{Q}_c^n \mathbb{G}_2^n) + \tau \mathbb{G}_2^n \mathcal{K}^n(\mathbf{Q}_c^n \mathbb{G}_1^n) \\ &= -\|\mathbb{G}_2^n\|_n^2 + \Upsilon_1 + \Upsilon_2 + \Upsilon_3 + \Upsilon_4. \end{aligned} \tag{5.11}$$

Note that the new norm $\|\mathbf{w}_h\|_n = \|(\mathbf{Q}_c^n)^{1/2} \mathbf{w}_h\|$ is used for convenience. Here the last four terms represent respectively the linear part, the nonlinear part, the numerical viscosity part and the time-marching part. For example,

$$\Upsilon_1 = \tau \mathbb{G}_1^n \mathcal{L}^n(\mathbf{Q}_c^n \mathbb{G}_2^n) + \tau \mathbb{G}_2^n \mathcal{L}^n(\mathbf{Q}_c^n \mathbb{G}_1^n). \tag{5.12}$$

We will estimate them separately below.

The term Υ_1 can be bounded along the same line as that for the linear part at the given time stage. However, we need an additional treatment to deal with the variation of the Jacobian matrix $\mathbf{f}'_{\mathbf{u}}(\mathbf{u}^{n,\ell})$ and the local setting of the local rotation matrix $\mathbf{Q}_c^{n,\ell}$. This again involves many applications of local shifting of the focus. A key point in the technique is to fully use the approximate skew-symmetric property (see Lem. 4.2). The complete proof is given in the appendix, and we just present the conclusion here.

Lemma 5.3. *There exists a bounding constant C , such that*

$$|\Upsilon_1| \leq \frac{1}{8} \sum_{\ell=0,1} [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}] \tau + C \Psi_{2,2}^n \tau. \tag{5.13}$$

Each term in \mathcal{Y}_2 is bounded easily by using Lemma 4.5. Together with the triangle inequalities to the norms $\|\mathbb{G}_1^n\|$ and $\|\mathbb{G}_2^n\|$, we have

$$|\mathcal{Y}_2| \leq C\Psi_{2,2}^n\tau. \quad (5.14)$$

There are some minor troubles to estimate \mathcal{Y}_3 , since the generalized viscosity matrix $\mathcal{A}^{n,\ell}$ involves different time stages. Thanks to assumption A2, we have the following result to dig out the complete information about $\|\mathbb{G}_2^n\|_n$. The proof is again postponed to the appendix.

Lemma 5.4. *There exist bounding constants K and C , such that*

$$|\mathcal{Y}_3| \leq K\lambda_{\max}\|\mathbb{G}_2^n\|_n^2 + \frac{1}{8} \sum_{\ell=0,1} [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}]_\tau + C\Psi_{2,2}^n\tau. \quad (5.15)$$

The last term \mathcal{Y}_4 is easily estimated from Lemma 4.9. Together with the projection property (3.6c) and the triangle inequality, we have

$$|\mathcal{Y}_4| \leq C\Psi_{2,2}^n\tau. \quad (5.16)$$

Next we turn to estimate separately the remaining terms, $\Pi_2^{(3)}$ and $\Pi_2^{(4)}$. It follows from Lemma 5.1, with different test functions $\mathbf{Q}_c^n\mathbb{G}_2^n$ and $\mathbf{Q}_c^n\mathbb{G}_3^n$, respectively, for $\ell = 2$, that

$$\Pi_2^{(3)} = 3(\mathbb{G}_3^n, \mathbf{Q}_c^n\mathbb{G}_2^n) = \tau\mathbb{G}_2^n\mathcal{K}^n(\mathbf{Q}_c^n\mathbb{G}_2^n), \quad \Pi_2^{(4)} = 3\|\mathbb{G}_3^n\|_n^2 = \tau\mathbb{G}_2^n\mathcal{K}^n(\mathbf{Q}_c^n\mathbb{G}_3^n). \quad (5.17)$$

Both estimates depend on the boundedness of $\mathbb{G}_2^n\mathcal{K}^n(\mathbf{Q}_c^n\mathbf{v}_h)$ for any $\mathbf{v}_h \in \mathbb{V}_h$. This property can be achieved along almost the same line as we have done in section 4. However, in the discussion on the linear part, we will also encounter the trouble on how to dig out the complete information about $\|\mathbb{G}_2^n\|_n$. The final conclusion is stated in the following lemma.

Lemma 5.5. *There holds for any $\mathbf{v}_h \in \mathbb{V}_h$, that*

$$\begin{aligned} |\tau\mathbb{G}_2^n\mathcal{K}^n(\mathbf{Q}_c^n\mathbf{v}_h)| &\leq \varepsilon\|\mathbf{v}_h\|_n^2 + \sum_{\ell=0,1,2} C_*h^{-1}\|\mathbf{e}^{n,\ell}\|_\infty\|\mathbf{v}_h\|^2\tau + C\Psi_{2,2}^n\tau \\ &\quad + K\varepsilon^{-1}\lambda_{\max} \left[\|\mathbb{G}_2^n\|_n^2 + \sum_{\ell=0,1,2} [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}]_\tau \right], \end{aligned} \quad (5.18)$$

where ε is any small positive constant.

We would like to postpone the proof to the appendix. Noticing (5.17), we take $\mathbf{v}_h = \mathbb{G}_2^n$ and $\varepsilon = 1/4$ in Lemma 5.5. Therefore,

$$|\Pi_2^{(3)}| \leq C\Psi_{2,2}^n\tau + \left[\frac{1}{4} + K\lambda_{\max}\right]\|\mathbb{G}_2^n\|_n^2 + K\lambda_{\max} \sum_{\ell=0,1,2} [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}]_\tau. \quad (5.19)$$

Here the triangle inequality to $\|\mathbb{G}_2^n\|$ is used. The estimate to $\Pi_2^{(4)}$ can be achieved by using the same lemma with $\mathbf{v}_h = \mathbb{G}_3^n$, and small enough parameter ε . Thus we can get from this process that

$$|\Pi_2^{(4)}| \leq C\Psi_{2,3}^n\tau + K\lambda_{\max} \left[\|\mathbb{G}_2^n\|_n^2 + \sum_{\ell=0,1,2} [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}]_\tau \right]. \quad (5.20)$$

Here we use the bigger bound $\Psi_{2,3}^n$, since \mathbb{G}_3^n involves the information about $\boldsymbol{\xi}^{n+1}$.

Finally, we collect all of the above estimates, namely, identity (5.11), inequalities (5.13)–(5.16), (5.19) and (5.20). This gives an estimate

$$\Pi_2 \leq - \left[\frac{3}{4} - K_1 \lambda_{\max} \right] \|\mathbb{G}_2^n\|_n^2 + \left[\frac{1}{4} + K_2 \lambda_{\max} \right] \sum_{\ell=0,1,2} [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}]_\tau + C \Psi_{2,3}^n \tau, \quad (5.21)$$

where K_1 and K_2 solely depend on the inverse constant μ . We have now estimated every term on the right-hand side of the energy equation (5.6).

5.5. The energy estimate

In this subsection we would like to complete the proof of Theorem 2.5. To deal with the nonlinearity of the flux function $\mathbf{f}(\mathbf{u})$, we would like first make the *a priori* assumption that, for all n , if $n\tau < T$, we have

$$\|\mathbf{e}^{n,\ell}\|_\infty \leq h, \quad \ell = 0, 1, 2, \quad (5.22)$$

which holds for h small enough. We note that, for the linear flux $\mathbf{f}(\mathbf{u}) = \mathbf{A}\mathbf{u}$ with the constant matrix \mathbf{A} , this *a priori* assumption is not necessary. Later, we will verify the correctness of this *a priori* assumption for the piecewise polynomials of degree $k \geq 2$.

The assumption (5.22) implies that $C(\mathbf{e}^{n,\ell}) \leq C$ for any n satisfying $n\tau < T$ and $\ell = 0, 1, 2$, where the bounding constant C is independent of n, h and τ . Substitute (5.9), (5.21) and (5.10) into the energy identity (5.6). By employing Lemma 5.2, we finally obtain for any n satisfying $n\tau < T$ that

$$3\|(\mathbf{Q}_c^{n+1})^{1/2} \boldsymbol{\xi}^{n+1}\|^2 - 3\|(\mathbf{Q}_c^n)^{1/2} \boldsymbol{\xi}^n\|^2 + \Theta_{\text{RKDG}}^n \leq C\|\boldsymbol{\xi}^n\|^2 \tau + Ch^{2k+1} \tau + C\tau^7, \quad (5.23)$$

where the bounding constant C is independent of n, h and τ , and

$$\Theta_{\text{RKDG}}^n = \left(\frac{1}{2} - K_1 \lambda_{\max} \right) \|\mathbb{G}_2^n\|_n^2 + \left(\frac{1}{2} - K_2 \lambda_{\max} \right) \sum_{\ell=0,1,2} [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}]_\tau, \quad (5.24)$$

with the constants K_1 and K_2 defined in (5.21). Suppose that the CFL condition (5.1) holds under the restriction

$$\lambda_{\max} \leq \min((2K_1)^{-1}, (2K_2)^{-1}), \quad (5.25)$$

there holds $\Theta_{\text{RKDG}}^n \geq 0$ to reflect the numerical stability of the RKDG3 method. Summing up the estimate (5.23) from 0 to n yields

$$\|\boldsymbol{\xi}^{n+1}\|^2 \leq C \sum_{n'=0}^n \|\boldsymbol{\xi}^{n'}\|^2 \tau + C \left[\|\boldsymbol{\xi}^0\|^2 + h^{2k+1} + \tau^6 \right], \quad \forall n : n\tau < T, \quad (5.26)$$

where the bounding constant $C > 0$ is independent of n, h and τ . Here we have used the uniform equivalence among the norms $\|\cdot\|$ and $\|(\mathbf{Q}_c^n)^{1/2} \cdot\| \equiv \|\cdot\|_n$ for any n , due to the uniform boundedness of \mathbf{Q}_c^n ; see assumption H3.

It follows from the setting of the initial solution that $\boldsymbol{\xi}^0 = 0$. An application of the discrete Gronwall lemma yields the error estimate for the fully-discrete DG scheme with TVDRK3 time-marching, in the form

$$\|\boldsymbol{\xi}^{n+1}\|^2 \leq Ch^{2k+1} + C\tau^6, \quad \forall n : n\tau < T, \quad (5.27)$$

where the bounding constant C is independent of n, h and τ . This inequality together with the projection property (3.6a) yields the conclusion of this theorem.

To complete the proof of this theorem for the generalized E-flux, we need to verify the *a priori* assumption (5.22), by using an induction process, as well as the inverse property (iii) and the projection property (3.6b).

Obviously there holds $\|\mathbf{e}^0\|_\infty \leq h$, owing to the setting of the initial solution. An application of Lemma 5.2 yields $\|\boldsymbol{\xi}^{0,1}\| \leq Ch^{k+1}$. Thus it follows that

$$\|\mathbf{e}^{0,1}\|_\infty \leq \mu h^{-\frac{1}{2}} \|\boldsymbol{\xi}^{0,1}\| + \|\boldsymbol{\eta}^{0,1}\|_\infty \leq Ch^{k+\frac{1}{2}} \leq h, \quad (5.28)$$

if h is small enough, since $k \geq 2$. Repeating this discussion shows $\|\mathbf{e}^{0,2}\|_\infty \leq h$. Next, we will prove the *a priori* assumption (5.22) by induction. Assume it holds on all time level from 0 to n . We can achieve the inequality (5.27) along the above analysis line, which follows

$$\|\mathbf{e}^{n+1}\|_\infty \leq C\mu h^{-1/2} \|\boldsymbol{\xi}^{n+1}\| + \|\boldsymbol{\eta}^{n+1}\|_\infty \leq Ch^k + h^{-1/2}\tau^3 \leq h, \quad (5.29)$$

if h is small enough. As we have done for the initial time level, we can employ Lemma 5.2 to get $\|\mathbf{e}^{n+1,1}\|_\infty \leq h$ and $\|\mathbf{e}^{n+1,2}\|_\infty \leq h$. Thus the assumption (5.22) is also true for $n+1$, and we complete the proof of the theorem.

Remark 5.6. The main technique used in this paper is the energy analysis. The energy technique does not require a uniform mesh and can be easily generalized to arbitrary triangulation in multi-dimensions and for linear equations with variable coefficients, as well as to some types of non-periodic boundary conditions. We have discussed in [28] the inflow boundary condition for the scalar hyperbolic equation.

Remark 5.7. For linear conservation laws, the error estimates hold for any piecewise polynomials, if the exact solution is smooth enough. However, for nonlinear conservation laws, the error estimates only hold for piecewise polynomials with degree $k \geq d/2 + 1$, where d is the spatial dimension. This restriction is solely used to ensure the *a priori* assumption (5.22). As a result, in this paper we do not consider the piecewise linear polynomials, which are rarely used together with the third order TVDRK time-marching.

5.6. Optimal error estimate for upwind fluxes

In this subsection, we would like to upgrade the error estimate to be optimal for the upwind numerical fluxes, as we have stated in Theorem 2.5. Since almost the same analysis line is used here, we would like in this section to only point out the new techniques and main differences to the above.

To obtain the optimal error estimate, we would need to use some standard tricks in the DG analysis, as we have done in [30, 31]. These consist of two main ingredients, both of which are carried out by the help of eigenvector decomposition. Let $\lambda_i(\mathbf{u})$ be the eigenvalues of the Jacobian matrix $\mathbf{f}'_{\mathbf{u}}(\mathbf{u})$, which has the left eigenvector $\mathbf{l}_i(\mathbf{u})$ and the right eigenvector $\mathbf{r}_i(\mathbf{u})$. It is assumed that $\lambda_i(\mathbf{u}), \mathbf{l}_i(\mathbf{u})$ and $\mathbf{r}_i(\mathbf{u})$ are smooth mappings which for strictly hyperbolic system follows from the regularity assumption of \mathbf{f} . Furthermore, we assume for any space point $x \in I$ that $\lambda_i(\mathbf{u}(x, t))$ does not change its sign for any $t \in [0, T]$.

The first one is the local Gauss–Radau projection, denoted by $\mathbb{Q}_h \mathbf{u}^{n,\ell}$, instead of the local L²-projection, where $\mathbf{u}^{n,\ell}$ is the reference stage solution as we have defined in section 3.1. In each element I_j , this process involves three steps.

1. Transform $\mathbf{u}^{n,\ell}$ to the characteristics fields. This is achieved by left multiplying this by the matrix whose rows are the left eigenvectors of $\mathbf{f}'_{\mathbf{u}}(\mathbf{u}_c^n)$.
2. Apply the scalar Gauss–Radau projection to each of the components of the transformed vector-valued function, which depends on the sign of the corresponding eigenvalue of $\mathbf{f}'_{\mathbf{u}}(\mathbf{u}_c^n)$. If the eigenvalue is not smaller than zero, we demand the exact collocation of the corresponding component on the right boundary point. Otherwise, we demand the exact collocation on the left boundary point.
3. The result is transformed back to the original space by left multiplying the matrix whose columns are the right eigenvectors $\mathbf{r}_i(\mathbf{u}_c^n)$.

For this projection, the projection properties given in Section 3.2 still hold. The advantage of this projection is that we can correct the order reduction in the element boundary.

The other one is the upwind setting of the reference vector for the numerical solution at each element boundary point, denoted by $(\mathbf{u}_h^{n,\ell})^{\text{upw}}$, which is the so-called upwind direction we have referred to before. This treatment

is similar as the above and use the eigenvector decomposition of the Jacobian matrix $\mathbf{f}'_{\mathbf{u}}(\mathbf{u}_b^n)$. At the element boundary point $x_{j+1/2}$, this process involves three steps.

1. Transform the two limits of $\mathbf{u}_h^{n,\ell}$ to the characteristics fields of $\mathbf{f}'_{\mathbf{u}}(\mathbf{u}_b^n)$.
2. Apply the scalar upwind setting to each of the components of the transformed vector, which depends on the sign of the corresponding eigenvalue of $\mathbf{f}'_{\mathbf{u}}(\mathbf{u}_b^n)$. If the eigenvalue is not smaller than zero, we take the left limit. Otherwise, we take the right limit.
3. The result is transformed back to the original space.

Then we change the definition of the linear part (3.9a), where $\{\mathbf{w}\}$ is replaced by \mathbf{w}^{upw} . It follows from the new linear part that the total numerical viscosity is shown explicitly in the term $[\mathbf{u}_h^{n,\ell}]^\top \langle \mathbf{H}(\mathbf{u}_b^{n,\ell}) \rangle [\mathbf{u}_h^{n,\ell}]$.

Also we change the definitions of the nonlinear part (3.9b) and the viscosity part (3.9c), where $\{\mathbf{f}(\mathbf{u}_h^{n,\ell})\}$ is replaced by $\mathbf{f}((\mathbf{u}_h^{n,\ell})^{\text{upw}})$. Note that we do not need to use the generalized numerical matrix. In this case, the new viscosity part is very small. The corresponding result is modified from the conclusion in Lemma 4.8, where only the last term on the right-hand side is left. It is worthy to point out that this new viscosity part is equal to zero for the special cases that each eigenvalue keeps the sign in the considered range for the exact solution and the numerical solution. If the sign changes locally, there must exist locally a zero eigenvalue. We refer to [31] for the detailed discussions for the scalar case. This is the key ingredient of the upwind numerical flux.

6. CONCLUDING REMARKS

In this paper we present quasi-optimal or optimal error estimates to the RKDG scheme with the explicit third order TVDRK time discretization, to solve a symmetrizable system of conservation laws with sufficiently smooth solution. Symmetrizable but not symmetric systems have several difficult issues. Although many numerical fluxes have been used successively in practical computation, the theoretical framework for error analysis has not been clear. In this paper we attempt to set an abstract framework to handle these difficulties, and to establish the so-called generalized numerical viscosity matrix. Another difficulty is related to the local variation of the numerical viscosity, which can be handled by the definition of the local rotation matrix. We remark that this local rotation matrix causes complexity in the analysis. The abstract framework established in this paper includes many commonly used numerical fluxes and is flexible to use for our analysis.

In this paper the exact solution of the conservation laws is assumed to be sufficiently smooth. However, generic solutions of nonlinear conservation laws contain discontinuities. In future work, we will consider piecewise smooth solutions and attempt to analyze the error in smooth regions away from discontinuities, extending the preliminary results for one-dimensional linear equations in [32].

APPENDIX

A.1. Verification for the Lax–Friedrichs flux

In this subsection we would like to show that the (global/local) Lax–Friedrichs flux is a generalized E-flux, and the corresponding generalized numerical viscosity matrix satisfies assumptions A2 and A3. The (global/local) Lax–Friedrichs flux is defined as follows:

$$\hat{\mathbf{f}}^{\text{LF}}(\mathbf{p}^-, \mathbf{p}^+) = \frac{1}{2} \left[\mathbf{f}(\mathbf{p}^-) + \mathbf{f}(\mathbf{p}^+) - \alpha [\mathbf{p}] \right], \quad \text{with } \alpha = \max_{\mathbf{p}} \varrho(\mathbf{f}'_{\mathbf{u}}(\mathbf{p})), \tag{A.1}$$

where the maximum is taken locally or globally.

Let \mathbf{s}_1 be any point in the super-rectangle with the vertexes \mathbf{p}^\pm , and define

$$\mathbf{Y}_1^{\text{LF}} = \frac{1}{2} \int_0^1 \left(\mathbf{Q}(\mathbf{s}) - \mathbf{Q}(\mathbf{s}_1) \right) (\mathbf{f}'_{\mathbf{u}}(\mathbf{s}) + \alpha \mathbf{I}) \, \text{d}\mathbf{s},$$

where \mathbf{I} is the identity matrix, and $\mathbf{s} = \mathbf{p}^- + s[\mathbf{p}]$. Every element in \mathbf{Y}_1^{LF} is bounded by $\mathcal{O}(\|\mathbf{p}\|_{\text{M}})$, since $\mathbf{v}'_{\mathbf{u}}$ is Lipschitz continuous and $\|\mathbf{s} - \mathbf{s}_1\|_{\text{M}} \leq \|\mathbf{p}\|_{\text{M}}$. A simple manipulation yields

$$[\mathbf{p}]^\top \mathbf{v}'_{\mathbf{u}}(s_1)(\mathbf{f}(\mathbf{p}^+) - \hat{\mathbf{f}}^{\text{LF}}(\mathbf{p}^-, \mathbf{p}^+)) + [\mathbf{p}]^\top \mathbf{Y}_1^{\text{LF}}[\mathbf{p}] = [\mathbf{p}]^\top \tilde{\mathbf{Y}}_1^{\text{LF}}[\mathbf{p}] \geq 0, \tag{A.2}$$

where

$$\tilde{\mathbf{Y}}_1^{\text{LF}} = \frac{1}{2} \int_0^1 \mathbf{Q}(s)(\mathbf{f}'_{\mathbf{u}}(s) + \alpha \mathbf{I}) ds,$$

is a symmetric positive semidefinite matrix, following (2.6) and (A.1). Similar discussion can be done for $\kappa = 2$. So the Lax–Friedrichs flux is a generalized E-flux.

Next we would like to verify assumption A2. We still consider $\kappa = 1$ as an example. In fact, it follows from the definition of the Lax–Friedrichs flux that

$$\left(\mathcal{D}\hat{\mathbf{f}}^{(1)}[\mathbf{p}^-, \mathbf{p}^+]\right)_{ij} = \frac{1}{2} \left[\frac{\partial f_i}{\partial p_j}(\tilde{\mathbf{p}}_{ij}) + \alpha \right], \quad i, j = 1, 2, \dots, m, \tag{A.3}$$

where $\tilde{\mathbf{p}}_{ij}$ is a mean point lying near $\{\mathbf{p}\}$. Since $\mathbf{s}_1 = \mathbf{s}_1(\mathbf{p}^-, \mathbf{p}^+)$ also lies near $\{\mathbf{p}\}$, it follows from the Lipschitz continuity of $\mathbf{f}'_{\mathbf{u}}$ and \mathbf{Q} , that $\mathcal{A}_1(\hat{\mathbf{f}}; \mathbf{p})$ is Lipschitz continuous. Similar discussion for $\mathcal{A}_2(\hat{\mathbf{f}}; \mathbf{p})$ implies assumption A2.

To verify assumption A3, we take $\mathcal{U} = \frac{1}{2} \mathbf{Q}^{1/2}[\alpha \mathbf{I} - \mathbf{K}(\{\mathbf{p}\})] \mathbf{Q}^{1/2}$, which is a symmetric positive semidefinite matrix. Noticing (2.17), a simple manipulation yields

$$\text{LHS of (2.16c)} = \left| -\frac{1}{2} \sum_{\kappa=1,2} [\mathbf{p}]^\top \left[\mathbf{Q}(\mathbf{s}_\kappa) - \mathbf{Q}(\{\mathbf{p}\}) \right] \left[\mathbf{f}(\mathbf{r}_\kappa) - \hat{\mathbf{f}}^{\text{LF}}(\mathbf{p}^-, \mathbf{p}^+) \right] \right|. \tag{A.4}$$

Noticing the Lipschitz continuity of \mathbf{Q} and $\hat{\mathbf{f}}^{\text{LF}}$, we can complete the verification of assumption A3.

A.2. Proof of Lemma 5.2

We can obtain this lemma by an induction process, to bound $\|\mathbb{G}_{\ell+1}^n\|$ by $\|\mathbb{G}_\ell^n\|$ and the others. This can be established by Lemma 5.1, with the test function $\mathbf{v}_h = \mathbf{Q}_c^n \mathbb{G}_{\ell+1}^n$.

The key here is to estimate the error functional $\mathbb{G}_\ell^n \mathcal{K}(\mathbf{Q}_c^n \mathbf{v}_h)$ for any test function in the finite element space. Similar analysis will be given in the proof of Lemma 5.5, so the detail is omitted here. The only difference is that after we bound the numerical viscosity part by Lemma 4.8, we amplify deeply those jumps by the inverse property (ii). Namely, for any $\mathbf{v}_h \in \mathbb{V}_h$,

$$\begin{aligned} |\mathcal{V}(\mathbf{u}_h; \mathbf{Q}_c^n \mathbf{v}_h)| &\leq K F_{\max} \mu h^{-1} \left[\varepsilon \|(\mathbf{Q}_c^n)^{1/2} \mathbf{v}_h\|^2 + \varepsilon^{-1} \|(\mathbf{Q}_c^n)^{1/2} \boldsymbol{\xi}\|^2 \right] \\ &\quad + C_* h^{-1} \|\mathbf{e}\|_\infty \left[\|\boldsymbol{\xi}\|^2 + \|\mathbf{v}_h\|^2 + h^{2k+2} \right]. \end{aligned} \tag{A.5}$$

Here ε is any given small positive constant.

Finally, we just need to use the CFL condition (5.1) and the assumption to simplify the obtained conclusion. This will yield the conclusion of this lemma. \square

A.3. Proof of Lemma 5.3

We can prove this lemma from the properties of the linear part and the equivalent expression $\mathcal{Y}_1 = \Phi_{\text{DG}} + \Phi_{\text{prj}} + \Phi_{\text{RK}}$, where

$$\begin{aligned} \Phi_{\text{DG}} &= \tau \mathcal{L}_{\text{main}}(\mathbf{f}'_{\mathbf{u}}(\mathbf{u}^n); \mathbb{G}_1^n, \mathbf{Q}_c^n \mathbb{G}_2^n) + \tau \mathcal{L}_{\text{main}}(\mathbf{f}'_{\mathbf{u}}(\mathbf{u}^n); \mathbb{G}_2^n, \mathbf{Q}_c^n \mathbb{G}_1^n), \\ \Phi_{\text{prj}} &= -\tau \mathcal{L}_{\text{main}}(\mathbf{f}'_{\mathbf{u}}(\mathbf{u}^n); \mathbb{G}_1^n \boldsymbol{\eta}^n, \mathbf{Q}_c^n \mathbb{G}_2^n) - \tau \mathcal{L}_{\text{main}}(\mathbf{f}'_{\mathbf{u}}(\mathbf{u}^n); \mathbb{G}_2^n \boldsymbol{\eta}^n, \mathbf{Q}_c^n \mathbb{G}_1^n), \\ \Phi_{\text{RK}} &= \tau \mathcal{L}_{\text{main}}(\mathbf{N}^{n,1}; \mathbf{e}^{n,1}, \mathbf{Q}_c^n \mathbb{G}_2^n - \mathbf{Q}_c^n \mathbb{G}_1^n) + 2\tau \mathcal{L}_{\text{main}}(\mathbf{N}^{n,2}; \mathbf{e}^{n,2}, \mathbf{Q}_c^n \mathbb{G}_1^n). \end{aligned}$$

Here $\mathbf{N}^{n,\ell} = \mathbf{f}'_{\mathbf{u}}(\mathbf{u}^{n,\ell}) - \mathbf{f}'_{\mathbf{u}}(\mathbf{u}^n)$ for $\ell = 1, 2$. Below we will estimate them one by one. It follows from Lemma 4.2 that

$$|\Phi_{\text{DG}}| \leq C_* \|\mathbb{G}_1^n\| \|\mathbb{G}_2^n\| \tau \leq C_* \left[\|\mathbb{G}_1^n\|^2 + \|\mathbb{G}_2^n\|^2 \right] \tau. \quad (\text{A.6})$$

This is the key step in the proof. Each term in Φ_{prj} can be bounded by Lemma 4.3 and Proposition 4.7. That is to say

$$\begin{aligned} |\Phi_{\text{prj}}| &\leq \varepsilon \sum_{\ell=1,2} [\mathbb{G}_\ell^n]^\top \langle \mathbf{H}(\mathbf{u}_b^n) \rangle [\mathbb{G}_\ell^n] \tau + \varepsilon \sum_{\ell=1,2} \|\mathbb{G}_\ell^n\|^2 \tau + Ch^{2k+1} \tau \\ &\leq K\varepsilon(1+\varepsilon) \sum_{\ell=0,1} [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}] \tau + \sum_{\ell=0,1} C(\varepsilon; \mathbf{e}^{n,\ell}) (\|\boldsymbol{\xi}^{n,\ell}\|^2 + h^{2k+1}) \tau, \end{aligned} \quad (\text{A.7})$$

where the triangle inequalities are used for $\|\mathbb{G}_1^n\|$ and $\|\mathbb{G}_2^n\|$. Since $\|\mathbf{N}^{n,\ell}\|_{\text{M}} = \mathcal{O}(\tau) = \mathcal{O}(h)$ holds for all x , along almost the same line as that in Lemmas 4.1 and 4.3, we have

$$|\Phi_{\text{RK}}| \leq C_* \left[\|\boldsymbol{\xi}^n\|^2 + \|\boldsymbol{\xi}^{n,1}\|^2 + \|\boldsymbol{\xi}^{n,2}\|^2 + \|\boldsymbol{\xi}^{n+1}\|^2 + h^{2k+2} + \tau^6 \right] \tau. \quad (\text{A.8})$$

Finally we collect up the above estimates and take small enough ε to complete the proof of this lemma. \square

A.4. Proof of Lemma 5.4

Recalling the definition of generalized E-fluxes, we need to move the focus of \mathbf{Q} from the element center \mathbf{u}_c^n to the rotation position $\mathbf{s}_\kappa^{n,\ell}$, for $\kappa = 1, 2$. After this is done, the term \mathcal{Y}_3 is changed to $\tilde{\mathcal{Y}}_3$. The difference between them is easy to bound. Since at every element boundary point there holds $\|\mathbf{Q}(\mathbf{u}_c^n) - \mathbf{Q}(\mathbf{s}_\kappa^{n,\ell})\|_{\text{M}} \leq C[h + \tau + \|\mathbf{e}^{n,\ell}\|_\infty]$, as stated in Section 4.1.2, the inverse properties (i) and (ii) imply that

$$|\mathcal{Y}_3 - \tilde{\mathcal{Y}}_3| \leq C\Psi_{2,2}^n \tau, \quad (\text{A.9})$$

due to $\tau = \mathcal{O}(h)$, where $\Psi_{2,2}^n$ has been given in (5.8). Hence we only need to estimate the new term $\tilde{\mathcal{Y}}_3$, which has the decomposition $\tilde{\mathcal{Y}}_3 = \tilde{\mathcal{Y}}_3(\boldsymbol{\eta}) - \tilde{\mathcal{Y}}_3(\mathbf{u}_h)$. Here

$$\tilde{\mathcal{Y}}_3(\mathbf{q}) = \tau [\mathbb{G}_2^n \mathbf{q}^n]^\top \mathbb{G}_1^n (\mathcal{A}^n [\mathbf{u}_h^n]) + \tau [\mathbb{G}_1^n \mathbf{q}^n]^\top \mathbb{G}_2^n (\mathcal{A}^n [\mathbf{u}_h^n]). \quad (\text{A.10})$$

Below we will estimate them separately.

The typical term in $\tilde{\mathcal{Y}}_3(\boldsymbol{\eta})$ is $[\mathbb{G}_\kappa^n \boldsymbol{\eta}]^\top \mathcal{A}^{n,\ell} [\mathbf{u}_h^{n,\ell}]$, where $\kappa = 1, 2$ and $\ell = 0, 1, 2$. Recall that all the elements in $\mathcal{A}^{n,\ell} \equiv \mathcal{A}(\hat{\mathbf{f}}; \mathbf{u}_h^{n,\ell})$ are bounded uniformly. Using the Cauchy–Schwarz inequality and the inverse property (ii), it is easy to see that

$$|\tau [\mathbb{G}_\kappa^n \boldsymbol{\eta}]^\top \mathcal{A}^{n,\ell} [\mathbf{u}_h^{n,\ell}]| \leq C \|\mathbb{G}_\kappa^n \boldsymbol{\eta}^n\|_{\Gamma_h} \left[\|\boldsymbol{\xi}^{n,\ell}\|_{\Gamma_h} + \|\boldsymbol{\eta}^{n,\ell}\|_{\Gamma_h} \right] \tau \leq C \left[\|\boldsymbol{\xi}^{n,\ell}\| + h^{2k+2} \right] \tau,$$

due to $\tau = \mathcal{O}(h)$, where the projection properties (3.6a) and (3.6c) are used also. Hence we will have

$$|\tilde{\mathcal{Y}}_3(\boldsymbol{\eta})| \leq C \left[\sum_{\ell=0,1,2} \|\boldsymbol{\xi}^{n,\ell}\|^2 + h^{2k+2} \right] \tau \leq C\Psi_{2,2}^n \tau. \quad (\text{A.11})$$

In order to estimate sharply the term $\tilde{\mathcal{Y}}_3(\mathbf{u}_h)$, we have to dig out the complete information about \mathbb{G}_2 , and cope with the difficulties resulting from the changing of $\mathcal{A}^{n,\ell}$ at different time stages. Define $\mathcal{B}^{n,\ell} = \mathcal{A}^{n,\ell} - \mathcal{A}^n$. A simple manipulation yields that

$$\tilde{\mathcal{Y}}_3(\mathbf{u}_h) = -\tilde{\mathcal{Y}}_3^{(1)} + \tilde{\mathcal{Y}}_3^{(2)} + \tilde{\mathcal{Y}}_3^{(3)},$$

where

$$\tilde{\Upsilon}_3^{(1)} = \tau[\mathbf{u}_h^n]^\top \mathcal{A}^n [\mathbb{G}_2^n \mathbf{u}_h^n] + \tau[\mathbb{G}_2^n \mathbf{u}_h^n]^\top \mathcal{A}^n [\mathbf{u}_h^n], \quad (\text{A.12a})$$

$$\tilde{\Upsilon}_3^{(2)} = \tau[\mathbb{G}_2^n \mathbf{u}_h^n]^\top \mathcal{A}^{n,1} [\mathbf{u}_h^{n,1}] + \tau[\mathbf{u}_h^{n,1}]^\top \mathcal{A}^{n,1} [\mathbb{G}_2^n \mathbf{u}_h^n], \quad (\text{A.12b})$$

$$\tilde{\Upsilon}_3^{(3)} = 2\tau[\mathbb{G}_1^n \mathbf{u}_h^n]^\top \mathcal{B}^{n,2} [\mathbf{u}_h^{n,2}] - \tau[\mathbb{G}_1^n \mathbf{u}_h^n]^\top \mathcal{B}^{n,1} [\mathbf{u}_h^{n,1}] - \tau[\mathbf{u}_h^{n,1}]^\top \mathcal{B}^{n,1} [\mathbb{G}_2^n \mathbf{u}_h^n]. \quad (\text{A.12c})$$

The first two terms include every term in the same sup-index between the matrices at least once. Here we do not use the symmetric property of $\mathcal{A}^{n,\ell}$.

Since $\mathcal{A}^{n,\ell} = \mathcal{S}^{n,\ell} - \mathcal{M}^{n,\ell}$, an application of Young's inequality yields

$$\begin{aligned} |\tilde{\Upsilon}_3^{(1)}| &\leq \varepsilon[\mathbf{u}_h^n]^\top \mathcal{S}^n [\mathbf{u}_h^n] \tau + K\varepsilon^{-1}[\mathbb{G}_2^n]^\top \mathcal{S}^n [\mathbb{G}_2^n] \tau + C\Psi_{2,2}^n \tau \\ &\leq \varepsilon[\mathbf{u}_h^n]^\top \mathcal{S}^n [\mathbf{u}_h^n] \tau + K\varepsilon^{-1} \mu F_{\max} \tau h^{-1} \|\mathbb{G}_2^n\|_n^2 + C\Psi_{2,2}^n \tau. \end{aligned} \quad (\text{A.13})$$

The above analysis is almost the same as that in Lemma 4.8. The second term $\tilde{\Upsilon}_3^{(2)}$ can be bounded similarly. It follows from assumption [A2](#) and the continuity of the reference stage solutions that

$$\|\mathcal{B}^{n,\ell}\|_{\mathbb{M}} \leq C_\star \|(\mathbf{u}_h^{n,\ell} - \mathbf{u}_h^n)^\pm\|_{\mathbb{M}} \leq C_\star (\|\mathbf{e}^{n,\ell}\|_\infty + \|\mathbf{e}^n\|_\infty + \tau), \quad \ell = 1, 2.$$

Then the inverse property (ii) implies $|\tilde{\Upsilon}_3^{(3)}| \leq C\Psi_{2,2}^n \tau$.

Finally, collecting up the above conclusions completes the proof of this lemma. \square

A.5. Proof of Lemma 5.5

The proof is almost the same as that for the error functional with respect to the single DG discretization (see Sect. 4). The key here is the estimate to the linear part, where the information about $\mathbb{G}_2^n \equiv \mathbb{G}_2^n \boldsymbol{\xi}^n$ will be found explicitly.

Along the same line as in Lemma 5.3, we fix the flow speed to be the same and rewrite the linear part in the equivalent form $\mathbb{G}_2^n \mathcal{L}^n(\mathbf{Q}^n \mathbf{v}_h) = \Omega_1 - \Omega_2 + \Omega_3$, where

$$\Omega_1 = \mathcal{L}_{\text{main}}(\mathbf{f}'_{\mathbf{u}}(\mathbf{u}^n); \mathbb{G}_2^n, \mathbf{Q}_c^n \mathbf{v}_h), \quad (\text{A.14a})$$

$$\Omega_2 = \mathcal{L}_{\text{main}}(\mathbf{f}'_{\mathbf{u}}(\mathbf{u}^n); \mathbb{G}_2^n \boldsymbol{\eta}^n, \mathbf{Q}_c^n \mathbf{v}_h), \quad (\text{A.14b})$$

$$\Omega_3 = 2\mathcal{L}_{\text{main}}(\mathbf{N}^{n,2}; \mathbf{e}^{n,2}, \mathbf{Q}_c^n \mathbf{v}_h) - \mathcal{L}_{\text{main}}(\mathbf{N}^{n,1}; \mathbf{e}^{n,1}, \mathbf{Q}_c^n \mathbf{v}_h). \quad (\text{A.14c})$$

Here $\mathbf{N}^{n,\ell} = \mathbf{f}'_{\mathbf{u}}(\mathbf{u}^{n,\ell}) - \mathbf{f}'_{\mathbf{u}}(\mathbf{u}^n)$ for $\ell = 1, 2$, as stated in Section A.3. There holds $\|\mathbf{N}^{n,\ell}\|_\infty = \mathcal{O}(\tau)$.

Using Lemmas 4.1, it follows from Cauchy-Schwarz inequality and Young's inequality, that

$$|\tau\Omega_1| \leq K\varepsilon^{-1} \lambda_{\max}^2 \|\mathbb{G}_2^n\|_n^2 + \varepsilon \|\mathbf{v}_h\|_n^2 + C_\star \left[\|\mathbb{G}_2^n\|^2 + \|\mathbf{v}_h\|^2 \right] \tau. \quad (\text{A.15})$$

Along the same line, we can get $|\tau\Omega_2| \leq \varepsilon \|\mathbf{v}_h\|_n^2 + C_\star \|\mathbf{v}_h\|^2 \tau + Ch^{2k+2} \tau^2$. Furthermore, we can estimate Ω_3 along the same line as that for Lemmas 4.1 and 4.3. Using the inverse properties (i) and (ii), together with approximation property (3.6a), we have $|\tau\Omega_3| \leq C_\star (\|\mathbf{v}_h\|^2 + \|\boldsymbol{\xi}^{n,1}\|^2 + \|\boldsymbol{\xi}^{n,2}\|^2 + h^{2k+2}) \tau$. Finally, the sum of the three inequalities above gives the result to the linear part. Namely,

$$\tau \mathbb{G}_2^n \mathcal{L}^n(\mathbf{Q}_c^n \mathbf{v}_h) \leq \varepsilon \|\mathbf{v}_h\|_n^2 + K\varepsilon^{-1} \lambda_{\max}^2 \|\mathbb{G}_2^n\|_n^2 + C_\star \|\mathbf{v}_h\|^2 \tau + C\Psi_{2,2}^n \tau. \quad (\text{A.16})$$

The next three parts are easy to estimate. Lemma 4.5 gives the estimate to the nonlinear part in the form

$$\tau \mathbb{G}_2^n \mathcal{N}^n(\mathbf{Q}_c^n \mathbf{v}_h) \leq C \|\mathbf{v}_h\|^2 \tau + C\Psi_{2,2}^n \tau. \quad (\text{A.17})$$

In order to estimate the viscosity part, the focus of the local rotation matrix \mathbf{Q} will be moved from \mathbf{u}_c^n to $\mathbf{u}_b^{n,\ell}$, and then be moved back. Using mainly Lemma 4.8, we have

$$\begin{aligned} \tau \mathbb{G}_2^n \mathcal{V}^n(\mathbf{Q}_c^n \mathbf{v}_h) &\leq \varepsilon \|\mathbf{v}_h\|_n^2 + C\varepsilon^{-1} F_{\max} \mu \tau h^{-1} \sum_{\ell=0,1,2} \tau [\mathbf{u}_h^{n,\ell}]^\top \mathcal{S}^{n,\ell} [\mathbf{u}_h^{n,\ell}] \\ &\quad + \sum_{\ell=0,1,2} C_* h^{-1} \|e^{n,\ell}\|_\infty \left[\sum_{\ell=0,1,2} \|\boldsymbol{\xi}^{n,\ell}\|^2 + \|\mathbf{v}_h\|^2 + h^{2k+2} \right] \tau. \end{aligned} \tag{A.18}$$

Note that we have used a new parameter ε . It is easy to bound the time-marching part, in the form

$$\tau \mathbb{G}_2^n \mathcal{T}^n(\mathbf{Q}_c^n \mathbf{v}_h) \leq \varepsilon \|\mathbf{v}_h\|^2 \tau + C(h^{2k+2} + \tau^6) \tau. \tag{A.19}$$

Here we have also lost a factor τ , since $\tau < 1$ as we have assumed.

Finally, we collect up the above four estimates and complete the proof of this lemma, since $\lambda_{\max} \leq 1$. \square

REFERENCES

- [1] R.A. Adams, Sobolev Spaces. Academic Press, New York (1975).
- [2] E. Burman, A. Ern and M.A. Fernandez, Explicit Runge–Kutta schemes and finite elements with symmetric stabilization for first-order linear PDE systems. *SIAM. J. Numer. Anal.* **48** (2010) 2019–2042.
- [3] P.G. Ciarlet, Finite Element Method for Elliptic Problems. North–Holland, Amsterdam (1978).
- [4] B. Cockburn and J. Guzmán, Error estimates for the Runge–Kutta discontinuous Galerkin method for the transport equation with discontinuous initial data. *SIAM. J. Numer. Anal.* **46** (2008) 1364–1398.
- [5] B. Cockburn and C.-W. Shu, TVB Runge–Kutta local projection discontinuous Galerkin finite element method for scalar conservation laws II: General framework. *Math. Comput.* **52** (1989) 411–435.
- [6] B. Cockburn and C.-W. Shu, The Runge–Kutta local projection P^1 -discontinuous Galerkin method for scalar conservation laws. *RAIRO Modél. Math. Anal. Numér.* **25** (1991) 337–361.
- [7] B. Cockburn and C.-W. Shu, The Runge–Kutta discontinuous Galerkin finite element method for conservation laws V: Multidimensional systems *J. Comput. Phys.* **141** (1998a) 199–224.
- [8] B. Cockburn and C.-W. Shu, The local discontinuous Galerkin method for time-dependent convection-diffusion systems *SIAM. J. Numer. Anal.* **35** (1998b) 2440–2463.
- [9] B. Cockburn and C.-W. Shu, Runge–Kutta discontinuous Galerkin methods for convection-dominated problems *J. Sci. Comput.* **16** (2001) 173–261.
- [10] B. Cockburn, S. Hou, and C.-W. Shu, TVB Runge–Kutta local projection discontinuous Galerkin finite element method for conservation laws IV: The multidimensional case. *Math. Comp.* **54** (1990) 545–581.
- [11] B. Cockburn, C. Johnson, C.-W. Shu and E. Tadmor, An introduction to the discontinuous Galerkin method for convection-dominated problems, in *Advanced Numerical Approximation of Nonlinear Hyperbolic Equations*, edited by Quarteroni. Vol. 1697 of *Lect. Notes Math.* Springer, Berlin (1998) 151–268.
- [12] B. Cockburn, S.-Y. Lin, and C.-W. Shu, TVB Runge–Kutta local projection discontinuous Galerkin finite element method for conservation laws III: One dimensional systems. *J. Comput. Phys.* **84** (1989) 90–113.
- [13] G.H. Golub and C.F. Van Loan, Matrix Computations. Posts and Telecom Press (2011).
- [14] A. Harten, On the symmetric form of systems of conservation laws with entropy. *J. Comput. Phys.* **49** (1983a) 151–164.
- [15] A. Harten, High resolution schemes for hyperbolic conservation laws. *J. Comput. Phys.* **49** (1983b) 357–393.
- [16] S.-M. Hou and X.-D. Liu, Solutions of multidimensional hyperbolic systems of conservation laws by square entropy condition satisfying discontinuous Galerkin method. *J. Sci. Comput.* **31** (2007) 127–151.
- [17] G.-S. Jiang and C.-W. Shu, On cell entropy inequality for discontinuous Galerkin methods. *Math. Comp.* **62** (1994) 531–538.
- [18] C. Johnson and J. Pitkäranta, An analysis of the discontinuous Galerkin method for a scalar hyperbolic equation. *Math. Comp.* **46** (1986) 1–26.
- [19] R. Kress, Numerical analysis. Springer-Verlag (1998).
- [20] P. Lesaint and P.A. Raviart, On a finite element method for solving the neutron transport equation, in *Mathematical Aspects of Finite Elements in Partial Differential Equations*, edited by C. de Boor. Academic Press, New York (1974) 89–145.
- [21] J. Luo, *A priori error estimates to Runge–Kutta discontinuous Galerkin finite element method for symmetrizable system of conservation laws with sufficiently smooth solutions*. Ph.D. thesis, Nanjing University (2013).
- [22] S. Osher, Riemann solvers, the entropy condition, and difference approximations. *SIAM. J. Numer. Anal.* **21** (1984) 217–235.

- [23] W.H. Reed and T.R. Hill, *Triangular mesh methods for the neutron transport equation*. Los Alamos Scientific Laboratory report LA-UR-73-479, Los Alamos, NM (1973).
- [24] P.L. Roe, Approximate Riemann solvers, parameter vectors, and difference schemes. *J. Comput. Phys.* **43** (1981) 357–372.
- [25] C.W. Schulz-Rinne, Classification of the Riemann problem for two-dimensional gas dynamics. *SIAM J. Math. Anal.* **24** (1993) 76–88.
- [26] C.-W. Shu and S. Osher, Efficient implementation of essentially non-oscillatory shock capturing schemes. *J. Comput. Phys.* **77** (1988) 439–471.
- [27] E.F. Toro, *Riemann solvers and numerical methods for fluid dynamics: a practical introduction*. Springer (2009).
- [28] Q. Zhang, Third order explicit Runge–Kutta discontinuous Galerkin method for linear conservation law with inflow boundary condition. *J. Sci. Comput.* **46** (2010) 294–313.
- [29] Q. Zhang and C.-W. Shu, Error estimates to smooth solution of Runge–Kutta discontinuous Galerkin methods for scalar conservation laws. *SIAM. J. Numer. Anal.* **42** (2004) 641–666.
- [30] Q. Zhang and C.-W. Shu, Error estimates to smooth solution of Runge–Kutta discontinuous Galerkin methods for symmetrizable system of conservation laws. *SIAM. J. Numer. Anal.* **44** (2006) 1702–1720.
- [31] Q. Zhang and C.-W. Shu, Stability analysis and a priori error estimates to the third order explicit Runge–Kutta discontinuous Galerkin Method for scalar conservation laws. *SIAM. J. Numer. Anal.* **48** (2010) 1038–1063.
- [32] Q. Zhang and C.-W. Shu, Error estimates for the third order explicit Runge–Kutta discontinuous Galerkin method for a linear hyperbolic equation in one-dimension with discontinuous initial data. *Numer. Math.* **126** (2014) 703–740.