

Applications of Mathematics

Abdulatif Badenjki; Gerald G. Warnecke

Theoretical and numerical studies of the $P_N P_M$ DG schemes in one space dimension

Applications of Mathematics, Vol. 64 (2019), No. 6, 599–635

Persistent URL: <http://dml.cz/dmlcz/147924>

Terms of use:

© Institute of Mathematics AS CR, 2019

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

THEORETICAL AND NUMERICAL STUDIES OF THE $P_N P_M$ DG
SCHEMES IN ONE SPACE DIMENSION

ABDULATIF BADENJKI, GERALD WARNECKE, Magdeburg

Received August 18, 2018. Published online November 20, 2019.

Abstract. We give a proof of the existence of a solution of reconstruction operators used in the $P_N P_M$ DG schemes in one space dimension. Some properties and error estimates of the projection and reconstruction operators are presented. Then, by applying the $P_N P_M$ DG schemes to the linear advection equation, we study their stability obtaining maximal limits of the Courant numbers for several $P_N P_M$ DG schemes mostly experimentally. A numerical study explains how the stencils used in the reconstruction affect the efficiency of the schemes.

Keywords: $P_N P_M$ DG scheme; piecewise polynomial; projection; reconstruction; least square; local continuous space time Galerkin method; discontinuous Galerkin; advection equation; conservation law; von Neumann stability analysis; time discretization

MSC 2010: 65M12, 65M60, 33C45

1. INTRODUCTION

A conservation law [6] is a system of hyperbolic PDEs that states that the rate of change of a physical state or conserved quantity is governed by a flux function. The initial value problem for a scalar conservation law in one space dimension is given by

$$\begin{aligned}v_t(t, x) + f(v(t, x))_x &= 0 \quad \text{with } (t, x) \in (0, \infty) \times \mathbb{R}, \\v(0, x) &= v_0(x) \quad \text{for } x \in \mathbb{R},\end{aligned}$$

where the dependent variable v is the conserved quantity. The differentiable function $f: \mathbb{R} \rightarrow \mathbb{R}$ is the flux function and $v_0: \mathbb{R} \rightarrow \mathbb{R}$ describes the initial data of v .

In applications one is interested in systems of such equations in one, two or three space dimensions. Since generally exact solutions to these systems are not known, numerical methods are important for solving such initial value problems. The aim of

this paper is to make some basic comparisons between numerical schemes from a large class of schemes. For this purpose it suffices to consider special scalar problems in one space dimension for which the exact solution is known. Therefore, we restrict ourselves accordingly.

Let $M, N \in \mathbb{N}_0$ satisfy $N \leq M$. The $P_N P_M$ DG schemes are a class of arbitrary high order schemes originally developed by Dumbser et al. [5] for systems in two and three space dimensions. These schemes use a reconstruction operator applied to the DG scheme at each time step to increase the accuracy. By varying the parameters N and M as well as the stencils for the reconstruction one obtains a very large number of schemes as M is increased. In this paper these schemes are applied to scalar conservation laws in one space dimension in order to systematically study some basic properties of these schemes.

In order to describe these schemes take $T > 0$ and $I = [a, b]$. The conserved functions will be from the function space

$$L^\infty([0, T], L^2(I)) = \left\{ v: [0, T] \times I \rightarrow \mathbb{R}; \operatorname{ess\,sup}_{t \in [0, T]} \|v(t, \cdot)\| < \infty \right\},$$

where $\|\cdot\|$ is the L^2 norm on I . For $f \in L^2(I)$ we have $\|f\| = (\int_I |f(x)|^2 dx)^{1/2}$.

Let $Z \in \mathbb{N}$. We discretize I by equally distant points $a = x_{1/2} < x_{3/2} < \dots < x_{Z+1/2} = b$ into Z subintervals $I_j = [x_{j-1/2}, x_{j+1/2}]$ for $j = 1, \dots, Z$ with the constant mesh size $h = x_{j+1/2} - x_{j-1/2} = (b - a)/Z$. We define spaces of the polynomials

$$P_{N, I_j} := \{p: I_j \rightarrow \mathbb{R}, p \text{ is a polynomial of maximal degree } N\}.$$

The numerical solutions will be taken from the following space of the piecewise polynomials:

$$P_{N, I, Z} := \{p: p|_{I_j} \in P_{N, I_j} \forall j = 1, \dots, Z\}.$$

Let $Z_1 \in \mathbb{N}$. We discretize the time interval $[0, T]$ by considering the times $0 = t_0 < t_1 < \dots < t_{Z_1} = T$ and define subintervals $T_n = [t_n, t_{n+1}[$ and a constant time step $\Delta t = t_{n+1} - t_n$ for $n = 0, \dots, Z_1 - 1$. First we make the *projection* of the initial data $v_0 \rightarrow u^0$, where $u^0 \in P_{N, I, Z}$ is a piecewise polynomial of degree N . Using this piecewise polynomial we iterate for $n \in \{0, \dots, Z_1 - 1\}$ the following steps:

- (1) The *reconstruction* $u^n \rightarrow w^n$, where $w^n \in P_{M, I, Z}$ is a piecewise polynomial of degree M .
- (2) The *time evolution* $w^n \rightarrow U^n$, where $U^n \in P_{M, T_n \times I, Z}$ is a piecewise polynomial of degree M in time and space. In this step the local continuous space-time Galerkin method introduced by Dumbser et al. [5] is used for the time evolution. The values in time are used to compute the higher order fluxes in the next step.

- (3) The *DG scheme* of order $N + 1$ giving $U^n \rightarrow u^{n+1}$ where $u^{n+1} \in P_{N,I,Z}$ is a piecewise polynomial of degree N . It is the new numerical solution after the time step Δt .

A central result of this paper is a proof of the unique solvability of the reconstruction step. Our choices of the stencil sizes for the reconstruction operators always generate systems of equations with full column rank. Furthermore, the reconstruction operators give approximations of the data considered, but as a special case they recover the same data when these originally are polynomials of the same degree.

Courant numbers are important for the stability of explicit schemes for conservation laws. We computationally explore maximal limits of these numbers for the $P_N P_M$ DG schemes by applying the von Neumann analysis and using an experimental procedure. We obtain a wide variety of stability limits, including some unstable cases for which we have only one value $\lambda = 1$ that gives a stable solution. Moreover, there are some semi-stable cases with a minimal bound on the time step and some cases with a stability interval larger than $]0, 1]$. This study of the stability uses the application of the $P_N P_M$ DG schemes to the linear advection equation. It is common to do this in theory and practice.

In the numerical literature some authors, e.g. Dumbser [5], usually determine these limits using the linear cases. Then they use a rate 0.8 or 0.7 of these limits for nonlinear cases. In this manner, one can try to extend the use of our computed limits of stability to nonlinear cases.

The reconstruction operator needs for its definition stencils related to the discretization. We study the numerical effect of the size and form of these stencils on the efficiency of the $P_N P_M$ DG schemes. A summary of the numerical results is given in the conclusion.

The paper is arranged as follows. In Section 2 we recall projection operators and the error estimate for these operators. In Section 3 we introduce the reconstruction operators in detail and prove the existence of solutions to the reconstruction problem. This proof is a first general proof of this fact which previously had been obtained for the special cases $M = 2N + 1$ only, see [7]. We consider two cases of the solution, either a unique exact solution or a unique solution obtained by using the least squares approach in the overdetermined case. Furthermore, some properties of the reconstruction and an error estimate are given. Then in Section 4, the continuous Galerkin scheme is recalled. It evolves the data in the time up to the same order of the space for step 2 above. Then, we display the 3rd step above in Section 5. Moreover, the $P_N P_M$ DG schemes are applied to the advection equation in Section 6. We study the stability and give maximal limits of the Courant numbers. Finally, numerical results for the effect of the stencils on the efficiency are given.

2. THE PROJECTION ONTO PIECEWISE POLYNOMIALS

The mutually orthogonal Legendre polynomials of degree $i \in \mathbb{N}_0$ on the reference interval $J = [-1, 1]$ can be determined by the Rodrigues formula

$$\mathcal{L}_i(s) = \frac{(-1)^i}{2^i i!} \frac{d^i}{ds^i} \{(1 - s^2)^i\}$$

for $s \in J$, see e.g. Stegun [10]. On J they satisfy the orthogonality condition

$$\int_{-1}^1 \mathcal{L}_m(s) \mathcal{L}_n(s) ds = \frac{2}{2n + 1} \delta_{mn}$$

where δ_{mn} is the Kronecker delta, and satisfy also

$$(2.1) \quad \mathcal{L}_i(-s) = (-1)^i \mathcal{L}_i(s) \quad \text{and} \quad \mathcal{L}_i(1) = 1,$$

see e.g. Koornwinder et al. [9], Table 18.6.1. For example, the first four polynomials are 1, s , $(3s^2 - 1)/2$, $(5s^3 - 3s)/2$.

Let I_j for $j = 1, \dots, Z$ be our discrete intervals with constant length h and midpoints x_j . We define linear reference transformations $\gamma_j: I_j \rightarrow J$ by $\gamma_j(x) := 2h^{-1}(x - x_j)$ for $x \in I_j$. Using these transformations, we obtain the transformed piecewise Legendre basis functions $\Phi_{i,j}: I \rightarrow \mathbb{R}$ by

$$(2.2) \quad \Phi_{i,j}(x) = \begin{cases} \mathcal{L}_i(\gamma_j(x)), & x \in I_j, \\ 0, & x \in I \setminus I_j, \end{cases} \quad j = 1, \dots, Z, \quad i = 0, \dots, N.$$

The set $B_{N,Z} := \{\Phi_{i,j}; j = 1, \dots, Z, i = 0, \dots, N\}$ is an orthogonal basis of the solution space $P_{N,I,Z}$.

Let $v \in L^2(I)$. Using these basis functions $\Phi_{i,j}$ the coefficients

$$(2.3) \quad \hat{u}_{i,j} = \frac{2i + 1}{h} \int_{I_j} v(x) \Phi_{i,j}(x) dx$$

give the L^2 projection operator $\Pi_{N,Z}: L^2(I) \rightarrow P_{N,I,Z}$ with $\Pi_{N,Z}(v) = u$ by the formula

$$\Pi_{N,Z}(v)(x) = u(x) = \sum_{j=1}^Z \sum_{i=0}^N \hat{u}_{i,j} \Phi_{i,j}(x) \quad \text{for } x \in I.$$

The operator $\Pi_{N,Z}$ has the following well-known properties. It is linear, and *idempotent*, i.e. a *projection*; this means that $\Pi_{N,Z}(\Pi_{N,Z}(v)) = \Pi_{N,Z}(v)$. This implies that, for $p \in P_{N,I}$, $\Pi_{N,Z}(p) = p$. It is an *orthogonal projection* of v on $P_{N,I,Z}$ with respect

to the L^2 scalar product, i.e. $\langle v - \Pi_{N,Z}(v), \Phi_{i,j} \rangle = 0$ for all $\Phi_{i,j} \in B_{N,Z}$. Also, it is the *best approximation* using piecewise polynomials, i.e. it satisfies $\|v - \Pi_{N,Z}(v)\|_{L^2(I)} \leq \|v - p\|_{L^2(I)}$ for all $p \in P_{N,I,Z}$. Also the following boundedness estimate holds: $\|\Pi_{N,Z}(v)\|_{L^2(I)} \leq \|v\|_{L^2(I)}$. Taking $\hat{\mathbf{u}}_j := (\hat{u}_{0,j}, \dots, \hat{u}_{N,j})^\top$ and using the Euclidean norm $\|\cdot\|_e$ in \mathbb{R}^{N+1} , the solution satisfies the estimates

$$(2.4) \quad \frac{h}{2N+1} \|\hat{\mathbf{u}}_j\|_e^2 \leq \|u_j\|_{L^2(I_j)}^2 \leq h \|\hat{\mathbf{u}}_j\|_e^2.$$

For the projection error estimates we use functions in the Sobolev spaces of order $r \in \mathbb{N}_0$, $W^{r,2}(I) = \{v \in L^2(I) : D^{(r)}v \in L^2(I)\}$. These spaces are associated with the norm $\|\cdot\|_{W^{r,2}(I)}$ and the seminorm $|\cdot|_{W^{r,2}(I)}$ where

$$(2.5) \quad \|v\|_{W^{r,2}(I)} = \sqrt{\sum_{i=0}^r (\|D^{(i)}v\|_{L^2(I)})^2}, \quad |v|_{W^{r,2}(I)} = \|D^{(r)}v\|_{L^2(I)},$$

for all $v \in W^{r,2}(I)$, see e.g. Adams [1]. For each $v \in W^{N+1,2}(I)$, the following error estimates hold:

$$(2.6) \quad \begin{aligned} \|\Pi_{N,Z}(v) - v\|_{L^2(I_j)}^2 &\leq C_2^2 h^{2N+2} |v|_{W^{N+1,2}(I_j)}^2, \\ \|\Pi_{N,Z}(v) - v\|_{L^2(I)} &\leq C_2 h^{N+1} |v|_{W^{N+1,2}(I)}, \\ \|\Pi_{N,Z}(v) - v\|_{L^1(I)} &\leq C_3 h^{N+1} |v|_{W^{N+1,2}(I)}. \end{aligned}$$

In a special case, when v is a bounded function with a discontinuity in I and h is the length of I , we have for all $1 \leq p < \infty$ the estimate

$$\|\Pi_{N,Z}(v) - v\|_{L^p(I)} \leq C_3 \|v\|_{L^\infty(I)} h^{1/p}.$$

For the proofs of all these properties and estimates, see [2], Chapter 2.

3. THE RECONSTRUCTION OPERATORS

The reconstruction is an approximation by higher degree polynomials obtained from a set of neighboring lower degree polynomials. This approximation has as a building block a chosen reconstruction stencil of neighboring elements. With $R, L \geq 0$ and $n_e := 1 + L + R$, we define the stencil $S_{I_j, n_e, L} = \bigcup_{c=-L}^R I_{j+c}$ which is related to the element I_j . It is constituted of I_j together with L elements to the left and R elements to the right of I_j , see Figure 1. Using all such stencils, for

$j = 1, \dots, Z$, we obtain the corresponding extended interval $I_{\text{ex}} := \bigcup_{j=1}^Z S_{I_j, n_e, L}$. It has $Z_{\text{ex}} := Z + L + R$ elements, and $I \subset I_{\text{ex}}$.

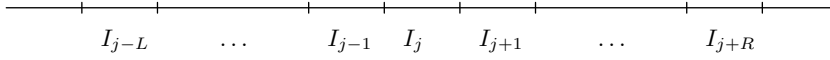


Figure 1. The stencil $S_{I_j, n_e, L}$.

For example, taking $n_e = 3$, we have three stencils $S_{I_j, 3, L}$ with $L = 0, 1, 2$, see Figure 2. The extended interval I_{ex} has $Z + 2$ elements. Figure 3 shows the extended interval I_{ex} for the case $L = R = 1$.

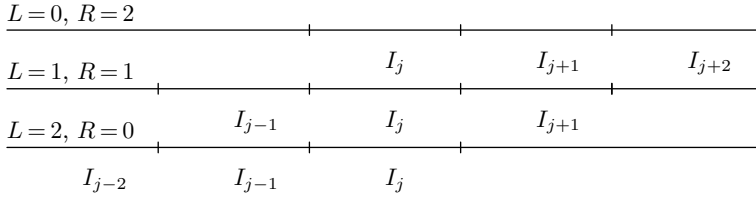


Figure 2. The stencils $S_{I_j, 3, 0}$, $S_{I_j, 3, 1}$, and $S_{I_j, 3, 2}$.

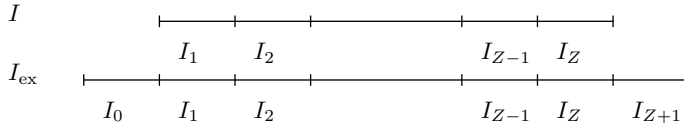


Figure 3. The extended interval I_{ex} for the case $L = R = 1$.

3.1. Definition of the reconstruction operator. Now, let $v \in L^2(I_{\text{ex}})$ be any given function, $u = \Pi_{N, Z_{\text{ex}}}(v) \in P_{N, I_{\text{ex}}, Z_{\text{ex}}}$ the projection of this function to the piecewise polynomials of degree N , and let $S_{I_j, n_e, L}$ with $L \in \{0, \dots, n_e - 1\}$ be the stencil for the reconstruction. The reconstruction operator $\mathfrak{R}_{N, M, S, Z}: P_{N, I_{\text{ex}}, Z_{\text{ex}}} \rightarrow P_{M, I, Z}$ will be defined to give a piecewise polynomial of degree $M \geq N$ that has the form

$$(3.1) \quad \mathfrak{R}_{N, M, S, Z}(u(x)) = w(x) := \sum_{j=1}^Z w_j(x) = \sum_{j=1}^Z \sum_{i=0}^M \hat{w}_{i, j} \Phi_{i, j}(x).$$

Only during the reconstruction step, each basis function $\Phi_{i, j}$ with $i = 0, \dots, M$ is extended from I_j over the whole stencil, by ignoring the condition of the definition (2.2) $\Phi_{i, j}(x) = 0$ for $x \in I \setminus I_j$. We use the notation $\Phi_{i, j}^e$ for the extended basis functions. For example, as shown in Figure 4, with the stencil $S_{I_j, 3, 1}$, we have

$$\Phi_{1, j}(x) = \begin{cases} \frac{2}{h}(x - x_j), & x \in I_j, \\ 0, & x \in I \setminus I_j, \end{cases} \quad \Phi_{1, j}^e(x) = \begin{cases} \frac{2}{h}(x - x_j), & x \in S_{I_j, n_e, L}, \\ 0, & x \in I_{\text{ex}} \setminus S_{I_j, n_e, L}. \end{cases}$$

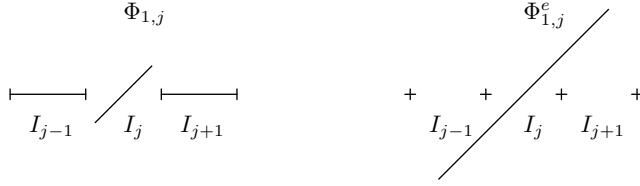


Figure 4. The basis function $\Phi_{1,j}$ (left) and its extension $\Phi_{1,j}^e$ (right) onto the stencil $S_{I_j,3,1}$.

Let $w_j := w|_{I_j}$ be the term of w on I_j , for $j = 1, \dots, Z$ fixed, let $S_{I_j, n_e, L}$ be a stencil of size $n_e = L + R + 1$ with $L, R \geq 0$. We set $c = -L, \dots, R$ and consider $I_{j+c} := [x_{j-1/2} + ch, x_{j+1/2} + ch[$ to be an element of this stencil. We compute the coefficients $\widehat{w}_{i,j}$ by assuming that the error in L^2 norm of computing w_j must be minimal on the stencil. This leads to the normal equations

$$\sum_{l=0}^M \widehat{w}_{l,j} \langle \Phi_{l,j}^e, \Phi_{k,j+c} \rangle_{j+c} = \sum_{i=0}^N \widehat{u}_{i,j+c} \langle \Phi_{i,j+c}, \Phi_{k,j+c} \rangle_{j+c}.$$

By orthogonality, this leads to the system of equations

$$(3.2) \quad \sum_{l=0}^M \widehat{w}_{l,j} \langle \Phi_{l,j}^e, \Phi_{k,j+c} \rangle_{j+c} = \frac{h}{2k+1} \widehat{u}_{k,j+c}.$$

This system consists of $n_e(N+1)$ equations with $M+1$ unknowns. We neglect the case of underdetermined systems, i.e., we require the condition

$$(3.3) \quad n_e(N+1) \geq (M+1)$$

to be satisfied when choosing any stencil. Two examples are shown in Appendix A. In the underdetermined case one could consider the unique solution of minimal Euclidean norm.

Due to the orthogonality of the Legendre basis functions we obtain the equalities

$$(3.4) \quad \widehat{w}_{i,j} = \widehat{u}_{i,j}, \quad i = 0, \dots, N, \quad j = 1, \dots, Z.$$

Then we can write $w(x) = u(x) + \sum_{j=1}^Z \sum_{i=N+1}^M \widehat{w}_{i,j} \Phi_{i,j}(x)$. We only have to determine the remaining degrees of freedom $\widehat{w}_{N+1,j}, \dots, \widehat{w}_{M,j}$ for $j = 1, \dots, Z$.

If $M = N$, the equalities (3.4) cover all $\widehat{w}_{i,j}$ for $i = 0, \dots, M$ and the reconstruction operator becomes the restriction $\mathfrak{R}_{M,M,S,Z}(u) = w = u|_I$. Then the $P_N P_N$ DG schemes are equivalent to the classical DG schemes, see [5].

3.2. The matrix form. For the left-hand side in (3.2) we define the coefficient vectors $\widehat{\mathbf{w}}_j := (\widehat{w}_{0,j}, \dots, \widehat{w}_{M,j})^\top \in \mathbb{R}^{M+1}$ and the matrices $\mathbf{M}_{j,c}$ as

$$\mathbf{M}_{j,c} := \begin{pmatrix} \langle \Phi_{0,j}^e, \Phi_{0,j+c} \rangle_{j+c} & \cdots & \langle \Phi_{M,j}^e, \Phi_{0,j+c} \rangle_{j+c} \\ \vdots & \vdots & \vdots \\ \langle \Phi_{0,j}^e, \Phi_{N,j+c} \rangle_{j+c} & \cdots & \langle \Phi_{M,j}^e, \Phi_{N,j+c} \rangle_{j+c} \end{pmatrix} \in \mathbb{R}^{(N+1) \times (M+1)}.$$

For the right-hand side we define the vectors $\widehat{\mathbf{u}}_{j+c} := (\widehat{u}_{0,j+c}, \dots, \widehat{u}_{N,j+c})^\top \in \mathbb{R}^{N+1}$ and the matrices \mathbf{A}_{j+c} as

$$\mathbf{A}_{j+c} = \begin{pmatrix} h & 0 & \cdots & 0 \\ 0 & \frac{h}{3} & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \frac{h}{2N+1} \end{pmatrix} \in \mathbb{R}^{(N+1) \times (N+1)}.$$

Then we can write the elemental matrix forms

$$(3.5) \quad \mathbf{M}_{j,c} \cdot \widehat{\mathbf{w}}_j = \mathbf{A}_{j+c} \cdot \widehat{\mathbf{u}}_{j+c} \quad \text{for } c = -L, \dots, R.$$

Taking $\mathbf{y}_{j,c} := \mathbf{A}_{j+c} \cdot \widehat{\mathbf{u}}_{j+c}$, these forms become $\mathbf{M}_{j,c} \cdot \widehat{\mathbf{w}}_j = \mathbf{y}_{j,c}$ for $c = -L, \dots, R$. Also defining vectors $\mathbf{y}_j := (\mathbf{y}_{j,-L}, \dots, \mathbf{y}_{j,R})^\top \in \mathbb{R}^{n_e(N+1)}$ and matrices $\mathbf{M}_j := (\mathbf{M}_{j,-L}, \dots, \mathbf{M}_{j,R})^\top \in \mathbb{R}^{n_e(N+1) \times (M+1)}$, we can merge the elemental matrix forms into the full matrix form

$$(3.6) \quad \mathbf{M}_j \cdot \widehat{\mathbf{w}}_j = \mathbf{y}_j,$$

which is related to the stencil $S_{I_j, n_e, L}$.

Lemma 3.1. *Suppose that $(p_n)_{n \in \mathbb{N}_0}$ is a sequence of orthogonal polynomials on the interval $[a, b]$ with p_n of degree n . Then, for each $k \in \{0, \dots, n\}$, the polynomial p_k has k simple zeros that lie in $]a, b[$.*

Proof. We consider p_n with the zeros $x_1^n, \dots, x_n^n \in \mathbb{C}$. We have $p_0 = c \neq 0$ and

$$0 = \langle p_0, p_n \rangle = c \int_a^b (x - x_1^n) \cdots (x - x_n^n) dx.$$

This means that p_n must have at least one real zero, e.g. x_* , in $]a, b[$, at which the polynomial p_n changes its sign. This zero must have an odd multiplicity. Let

$\Psi = \{x \in]a, b[: x \in \{x_1^n, \dots, x_n^n\} \text{ with odd multiplicities}\}$. Then we know that $\Psi \neq \emptyset$, since at least $x_* \in \Psi$. We set $\pi(x) := \prod_{t \in \Psi} (x - t)$. The function π has only simple zeros in $]a, b[$, since it is a product of different linear factors. Then the function $p_n \cdot \pi$ has in $]a, b[$ real zeros with even multiplicities only. This implies that $p_n \cdot \pi$ has no sign change on $]a, b[$. Thus we obtain $\langle p_n, \pi \rangle \neq 0$. Now we assume that $\pi \in P_l$ with $l < n$, i.e. $\pi = \sum_{j=0}^l a_j p_j$. Then $\langle p_n, \pi \rangle = \sum_{j=0}^l a_j \langle p_j, p_n \rangle = 0$. This is a contradiction to $\langle p_n, \pi \rangle \neq 0$. This means that $\pi \in \text{span}(p_n)$, thus $\pi = \lambda p_n$, for some $\lambda \in \mathbb{R}$. Consequently, p_n has only simple zeros all of which lie in $]a, b[$. \square

Corollary 3.1. *Suppose that on the interval $I = [a, b]$ we have a set $\{p_0, \dots, p_N\}$ of $N + 1$ orthogonal polynomials with $p_k \in P_{k,I}$ for $k = 0, \dots, N$. Suppose that $p \in P_{M,I}$ is a polynomial of degree $M > N$ which is orthogonal to all p_k . Then it follows, by a proof analogous to that of Lemma 3.1, that p has at least $N + 1$ different zeros on $]a, b[$.*

Now we prove the existence of the solution by depending on the rank of the matrix \mathbf{M}_j .

Theorem 3.1. *The matrix \mathbf{M}_j has a full column rank $M + 1$.*

Proof. We consider the homogeneous matrix form $\mathbf{M}_j \cdot \widehat{\mathbf{w}}_j = \mathbf{0}$ of (3.6). This system means that the coefficient vector of the reconstruction polynomial w_j , see (3.1), satisfies

$$\mathbf{M}_{j,c} \cdot \widehat{\mathbf{w}}_j = \begin{pmatrix} \langle w_j, \Phi_{0,j+c} \rangle_{j+c} \\ \vdots \\ \langle w_j, \Phi_{N,j+c} \rangle_{j+c} \end{pmatrix} = \mathbf{0}.$$

Therefore, the polynomial w_j of degree M is orthogonal to the $N + 1$ basis functions $\Phi_{i,j+c}$ on all elements I_{j+c} of the stencil $S_{I_j, n_e, L}$, with $i = 0, \dots, N$ and $c = -L, \dots, R$. According to Corollary 3.1, there are at least $N + 1$ different zeros of w on each element I_{j+c} . This gives $n_e(N + 1)$ different zeros on the whole stencil. Also according to the condition (3.3) we find $n_e(N + 1) \geq (M + 1) > M$. Therefore, w is the zero polynomial. This proves the injectivity of the reconstruction and implies that the matrix \mathbf{M}_j has the full column rank $M + 1$. \square

3.3. The solution of the reconstruction problem. For $c = 0$, the equations (3.2) directly give the equalities (3.4), due to the orthogonality of the basis functions on I_j . Thus we can ignore the equations related to I_j in the system (3.2). We now consider the corresponding reduced system. Defining vectors $\widehat{\mathbf{y}}_j := (\mathbf{y}_{j,-L}, \dots, \mathbf{y}_{j,-1}, \mathbf{y}_{j,1}, \dots, \mathbf{y}_{j,R})^\top \in \mathbb{R}^{(n_e-1)(N+1)}$ and matrices $\widetilde{\mathbf{M}}_j :=$

$(\mathbf{M}_{j,-L}, \dots, \mathbf{M}_{j,-1}, \mathbf{M}_{j,1}, \dots, \mathbf{M}_{j,R})^\top \in \mathbb{R}^{(n_e-1)(N+1) \times (M+1)}$, the reduced system is given by

$$\widetilde{\mathbf{M}}_j \cdot \widehat{\mathbf{w}}_j = \widehat{\mathbf{y}}_j.$$

The vector $\widehat{\mathbf{w}}_j$ can be divided into two vectors, $\widehat{\mathbf{u}}_j := (\widehat{u}_{0,j}, \dots, \widehat{u}_{N,j})^\top \in \mathbb{R}^{N+1}$ of the known coefficients and $\widehat{\mathbf{x}}_j := (\widehat{w}_{N+1,j}, \dots, \widehat{w}_{M,j})^\top \in \mathbb{R}^{M-N}$ of unknown coefficients. Moreover, the first $N+1$ columns in each matrix $\widetilde{\mathbf{M}}_j$ are related to the known coefficients. Thus the matrices $\widetilde{\mathbf{M}}_j$ can be divided into two parts, in the form $\widetilde{\mathbf{M}}_j = (\widetilde{\mathbf{M}}_{j,1}, \widetilde{\mathbf{M}}_{j,2})$ where $\widetilde{\mathbf{M}}_{j,1} \in \mathbb{R}^{(n_e-1)(N+1) \times (N+1)}$ and $\widetilde{\mathbf{M}}_{j,2} \in \mathbb{R}^{(n_e-1)(N+1) \times (M-N)}$.

We can rewrite the last system as $(\widetilde{\mathbf{M}}_{j,1}, \widetilde{\mathbf{M}}_{j,2}) \cdot \begin{pmatrix} \widehat{\mathbf{u}}_j \\ \widehat{\mathbf{x}}_j \end{pmatrix} = \widehat{\mathbf{y}}_j$, or

$$(3.7) \quad \widetilde{\mathbf{M}}_{j,2} \cdot \widehat{\mathbf{x}}_j = \widehat{\mathbf{y}}_j - \widetilde{\mathbf{M}}_{j,1} \cdot \widehat{\mathbf{u}}_j.$$

Since the matrix $\widetilde{\mathbf{M}}_{j,2}$ is a submatrix from \mathbf{M}_j and \mathbf{M}_j has, according to Theorem 3.1, the full column rank, $\widetilde{\mathbf{M}}_{j,2}$ has also the full column rank. Thus we conclude the following cases:

(1) If $n_e = (M+1)/(N+1)$, then $(n_e-1)(N+1) = (M-N)$, then the matrix $\widetilde{\mathbf{M}}_{j,2}$ is square and thus it is invertible. Then we get the unique solution

$$(3.8) \quad \widehat{\mathbf{x}}_j = \widetilde{\mathbf{M}}_{j,2}^{-1} \cdot (\widehat{\mathbf{y}}_j - \widetilde{\mathbf{M}}_{j,1} \cdot \widehat{\mathbf{u}}_j).$$

(2) If $n_e > (M+1)/(N+1)$, then $(n_e-1)(N+1) > (M-N)$, then according to the least square solution method¹, see e.g. Strang [11], p. 200, we consider $\widetilde{\mathbf{M}}_{j,2}^\top \cdot \widetilde{\mathbf{M}}_{j,2}$ which is invertible. Moreover, the system (3.7) is over-determined. Now with $\mathbf{A} = \widetilde{\mathbf{M}}_{j,2}$, $\widetilde{\mathbf{x}} = \widehat{\mathbf{x}}_j$, and $\mathbf{b} = \widehat{\mathbf{y}}_j - \widetilde{\mathbf{M}}_{j,1} \cdot \widehat{\mathbf{u}}_j$, the normal equations are

$$(\widetilde{\mathbf{M}}_{j,2}^\top \widetilde{\mathbf{M}}_{j,2}) \widetilde{\mathbf{x}}_j = \widetilde{\mathbf{M}}_{j,2}^\top (\widehat{\mathbf{y}}_j - \widetilde{\mathbf{M}}_{j,1} \cdot \widehat{\mathbf{u}}_j),$$

and the least square solution is given by

$$(3.9) \quad \widehat{\mathbf{x}}_j = (\widetilde{\mathbf{M}}_{j,2}^\top \cdot \widetilde{\mathbf{M}}_{j,2})^{-1} \cdot \widetilde{\mathbf{M}}_{j,2}^\top \cdot (\widehat{\mathbf{y}}_j - \widetilde{\mathbf{M}}_{j,1} \cdot \widehat{\mathbf{u}}_j).$$

(3) If $n_e < (M+1)/(N+1)$, then $(n_e-1)(N+1) < (M-N)$, then the system (3.7) is underdetermined and we ignore this case in our study.²

¹ For an overdetermined problem $\mathbf{A}\mathbf{x} = \mathbf{b}$ with $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$, $m > n$, and $\text{rank } \mathbf{A} = n$, the quadratic minimization problem $\widetilde{\mathbf{x}} = \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$ with the Euclidean norm has a unique solution, provided that the n columns of \mathbf{x} are linearly independent, given by solving the normal equations $(\mathbf{A}^\top \mathbf{A}) \widetilde{\mathbf{x}} = \mathbf{A}^\top \mathbf{b}$. Moreover, the least-squares solution is given by $\widetilde{\mathbf{x}} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$. For details see e.g. Strang [11].

² One could consider the solution of the smallest Euclidean norm which is given by $\widehat{\mathbf{x}}_j = \widetilde{\mathbf{M}}_{j,2}^\top \cdot (\widetilde{\mathbf{M}}_{j,2} \cdot \widetilde{\mathbf{M}}_{j,2}^\top)^{-1} \cdot (\widehat{\mathbf{y}}_j - \widetilde{\mathbf{M}}_{j,1} \cdot \widehat{\mathbf{u}}_j)$. For details see Strang [11], p. 405.

3.4. The solutions as linear combinations. We conclude from the formulas (3.8) and (3.9) that for one fixed element I_j each of the coefficients $\widehat{w}_{i,j}$ for $i = 0, \dots, M$ of the term w_j can be written as a linear combination of all coefficients $\widehat{u}_{k,j+c}$ with $k = 0, \dots, N$ and $c = -L, \dots, R$. This means that there are constants $c_{i,k,j+c} \in \mathbb{R}$ such that

$$\widehat{w}_{i,j} = \sum_{c=-L}^R \sum_{k=0}^N c_{i,k,j+c} \widehat{u}_{k,j+c}.$$

We define the vectors $\mathbf{c}_i \in \mathbb{R}^{n_e(N+1)}$ for $i = 0, \dots, M$, and $\widehat{\mathbf{u}}_{j,s} \in \mathbb{R}^{n_e(N+1)}$ for $j = 1, \dots, Z$, as follows:

$$\begin{aligned} \mathbf{c}_i &:= (c_{i,0,j-L}, c_{i,1,j-L}, \dots, c_{i,N,j-L}, c_{i,0,j-L+1}, \dots, c_{i,N,j-L+1}, \dots, c_{i,N,j+R}), \\ \widehat{\mathbf{u}}_{j,s} &:= (\widehat{u}_{0,j-L}, \widehat{u}_{1,j-L}, \dots, \widehat{u}_{N,j-L}, \widehat{u}_{0,j-L+1}, \dots, \widehat{u}_{N,j-L+1}, \dots, \widehat{u}_{N,j+R}). \end{aligned}$$

The vectors \mathbf{c}_i are identical for different j . They only depend on the form of the stencil. With these vectors we can write $\widehat{w}_{i,j} = \mathbf{c}_i \cdot \widehat{\mathbf{u}}_{j,s}$. By defining the matrix $\mathbf{C} := (\mathbf{c}_0, \dots, \mathbf{c}_M)^\top \in \mathbb{R}^{(M+1) \times n_e(N+1)}$, we can write

$$(3.10) \quad \widehat{\mathbf{w}}_j = \mathbf{C} \cdot \widehat{\mathbf{u}}_{j,s}.$$

In the same way as we studied the matrix $\widetilde{\mathbf{M}}_{j,2}$, we have the following cases:

- (1) If $(n_e - 1)(N + 1) > M - N$, then the matrix $\mathbf{C}^\top \mathbf{C}$ is invertible, it is positive definite.
- (2) If $(n_e - 1)(N + 1) = M - N$, then \mathbf{C} is invertible, and thus the product $\mathbf{C}^\top \mathbf{C}$ is positive definite.

Whether \mathbf{C} is a square matrix or not, the product $\mathbf{C}^\top \mathbf{C}$ will always be positive definite. Then by using the Euclidean norm $\|\cdot\|_e$ and the spectral norm³ $\|\mathbf{C}\|_2 = \sqrt{\beta_{\max}(\mathbf{C}^\top \mathbf{C})}$ we have

$$(3.11) \quad \|\widehat{\mathbf{w}}_j\|_e^2 \leq \|\mathbf{C}\|_2^2 \|\widehat{\mathbf{u}}_{j,s}\|_e^2 = \|\mathbf{C}\|_2^2 (\|\widehat{\mathbf{u}}_{j-L}\|_e^2 + \dots + \|\widehat{\mathbf{u}}_{j+R}\|_e^2).$$

The last inequality is needed for proving the boundedness later.

3.5. Properties of the reconstruction operator. The reconstruction operator is linear. This means that, for all $p, q \in P_{N, I_{\text{ex}}, Z_{\text{ex}}}$ and $\varrho \in \mathbb{R}$, $\mathfrak{R}_{N, M, S, Z}(p + q) =$

³ Let $\mathbf{Q} \in \mathbb{R}^{n \times m}$. The spectral norm of \mathbf{Q} is the largest singular value of \mathbf{Q} , i.e. the square root of the largest eigenvalue $\beta_{\max}(\mathbf{Q}^\top \mathbf{Q})$ of the positive semidefinite matrix $\mathbf{Q}^\top \mathbf{Q}$:

$$\|\mathbf{Q}\|_2 = \sqrt{\beta_{\max}(\mathbf{Q}^\top \mathbf{Q})}.$$

$\mathfrak{R}_{N,M,S,Z}(p) + \mathfrak{R}_{N,M,S,Z}(q)$ and $\mathfrak{R}_{N,M,S,Z}(\varrho p) = \varrho \mathfrak{R}_{N,M,S,Z}(p)$. Moreover, the reconstruction operators have the following properties.

3.5.1. Conservativity. The conservation property $\langle \mathfrak{R}_{N,M,S,Z}(u), u|_I \rangle = \langle u|_I, u|_I \rangle$ holds for all $u \in P_{N,I_{\text{ex}},Z_{\text{ex}}}$ where $\langle \cdot, \cdot \rangle$ is the scalar product on $L^2(I)$.

Proof. Let $w = \mathfrak{R}_{N,M,S,Z}(u)$ be the reconstructed polynomial. We use the equalities (3.4) which hold for the terms w_j and u_j of w and u , respectively, on the elements I_j . We have by using the orthogonality of the basis functions $\langle w, u|_I \rangle = \sum_{j=1}^Z \langle w_j, u_j \rangle_j = \sum_{j=1}^Z \langle u_j, u_j \rangle_j = \langle u|_I, u|_I \rangle$, where $\langle \cdot, \cdot \rangle_j$ is the L^2 scalar product on the element I_j . \square

3.5.2. Consistency of the reconstruction for $p \in P_M$ and an identity property. We prove that in the special case when the given function v is a polynomial $p \in P_M = P_{M,\mathbb{R}}$ of degree M , the system (3.6), $\mathbf{M}_j \cdot \widehat{\mathbf{w}}_j = \mathbf{y}_j$ has a consistent right-hand side.

Theorem 3.2. *Let $p \in P_M$, $u = \Pi_{N,Z_{\text{ex}}}(p)$ and $w = \mathfrak{R}_{M,M,S,Z}(u)$. The system (3.6) has a consistent right-hand side. This means that for $p \in P_M$ the least squares solution is an exact solution of the linear system (3.6) in the overdetermined case.*

Proof. Since $p \in P_M$ and the extended polynomials $\Phi_{0,j}^e, \dots, \Phi_{M,j}^e$ are a basis of P_M , there exist constants $\widehat{p}_{0,j}, \dots, \widehat{p}_{M,j} \in \mathbb{R}$ such that $p(x) = \sum_{m=0}^M \widehat{p}_{m,j} \Phi_{m,j}^e(x)$. For $k = 0, \dots, N$ and $c = -L, \dots, R$ let us now consider the entries of the system $\mathbf{A}_{j+c} \cdot \widehat{\mathbf{u}}_{j+c} = \mathbf{M}_{j,c} \cdot \widehat{\mathbf{w}}_j$ defined in (3.5). Using (2.3) as in Section 2 we obtain for a row of the system

$$\begin{aligned} \frac{h}{2k+1} \widehat{u}_{k,j+c} &= \int_{I_{j+c}} \Phi_{k,j+c}(x) p(x) \, dx = \sum_{m=0}^M \widehat{p}_{m,j} \int_{I_{j+c}} \Phi_{k,j+c}(x) \Phi_{m,j}^e(x) \, dx \\ &= \sum_{m=0}^M \widehat{p}_{m,j} \langle \Phi_{k,j+c}, \Phi_{m,j}^e \rangle_{j+c}. \end{aligned}$$

This is the k th row of the equation $\mathbf{y}_{j,c} = \mathbf{A}_{j+c} \cdot \widehat{\mathbf{u}}_{j+c} = \mathbf{M}_{j,c} \cdot \widehat{\mathbf{p}}_j$ for $\widehat{\mathbf{p}}_j := (\widehat{p}_{0,j}, \dots, \widehat{p}_{M,j})^\top$. This means that $\mathbf{y}_{j,c}$ is in the range of $\mathbf{M}_{j,c}$ for all $c = -L, \dots, R$. This holds also for the system (3.6), i.e., the vector \mathbf{y}_j is in the range of \mathbf{M}_j . Thus the system (3.6) has a consistent right-hand side. \square

Lemma 3.2. Let $p \in P_M$. For all $N \in \{0, \dots, M\}$, by using any stencil $S = S_{I_j, n_e, L}$ with size n_e satisfying $n_e \geq (M+1)/(N+1)$, we have

$$\mathfrak{R}_{N, M, S, Z}(II_{N, Z_{\text{ex}}}(p)) = p|_I.$$

So $\mathfrak{R}_{N, M, S, Z} \circ II_{N, Z_{\text{ex}}}$ is a quasi identity map for polynomials in P_M . It uses the known values of $p \in P_M$ on the extended interval I_{ex} .

Proof. Let $u = II_{N, Z_{\text{ex}}}(p) \in P_{N, I_{\text{ex}}, Z_{\text{ex}}}$ and $w = \mathfrak{R}_{M, M, S, Z}(u) \in P_{M, I, Z}$. Since P_M may be identified one to one with $P_{M, I_{\text{ex}}}$ the polynomial p can be expanded to I_{ex} piecewise. Then we have

$$p = \sum_{j=1}^{Z_{\text{ex}}} \sum_{i=0}^M \widehat{p}_{i,j} \Phi_{i,j}, \quad u = \sum_{j=1}^{Z_{\text{ex}}} \sum_{k=0}^N \widehat{u}_{k,j} \Phi_{k,j}, \quad w = \sum_{j=1}^Z \sum_{l=0}^M \widehat{w}_{l,j} \Phi_{l,j},$$

where all $\widehat{u}_{k,j}$ and $\widehat{p}_{i,j}$ are given by (2.3) and the $\widehat{w}_{l,j}$ are the solutions of the reconstruction equations. We have on the extended interval $\widehat{p}_{i,j} = \widehat{u}_{i,j}$ for $i = 0, \dots, N$ and $j = 1, \dots, Z_{\text{ex}}$ and on the original interval the reconstruction coefficients satisfy $\widehat{w}_{l,j} = \widehat{p}_{l,j} = \widehat{u}_{l,j}$ for $l = 0, \dots, N$ and $j = 1, \dots, Z$. We want to prove that $\widehat{w}_{l,j} = \widehat{p}_{l,j}$ for $l = N+1, \dots, M$ and $j = 1, \dots, Z$. By Theorem 3.1 the solution to the reconstruction system (3.6) is unique since \mathbf{M}_j has full column rank. This means that for $p \in P_M$ we have $\widehat{p}_{l,j} = \widehat{w}_{l,j}$ for $l = N+1, \dots, M$ and $j = 1, \dots, Z$. \square

3.5.3. Further relations between $\mathfrak{R}_{N, M, S, Z}$ and $II_{N, Z_{\text{ex}}}$.

Theorem 3.3. Let $p, q \in P_{N, I_{\text{ex}}, Z_{\text{ex}}}$ be such that $\mathfrak{R}_{N, M, S, Z}(p) = \mathfrak{R}_{N, M, S, Z}(q)$. Then we have $p|_I = q|_I$.

Proof. Let $P = \mathfrak{R}_{N, M, S, Z}(p) = \mathfrak{R}_{N, M, S, Z}(q) = Q$. Then we may write

$$\sum_{j=1}^Z \sum_{i=0}^M \widehat{P}_{i,j} \Phi_{i,j}(x) = \sum_{j=1}^Z \sum_{i=0}^M \widehat{Q}_{i,j} \Phi_{i,j}(x),$$

or $\sum_{j=1}^Z \sum_{i=0}^M (\widehat{P}_{i,j} - \widehat{Q}_{i,j}) \Phi_{i,j}(x) = 0$. Since the piecewise polynomials $\Phi_{i,j}$ are linearly independent basis functions in $P_{M, I, Z}$, we have $\widehat{P}_{i,j} = \widehat{Q}_{i,j}$ for all $i = 0, \dots, M$ and $j = 1, \dots, Z$. On the other hand, the equalities (3.4) give $\widehat{p}_{i,j} = \widehat{P}_{i,j}$ and $\widehat{q}_{i,j} = \widehat{Q}_{i,j}$ for all $i = 0, \dots, N$ and $j = 1, \dots, Z$. Thus we find $\widehat{p}_{i,j} = \widehat{q}_{i,j}$ for $i = 0, \dots, N$ and $j = 1, \dots, Z$. Then $p|_I = q|_I$. \square

Theorem 3.4. For any stencil $S_{I_j, n_e, L}$ with $n_e \geq (M+1)/(N+1)$, we have for any $u \in P_{N, I_{\text{ex}}, Z_{\text{ex}}}$

$$\Pi_{N, Z}(\mathfrak{R}_{N, M, S, Z}(u)) = u|_I.$$

Proof. Let $w = \mathfrak{R}_{N, M, S, Z}(u)$, and $q = \Pi_{N, Z}(w)$. We want to prove that $q = u|_I$. We have

$$u = \sum_{j=1}^{Z_{\text{ex}}} \sum_{i=0}^N \hat{u}_{i,j} \Phi_{i,j}, \quad w = \sum_{j=1}^Z \sum_{i=0}^M \hat{w}_{i,j} \Phi_{i,j},$$

and

$$q = \sum_{j=1}^Z \sum_{i=0}^N \hat{q}_{i,j} \Phi_{i,j}.$$

For $i = 0, \dots, N$ and $j = 1, \dots, Z$ according to (2.3), we have

$$\hat{q}_{i,j} = \frac{2i+1}{h} \langle w, \Phi_{i,j} \rangle_j = \frac{2i+1}{h} \sum_{k=1}^Z \sum_{l=0}^M \hat{w}_{l,k} \langle \Phi_{l,k}, \Phi_{i,j} \rangle_j.$$

Since the basis functions satisfy

$$\int_{I_j} \Phi_{m,j}(x) \Phi_{n,j'}(x) dx = \begin{cases} \frac{h}{2m+1} \delta_{mn}, & j = j', \\ 0, & j \neq j', \end{cases} \quad m, n = 0, \dots, N,$$

we have

$$\hat{q}_{i,j} = \frac{2i+1}{h} \left(\hat{w}_{i,j} \frac{h}{2i+1} \right) = \hat{w}_{i,j}.$$

According to (3.4), we have $\hat{w}_{i,j} = \hat{u}_{i,j}$, then $\hat{q}_{i,j} = \hat{u}_{i,j}$ for $i = 0, \dots, N$ and $j = 1, \dots, Z$. This implies finally that $q = u|_I$. \square

Theorem 3.5. For all $w \in P_{M, I_{\text{ex}}, Z_{\text{ex}}}$ and all $N \in \{0, \dots, M\}$, the relation $\mathfrak{R}_{N, M, S, Z}(\Pi_{N, Z_{\text{ex}}}(w)) = w|_I$ holds.

Proof. Follows directly from Lemma 3.2. \square

3.5.4. Boundedness of the reconstruction operator.

Theorem 3.6. Let $w = \mathfrak{R}_{N, M, S, Z}(u)$ and let w_j and u_j be the restrictions to I_j of w and u , respectively. Then the following inequalities hold:

$$(3.12) \quad \|u_j\|_{L^2(I_j)}^2 \leq \|w_j\|_{L^2(I_j)}^2 \leq C_6 \sum_{c=-L}^R \|u_{j+c}\|_{L^2(I_{j+c})}^2.$$

Proof. (1) We start with the first part of the inequality (3.12). From the conservation property 3.5.1 and from the Cauchy-Schwarz inequality, we have

$$\|u_j\|_{L^2(I_j)}^2 = \langle u_j, u_j \rangle_j = \langle w_j, u_j \rangle_j \leq \|u_j\|_{L^2(I_j)} \|w_j\|_{L^2(I_j)}.$$

If $\|u_j\|_{L^2(I_j)} = 0$ then trivially $\|u_j\|_{L^2(I_j)}^2 \leq \|w_j\|_{L^2(I_j)}^2$. If $\|u_j\|_{L^2(I_j)} > 0$, then after dividing by $\|u_j\|_{L^2(I_j)}$ we get $\|u_j\|_{L^2(I_j)} \leq \|w_j\|_{L^2(I_j)}$ and hence $\|u_j\|_{L^2(I_j)}^2 \leq \|w_j\|_{L^2(I_j)}^2$, thus the first part of (3.12) follows.

(2) By taking $\widehat{\mathbf{w}}_j = (\widehat{w}_{0,j}, \dots, \widehat{w}_{M,j})^\top$ and using the Euclidean norm, we have

$$(3.13) \quad \frac{h}{2M+1} \|\widehat{\mathbf{w}}_j\|_e^2 \leq \|w_j\|_{L^2(I_j)}^2 \leq h \|\widehat{\mathbf{w}}_j\|_e^2.$$

On the other hand, according to the formula (3.10), we have $\widehat{\mathbf{w}}_j = \mathbf{C} \widehat{\mathbf{u}}_{j,s}$. From (3.11) we get

$$\|\widehat{\mathbf{w}}_j\|_e^2 \leq \|\mathbf{C}\|_2^2 \|\widehat{\mathbf{u}}_{j,s}\|_e^2 = \|\mathbf{C}\|_2^2 (\|\widehat{\mathbf{u}}_{j-L}\|_e^2 + \dots + \|\widehat{\mathbf{u}}_{j+R}\|_e^2).$$

From (2.4) for $c = -L, \dots, R$, $\|\widehat{\mathbf{u}}_{j+c}\|_e^2 \leq (2N+1)h^{-1} \|u_{j+c}\|_{L^2(I_{j+c})}^2$ we have

$$\|\widehat{\mathbf{w}}_j\|_e^2 \leq \frac{(2N+1)\|\mathbf{C}\|_2^2}{h} (\|u_{j-L}\|_{L^2(I_{j-L})}^2 + \dots + \|u_{j+R}\|_{L^2(I_{j+R})}^2).$$

Substituting into (3.13), we obtain

$$\|w_j\|_{L^2(I_j)}^2 \leq (2N+1)\|\mathbf{C}\|_2^2 (\|u_{j-L}\|_{L^2(I_{j-L})}^2 + \dots + \|u_{j+R}\|_{L^2(I_{j+R})}^2).$$

Finally, we get $C_6 = (2N+1)\|\mathbf{C}\|_2^2$, noting that the coefficients in \mathbf{C} only depend on the basis polynomials and not on v_j or w_j . Now, we have

$$\|w_j\|_{L^2(I_j)}^2 \leq C_6 \sum_{c=-L}^R \|u_{j+c}\|_{L^2(I_{j+c})}^2,$$

which is the right inequality. \square

3.6. The error estimate of the reconstruction operator.

Theorem 3.7. *Suppose that the interval $I = [a, b]$ has a uniform partition of Z subintervals with constant mesh size $h = (b-a)/Z$. Let $N \leq M$, $S = S_{I_j, n_e, L}$ be a stencil with $n_e \geq (M+1)/(N+1)$, and $I_{\text{ex}} = \bigcup_{j=1}^Z S_{I_j, n_e, L}$ be the extended interval. Then, for each $v \in W^{M+1,2}(I_{\text{ex}})$, the following error estimates hold:*

$$(3.14) \quad \|\mathfrak{R}_{N,M,S,Z}(II_{N,Z_{\text{ex}}}(v)) - v\|_{L^2(I)} \leq C_7 h^{M+1} |v|_{W^{M+1,2}(I)},$$

where $|\cdot|_{W^{M+1,2}(I)}$ is the seminorm on $W^{M+1,2}(I)$ given in (2.5).

Proof. Let I_j be an element, with $j = 1, \dots, Z$ fixed. Using the triangle inequality, we obtain

$$(3.15) \quad \begin{aligned} & \|\mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(v)) - v|_I\|_{L^2(I_j)}^2 \\ & \leq \|\mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(v)) - \Pi_{M,Z_{\text{ex}}}(v)|_I\|_{L^2(I_j)}^2 \\ & \quad + \|\Pi_{M,Z_{\text{ex}}}(v)|_I - v|_I\|_{L^2(I_j)}^2. \end{aligned}$$

Due to the identity in Lemma 3.2 and the linearity of the reconstruction operator we have

$$\begin{aligned} & \|\mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(v)) - \Pi_{M,Z_{\text{ex}}}(v)|_I\|_{L^2(I_j)}^2 \\ & = \|\mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(v)) - \mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(\Pi_{M,Z_{\text{ex}}}(v)))\|_{L^2(I_j)}^2 \\ & = \|\mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(v) - \Pi_{N,Z_{\text{ex}}}(\Pi_{M,Z_{\text{ex}}}(v)))\|_{L^2(I_j)}^2. \end{aligned}$$

By virtue of inequality (3.12) and the linearity as well as boundedness of the projection operator we obtain

$$\begin{aligned} & \|\mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(v)) - \Pi_{M,Z_{\text{ex}}}(v)|_I\|_{L^2(I_j)}^2 \\ & \leq C_6 \sum_{c=-L}^R \|(\Pi_{N,Z_{\text{ex}}}(v) - \Pi_{N,Z_{\text{ex}}}(\Pi_{M,Z_{\text{ex}}}(v)))|_{I_{j+c}}\|_{L^2(I_{j+c})}^2 \\ & = C_6 \sum_{c=-L}^R \|(\Pi_{N,Z_{\text{ex}}}(v - \Pi_{M,Z_{\text{ex}}}(v)))|_{I_{j+c}}\|_{L^2(I_{j+c})}^2 \\ & = C_6 \sum_{c=-L}^R \|(v - \Pi_{M,Z_{\text{ex}}}(v))|_{I_{j+c}}\|_{L^2(I_{j+c})}^2. \end{aligned}$$

Substituting into (3.15) we get

$$\begin{aligned} & \|\mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(v)) - v|_I\|_{L^2(I_j)}^2 \\ & \leq C_6 \sum_{c=-L}^R \|(v - \Pi_{M,Z_{\text{ex}}}(v))|_{I_{j+c}}\|_{L^2(I_{j+c})}^2 + \|\Pi_{M,Z_{\text{ex}}}(v)|_I - v|_I\|_{L^2(I_j)}^2. \end{aligned}$$

Now the error estimate (2.6) gives

$$\begin{aligned} & \|\mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(v)) - v|_I\|_{L^2(I_j)}^2 \\ & \leq C_6 \sum_{c=-L}^R (C_2^2 h^{2M+2} |v|_{W^{M+1,2}(I_{j+c})}^2) + C_2^2 h^{2M+2} |v|_{W^{M+1,2}(I_j)}^2 \\ & = (1 + n_e C_6) C_2^2 h^{2M+2} |v|_{W^{M+1,2}(I_j)}^2. \end{aligned}$$

By taking $C_7 = C_2\sqrt{1 + n_e C_6}$ and by summing over all j , we get

$$\|\mathfrak{R}_{N,M,S,Z}(\Pi_{N,Z_{\text{ex}}}(v)) - v|_I\|_{L^2(I)}^2 \leq C_7^2 h^{2M+2} |v|_{W^{M+1,2}(I)}^2.$$

Finally, we take the square root and get the inequality (3.14). \square

4. THE LOCAL SPACE TIME GALERKIN SCHEME

This scheme evolves the reconstructed polynomials locally in time inside each element to the same order of accuracy as in space, by using the governing equations. We follow the procedure of Dumbser et al. [5] for the $P_N P_M$ DG schemes.

Recall that for $Z_1 \in \mathbb{N}$ we discretize the time interval $[0, T]$ by the times $0 = t_0 < t_1 < \dots < t_{Z_1} = T$, subintervals $T_n = [t_n, t_{n+1}[$ and a constant time step $\Delta t = t_{n+1} - t_n$ for $n = 0, \dots, Z_1 - 1$. Further, the space of polynomials of degree M in space and time on $T_n \times I_j$ is denoted as $P_{M, T_n \times I_j}$. We use basis functions $\theta_{i,j} \in P_{M, T_n \times I_j}$. The functions $\theta_{i,j}$ are nodal functions. This means that we choose some nodes on $T_n \times I_j$. Then we relate each node to a function such that this function equals to 1 at this node and equals to 0 at the others. The number of nodes should equal the number of degrees of freedom of these polynomials. These functions take their value to be zero outside of $T_n \times I_j$. For an arbitrary degree M , the number of functions $\theta_{i,j}$ is denoted by \mathcal{N} . It is given by $\mathcal{N} = (M + 1)(M + 2)/2$. Then we set the basis to be $\Theta_{M,j} = \{\theta_{1,j}, \dots, \theta_{\mathcal{N},j}\}$. The first $M + 1$ functions are taken to depend on the spatial variable x at $t = t_n$. They could be grouped together in a subbasis $\Theta_{M,j}^0$. All other basis functions vanish for $t = t_n$ and could be grouped together in a subbasis $\Theta_{M,j}^1$. This means $\Theta_{M,j} = \Theta_{M,j}^0 \cup \Theta_{M,j}^1$. For an example of these basis functions see Appendix B.

Let $v \in L^\infty([0, T], L^2(I))$. The general form of a conservation law is given by

$$(4.1) \quad v_t(t, x) + f(v(t, x))_x = 0 \quad \text{for } x \in I, t \in [0, T].$$

We multiply the equation by $\theta_{k,j}$ for $k = 1, \dots, \mathcal{N}$. Integrating over $T_n \times I_j$, we get

$$(4.2) \quad \int_{T_n} \int_{I_j} \theta_{k,j}(t, x) \frac{\partial}{\partial t} v(t, x) \, dx \, dt + \int_{T_n} \int_{I_j} \theta_{k,j}(t, x) \frac{\partial}{\partial x} f(v(t, x)) \, dx \, dt = 0.$$

We suppose that the solution v and the flux $f(v)$ are approximated on $T_n \times I$ by the formulas

$$(4.3) \quad v(t, x) := U^n(t, x) = \sum_{j=1}^Z \sum_{i=1}^{\mathcal{N}} \widehat{U}_{i,j}^n \theta_{i,j}(t, x),$$

$$f(v(t, x)) := F^n(t, x) = \sum_{j=1}^Z \sum_{i=1}^{\mathcal{N}} f(\widehat{U}_{i,j}^n) \theta_{i,j}(t, x),$$

for $(t, x) \in T_n \times I$. Now we introduce the scalar product

$$(4.4) \quad \langle g, h \rangle_{tx} = \int_{T_n} \int_{I_j} g(t, x) h(t, x) \, dx \, dt,$$

and use these approximations in (4.2) to write

$$\sum_{i=1}^{\mathcal{N}} \left\langle \theta_{k,j}, \frac{\partial}{\partial t} \theta_{i,j} \right\rangle_{tx} \widehat{U}_{i,j}^n + \sum_{i=1}^{\mathcal{N}} \left\langle \theta_{k,j}, \frac{\partial}{\partial x} \theta_{i,j} \right\rangle_{tx} f(\widehat{U}_{i,j}^n) = 0.$$

We introduce the matrix entries $G_{ki} := \langle \theta_{k,j}, \frac{\partial}{\partial t} \theta_{i,j} \rangle_{tx}$ and $H_{ki} := \langle \theta_{k,j}, \frac{\partial}{\partial x} \theta_{i,j} \rangle_{tx}$ for $1 \leq i, k \leq \mathcal{N}$. The values of these entries do not depend on j , since we use the same shifted basis for each j via a reference transformation. Introducing the vectors $\widehat{\mathbf{U}}_j^n := (\widehat{U}_{1,j}^n, \dots, \widehat{U}_{\mathcal{N},j}^n)^\top$ and $\widehat{\mathbf{F}}_j^n := (f(\widehat{U}_{1,j}^n), \dots, f(\widehat{U}_{\mathcal{N},j}^n))^\top$, we get the matrix form $\mathbf{G} \widehat{\mathbf{U}}_j^n + \mathbf{H} \widehat{\mathbf{F}}_j^n = 0$. The first $M+1$ degrees of freedom are related to the functions $\Theta_{M,j}^0$. We group them together into the subvector $\widehat{\mathbf{U}}_j^{n,0} \in \mathbb{R}^{M+1}$. All other degrees of freedom are grouped together into the subvector $\widehat{\mathbf{U}}_j^{n,1} \in \mathbb{R}^{\mathcal{N}-M-1}$. Analogously, we define the subvectors $\widehat{\mathbf{F}}_j^{n,0}$ and $\widehat{\mathbf{F}}_j^{n,1}$. Then the matrix form becomes $\mathbf{G} \begin{pmatrix} \widehat{\mathbf{U}}_j^{n,0} \\ \widehat{\mathbf{U}}_j^{n,1} \end{pmatrix} + \mathbf{H} \begin{pmatrix} \widehat{\mathbf{F}}_j^{n,0} \\ \widehat{\mathbf{F}}_j^{n,1} \end{pmatrix} = 0$. We write the matrices \mathbf{G} and \mathbf{H} as block matrices $\mathbf{G} = \begin{pmatrix} \mathbf{G}^{00} & \mathbf{G}^{01} \\ \mathbf{G}^{10} & \mathbf{G}^{11} \end{pmatrix}$ and $\mathbf{H} = \begin{pmatrix} \mathbf{H}^{00} & \mathbf{H}^{01} \\ \mathbf{H}^{10} & \mathbf{H}^{11} \end{pmatrix}$, respectively, where $\mathbf{G}^{00}, \mathbf{H}^{00} \in \mathbb{R}^{(M+1) \times (M+1)}$, $\mathbf{G}^{01}, \mathbf{H}^{01} \in \mathbb{R}^{(M+1) \times (\mathcal{N}-M-1)}$, $\mathbf{G}^{10}, \mathbf{H}^{10} \in \mathbb{R}^{(\mathcal{N}-M-1) \times (M+1)}$, and $\mathbf{G}^{11}, \mathbf{H}^{11} \in \mathbb{R}^{(\mathcal{N}-M-1) \times (\mathcal{N}-M-1)}$. Then we get

$$(4.5) \quad \begin{pmatrix} \mathbf{G}^{00} & \mathbf{G}^{01} \\ \mathbf{G}^{10} & \mathbf{G}^{11} \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{U}}_j^{n,0} \\ \widehat{\mathbf{U}}_j^{n,1} \end{pmatrix} + \begin{pmatrix} \mathbf{H}^{00} & \mathbf{H}^{01} \\ \mathbf{H}^{10} & \mathbf{H}^{11} \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{F}}_j^{n,0} \\ \widehat{\mathbf{F}}_j^{n,1} \end{pmatrix} = 0.$$

We determine the components of the vector $\widehat{\mathbf{U}}_j^{n,0}$ by projecting the reconstructed polynomial w^n at time $t = t_n$ onto the space spanned by the first nodal functions $\Theta_{M,j}^0$. This gives for $k = 1, \dots, M+1$ the system of equations

$$\int_{I_j} \theta_{k,j}(t_n, x) \sum_{i=1}^{\mathcal{N}} \widehat{U}_{i,j}^n \theta_{i,j}(t_n, x) \, dx = \int_{I_j} \theta_{k,j}(t_n, x) \sum_{l=0}^M \widehat{w}_{l,j}^n \Phi_{l,j}(x) \, dx.$$

Since $\theta_{i,j}(t_n, x) = 0$ for $i = M+2, \dots, \mathcal{N}$, the sum on the left-hand side reduces to the first $M+1$ terms. Thus we get

$$\sum_{i=1}^{M+1} \langle \theta_{k,j}(t_n, \cdot), \theta_{i,j}(t_n, \cdot) \rangle_j \widehat{U}_{i,j}^n = \sum_{l=0}^M \langle \theta_{k,j}(t_n, \cdot), \Phi_{l,j}(\cdot) \rangle_j \widehat{w}_{l,j}^n.$$

Again, this system can be written in a matrix-vector form. Since the functions $\theta_{i,j}$ for $i = 1, \dots, M + 1$ belong to the basis, they are linearly independent. This implies the linear independence of the columns of the scalar product matrix. Thus we have a unique solution and the first $M + 1$ coefficients $\widehat{U}_{i,j}^n$ are known from the reconstructed coefficients $\widehat{w}_{i,j}^n$.

Since we now have determined $M + 1$ known degrees of freedom, we no longer need the upper blocks in (4.5). Therefore, we cancel the first $M + 1$ rows of this system. We obtain the smaller system

$$(\mathbf{G}^{10} \mathbf{G}^{11}) \begin{pmatrix} \widehat{\mathbf{U}}_j^{n,0} \\ \widehat{\mathbf{U}}_j^{n,1} \end{pmatrix} + (\mathbf{H}^{10} \mathbf{H}^{11}) \begin{pmatrix} \widehat{\mathbf{F}}_j^{n,0} \\ \widehat{\mathbf{F}}_j^{n,1} \end{pmatrix} = 0.$$

In order to determine the vector $\widehat{\mathbf{U}}_j^{n,1}$ we have to solve the nonlinear equations

$$\mathbf{G}^{11} \widehat{\mathbf{U}}_j^{n,1} + \mathbf{H}^{11} \widehat{\mathbf{F}}_j^{n,1} = -\mathbf{H}^{10} \widehat{\mathbf{F}}_j^{n,0} - \mathbf{G}^{10} \widehat{\mathbf{U}}_j^{n,0}.$$

The quadratic matrix \mathbf{G}^{11} depends on the mesh size h but not on the time step or the equations to be solved. For all orders of accuracy, the matrix \mathbf{G}^{11} is invertible, since its columns are linearly independent. Therefore, after inverting, we obtain a fixed-point problem for the unknowns $\widehat{\mathbf{U}}_j^{n,1}$:

$$\widehat{\mathbf{U}}_j^{n,1} = (\mathbf{G}^{11})^{-1} [-\mathbf{H}^{11} \widehat{\mathbf{F}}_j^{n,1} - \mathbf{H}^{10} \widehat{\mathbf{F}}_j^{n,0} - \mathbf{G}^{10} \widehat{\mathbf{U}}_j^{n,0}].$$

We solve this system using the fixed-point iteration

$$\widehat{\mathbf{U}}_j^{n,1,i+1} = (\mathbf{G}^{11})^{-1} [-\mathbf{H}^{11} \widehat{\mathbf{F}}_j^{n,1,i} - \mathbf{H}^{10} \widehat{\mathbf{F}}_j^{n,0} - \mathbf{G}^{10} \widehat{\mathbf{U}}_j^{n,0}].$$

The superscript i denotes the iteration number. This approach works, since $(\mathbf{G}^{11})^{-1} \mathbf{H}^{11}$ turns out to be a contraction mapping, see Dumbser et al. [5], p. 8218. In our practical computations the fixed-point was determined after at most $M + 1$ iterations.

As suggested in [5] we begin iterating by using a stationary in time solution of (4.1) as an initial guess value for $\widehat{\mathbf{U}}_j^{n,1}$. The stationary equation is $v_t = 0$. The matrix form is $\mathbf{G} \widehat{\mathbf{U}}_j^n = 0$. Then we get the initial guess with $i = 0$, $\widehat{\mathbf{U}}_j^{n,1,0} = -(\mathbf{G}^{11})^{-1} \mathbf{G}^{10} \widehat{\mathbf{U}}_j^{n,0}$.

5. THE DG SCHEMES

We recall the basic steps leading to the DG schemes. Now we apply the DG schemes [3], [4] and use numerical fluxes whose arguments are the solutions U^n and F^n of the previous step.

Again we consider a conservation law in one space dimension $v_t + f_x(v) = 0$ for $(t, x) \in [0, T] \times I$. Let I_j be a space element with $j = 1, \dots, Z$ fixed and $T_n = [t_n, t_{n+1}[$. Now we multiply by an arbitrary smooth function $\chi \in L^2(I_j)$, integrate over $T_n \times I_j$, and use integration by parts in space to get

$$\begin{aligned} \int_{T_n} \int_{I_j} \chi(x) \frac{\partial}{\partial t} v(t, x) \, dx \, dt + \int_{T_n} \chi(x) f(v(t, x)) \Big|_{x_{j-1/2}}^{x_{j+1/2}} \, dt \\ - \int_{T_n} \int_{I_j} \frac{\partial \chi(x)}{\partial x} f(v(t, x)) \, dx \, dt = 0. \end{aligned}$$

We assume that the numerical solution u^n is a piecewise polynomial of degree N and is defined on $T_n \times I$ as

$$u^n(t, x) = \sum_{j=1}^Z \sum_{l=0}^N \hat{u}_{l,j}^n(t) \Phi_{l,j}(x),$$

where $\hat{u}_{l,j}^n \in \mathbb{R}$ are the unknowns and $\Phi_{l,j}$ are the Legendre basis functions (2.2). Substituting the numerical solution u^n for v and replacing the test function χ by the basis functions $\Phi_{k,j}$ for $k = 0, \dots, N$ leads to

$$\begin{aligned} \int_{T_n} \int_{I_j} \Phi_{k,j}(x) \frac{\partial}{\partial t} u^n(t, x) \, dx \, dt + \int_{T_n} \Phi_{k,j}(x) f(u^n(t, x)) \Big|_{x_{j-1/2}}^{x_{j+1/2}} \, dt \\ - \int_{T_n} \int_{I_j} \frac{\partial \Phi_{k,j}(x)}{\partial x} f(u^n(t, x)) \, dx \, dt = 0. \end{aligned}$$

The sum over the index j in the solution u^n is reduced only to one term u_j^n which has values on I_j and the other terms have value zero on I_j . By the index $k = 0, \dots, N$ we have $N + 1$ equations for the $N + 1$ unknown coefficients $\hat{u}_{k,j}^n(t)$. Now with

$$u_j^n(t, x) = \sum_{l=0}^N \hat{u}_{l,j}^n(t) \Phi_{l,j}(x) \text{ we have}$$

$$\begin{aligned} (5.1) \quad \int_{T_n} \int_{I_j} \Phi_{k,j}(x) \frac{\partial}{\partial t} \left(\sum_{l=0}^N \hat{u}_{l,j}^n(t) \Phi_{l,j}(x) \right) \, dx \, dt \\ + \int_{T_n} \Phi_{k,j}(x) f(u_j^n(t, x)) \Big|_{x_{j-1/2}}^{x_{j+1/2}} \, dt - \int_{T_n} \int_{I_j} \frac{\partial \Phi_{k,j}(x)}{\partial x} f(u_j^n(t, x)) \, dx \, dt = 0. \end{aligned}$$

The time derivative and the time integral in the first term are applied to $\widehat{u}_{l,j}^n$, since the functions $\Phi_{l,j}$ are independent of the time variable. This derivative can be replaced for $x \in I_j$ as follows:

$$\begin{aligned} \int_{T_n} \frac{\partial}{\partial t} \left(\sum_{l=0}^N \widehat{u}_{l,j}^n(t) \Phi_{l,j}(x) \right) dt &= \sum_{l=0}^N \left(\int_{T_n} \frac{d}{dt} \widehat{u}_{l,j}^n(t) dt \right) \Phi_{l,j}(x) \\ &= \sum_{l=0}^N (\widehat{u}_{l,j}^{n+1} - \widehat{u}_{l,j}^n) \Phi_{l,j}(x). \end{aligned}$$

Thus for the first term in (5.1) we have $\int_{I_j} \Phi_{k,j}(x) \sum_{l=0}^N (\widehat{u}_{l,j}^{n+1} - \widehat{u}_{l,j}^n) \Phi_{l,j}(x) dx$. Due to the orthogonality of the basis functions, this sum reduces to one term with index $k = l$ and the space integral has the value $h/(2k + 1)$. Thus the first term becomes $(\widehat{u}_{k,j}^{n+1} - \widehat{u}_{k,j}^n)h/(2k + 1)$.

According to (2.1) we have $\Phi_{k,j}(x_{j-1/2}) = (-1)^k$ and $\Phi_{k,j}(x_{j+1/2}) = 1$ for $k = 0, \dots, N$. Then for the second term in (5.1) we have

$$\int_{T_n} [f(u_j^n(t, x_{j+1/2})) - (-1)^k f(u_j^n(t, x_{j-1/2}))] dt.$$

Substituting these first and second terms in (5.1), we obtain the system

$$(5.2) \quad \begin{aligned} \frac{h}{2k+1} (\widehat{u}_{k,j}^{n+1} - \widehat{u}_{k,j}^n) + \int_{T_n} [f(u_j^n(t, x_{j+1/2})) - (-1)^k f(u_j^n(t, x_{j-1/2}))] dt \\ - \int_{T_n} \int_{I_j} \frac{\partial \Phi_{k,j}(x)}{\partial x} f(u_j^n(t, x)) dx dt = 0. \end{aligned}$$

For the values $f(u_j^n(t, x_{j+1/2}))$, we follow Dumbser et al. [5] and use numerical flux functions whose arguments are the solutions U_j^n , see (4.3), of the local Galerkin scheme, i.e. $f(u_j^n(t, x_{j+1/2})) \approx \mathcal{F}_{j+1/2}^n(t) := \mathcal{F}(U_j^n(t, x_{j+1/2}), U_{j+1}^n(t, x_{j+1/2}))$. For \mathcal{F} we may use any consistent numerical flux function, see the next section. Similarly $f(u_j^n(t, x_{j-1/2})) \approx \mathcal{F}_{j-1/2}^n(t) := \mathcal{F}(U_{j-1}^n(t, x_{j-1/2}), U_j^n(t, x_{j-1/2}))$. We insert also these solutions U_j^n in the formulas of the flux $f(u_j^n(t, x)) \approx \sum_{i=1}^{\mathcal{N}} f(\widehat{U}_{i,j}^n) \theta_{i,j}(t, x)$ for $(t, x) \in T_n \times I_j$.

Substituting into (5.2), we get using (4.3)

$$\begin{aligned} \frac{h}{2k+1} (\widehat{u}_{k,j}^{n+1} - \widehat{u}_{k,j}^n) + \int_{T_n} \mathcal{F}_{j+1/2}^n(t) dt - (-1)^k \int_{T_n} \mathcal{F}_{j-1/2}^n(t) dt \\ - \int_{T_n} \int_{I_j} \frac{\partial \Phi_{k,j}(x)}{\partial x} \left[\sum_{i=1}^{\mathcal{N}} f(\widehat{U}_{i,j}^n) \theta_{i,j}(t, x) \right] dx dt = 0. \end{aligned}$$

The coefficients $f(\widehat{U}_{i,j}^n)$ are constants, thus, using the scalar product $\langle \cdot, \cdot \rangle_{tx}$, given by (4.4), we get

$$\begin{aligned} \frac{h}{2k+1}(\widehat{u}_{k,j}^{n+1} - \widehat{u}_{k,j}^n) + \int_{T_n} \mathcal{F}_{j+1/2}^n(t) dt - (-1)^k \int_{T_n} \mathcal{F}_{j-1/2}^n(t) dt \\ - \sum_{i=1}^{\mathcal{N}} f(\widehat{U}_{i,j}^n) \left\langle \frac{\partial \Phi_{k,j}}{\partial x}, \theta_{i,j} \right\rangle_{tx} = 0. \end{aligned}$$

Finally, by rearranging the terms, we get the fully discrete one-step $P_N P_M$ DG scheme

$$(5.3) \quad \widehat{u}_{k,j}^{n+1} = \widehat{u}_{k,j}^n - \frac{2k+1}{h} \left(\int_{T_n} \mathcal{F}_{j+1/2}^n(t) dt - (-1)^k \int_{T_n} \mathcal{F}_{j-1/2}^n(t) dt \right. \\ \left. - \sum_{i=1}^{\mathcal{N}} f(\widehat{U}_{i,j}^n) \left\langle \frac{\partial \Phi_{k,j}}{\partial x}, \theta_{i,j} \right\rangle_{tx} \right).$$

These equations give the updates of $\widehat{u}_{k,j}^n$ from the time t_n to t_{n+1} . The numerical discrete solution updated at the new time t_{n+1} is $u^{n+1}(x) = \sum_{j=1}^Z \sum_{k=0}^N \widehat{u}_{k,j}^{n+1} \Phi_{k,j}(x)$ for $(t, x) \in T_n \times I$.

6. NUMERICAL STUDIES

Starting from this section, for brevity we say only the $P_N P_M$ schemes.

Let us consider the scalar linear advection equation $v_t(t, x) + av_x(t, x) = 0$, for $t \in [0, T]$ with $T > 0$, $x \in I = [\varepsilon_1, \varepsilon_2] \subset \mathbb{R}$, and $a > 0$. Suppose that the initial solution is $v_0 = v(0, \cdot) \in L^2(I)$. The exact solution is given by $v_e(t, x) = v_0(x - at)$. We apply the $P_N P_M$ schemes and use the modified Lax-Friedrichs flux given by Cockburn and Shu [4]. We have

$$\mathcal{F}_{LF,j+1/2}^n(t) = \frac{1}{2}[(a - |a|)U_{j+1}^n(t, x_{j+1/2}) + (a + |a|)U_j^n(t, x_{j+1/2})] = aU_j^n(t, x_{j+1/2}),$$

or

$$\mathcal{F}_{LF,j+1/2}^n(t) = aU_j^n(t, x_{j+1/2}) = a \sum_{i=1}^{\mathcal{N}} \widehat{U}_{i,j}^n \theta_{i,j}(t, x_{j+1/2})$$

and

$$\mathcal{F}_{LF,j-1/2}^n(t) = aU_{j-1}^n(t, x_{j-1/2}) = a \sum_{i=1}^{\mathcal{N}} \widehat{U}_{i,j-1}^n \theta_{i,j-1}(t, x_{j-1/2}).$$

Substituting into (5.3) we obtain

$$\begin{aligned} \widehat{u}_{k,j}^{n+1} = \widehat{u}_{k,j}^n &- \frac{(2k+1)a}{h} \sum_{i=1}^{\mathcal{N}} \left\{ \widehat{U}_{i,j}^n \int_{T_n} \theta_{i,j}(t, x_{j+1/2}) dt \right. \\ &\left. - (-1)^k \widehat{U}_{i,j-1}^n \int_{T_n} \theta_{i,j-1}(t, x_{j-1/2}) dt - \widehat{U}_{i,j}^n \left\langle \frac{\partial \Phi_{k,j}}{\partial x}, \theta_{i,j} \right\rangle_{tx} \right\}. \end{aligned}$$

As examples of these solutions we give some formulas in Appendix C.

There are several parameters which control the $P_N P_M$ schemes.

- (1) The orders N and M .
- (2) The size n_e of any stencil $S_{I_j, n_e, L}$, that must satisfy the condition $n_e \geq (M+1)/(N+1)$.
- (3) The index L of the stencil $S_{I_j, n_e, L}$ that indicates the form of the stencil.
- (4) The mesh size Z which gives the length h of the elements.
- (5) The maximal time value T with the time step $\Delta t \leq T$.
- (6) The Courant number $\lambda = (|a|\Delta t)/h$ which relates the time step Δt to the mesh length h , see [3].

In the following we study the stability and efficiency of the $P_N P_M$ schemes by studying three of these parameters, namely, the Courant number as well as the size and form of the stencils.

6.1. Stability analysis. The Courant number is important for the stability of the schemes. We determine maximal Courant numbers which are limits of the stability. We study the $P_N P_M$ schemes applying them to the linear advection equation $v_t + v_x = 0$ with $a = 1$.

6.1.1. Von Neumann analysis. We apply the von Neumann stability analysis [8] in the special case $N = 0$. The computational domain of the Fourier representations is the region $[-z, z]$ which is discretized into $2Z_f$ mesh elements with equidistant length element $h_f = z/Z_f$ and $z \in \mathbb{R}$ is the period of the initial data. We decompose the coefficients $\widehat{u}_{0,j}^n$ inside the element I_j , into a Fourier sum as $\widehat{u}_{0,j}^n = \sum_{l=-Z_f}^{Z_f} \mathcal{A}_l^n e^{ij\varphi_l}$ where \mathcal{A}_l^n is called the amplitude vector at time level t_n , $i = \sqrt{-1}$ is the imaginary unit, and φ_l is the wave number which is given by $\varphi_l = l\pi/Z_f$ with $l = -Z_f, \dots, Z_f$. This finite sum splits the time dependence from the spatial one, where the time evolution is included in the time dependence of the amplitude \mathcal{A}_l^n .

Now we substitute this finite sum into the scheme considered. Then, dividing by $e^{ij\varphi_l}$, we obtain a relation between the amplitude vectors \mathcal{A}_l^n and \mathcal{A}_l^{n+1} with

some space shifts $e^{\mp i\varphi_l}$. This relation can be written as $\mathcal{A}_l^{n+1} = D_l \mathcal{A}_l^n$, where D_l is called the amplification factor for $l = -Z_f, \dots, Z_f$.

The stability condition of the von Neumann analysis states that the Euclidean norm of the amplitude vector \mathcal{A}_l^n for any wave number φ_l does not grow in time. This condition is written as $|D_l| \leq 1$ for all φ_l .

For the P_0P_0 scheme, we can obtain the Courant numbers $\lambda = \Delta t/h$ that give stability exactly. We have $0 < \lambda \leq 1$.

For those P_NP_M schemes with $N = 0$, we determine the maximal Courant numbers numerically. The amplification factor D_l is a function of two variables $D_l = D_l(\varphi_l, \lambda)$. We take $Z_f = 3$, then $l = -3, \dots, 3$ and

$$\varphi_l \in \{-\pi, -2\pi/3, -\pi/3, 0, \pi/3, 2\pi/3, \pi\},$$

and define the variable $\lambda_s = s/10$ with $1 \leq s \leq 30$, which covers the interval $[1/10, 3]$. Then we compute the modulus of D_l at each value of φ_l and of λ , then we get a 7×30 matrix of these values. Each column is related to one value of λ_s . If all entries of the column are less than or equal to one then the value λ_s , to which this column is associated, gives a stable solution of the scheme.

For example, for the P_0P_1 scheme with the stencil $S_{I_j,2,0}$, we obtain the matrix

$$\left(\begin{array}{cccccccc} 0.98 & 0.92 & \dots & 0.62 & 1 & 1.42 & 1.88 & \dots \\ 0.98 & 0.95 & \dots & 0.80 & 1 & 1.25 & 1.55 & \dots \\ 0.99 & 0.99 & \dots & 0.98 & 1 & 1.03 & 1.07 & \dots \\ 1 & 1 & \dots & 1 & 1 & 1 & 1 & \dots \\ 0.99 & 0.99 & \dots & 0.98 & 1 & 1.03 & 1.07 & \dots \\ 0.98 & 0.95 & \dots & 0.80 & 1 & 1.25 & 1.55 & \dots \\ 0.98 & 0.92 & \dots & 0.62 & 1 & 1.42 & 1.88 & \dots \\ \uparrow & \uparrow & & \uparrow & \uparrow & \uparrow & & \\ \lambda_1 = 0.1 & \lambda_2 = 0.2 & & \lambda_9 = 0.9 & \lambda_{10} = 1 & \lambda_{11} = 1.1 & & \end{array} \right).$$

Note that, starting from the eleventh column, the entries are larger than one. This proves that the value $\lambda_{11} = 1.1$ gives an unstable solution. Thus the maximum value of λ_s which gives a stable solution is λ_{10} , thus $\lambda_{\max} \approx \lambda_{10} = 1$. On the other hand, all columns to the left of the eleventh have entries less than or equal to one, thus their λ_s give stable solutions. We obtain that the range of the Courant number for the P_0P_1 scheme using the stencil $S_{I_j,2,0}$ is the interval $\lambda \in (0, 1]$.

Table 1 includes the maximal limits λ_{\max} , which are computed numerically in this way, of the Courant numbers for the P_0P_M schemes with $N = 0$ and various orders $M = 1, \dots, 5$ with all cases of the stencils $S_{I_j, n_e, L}$ with $n_e = [(M + 1)/(N + 1)], \dots, 6$ and $L = 0, \dots, n_e - 1$. The symbol * indicates unstable

n_e	L	P_0P_1	P_0P_2	P_0P_3	P_0P_4	P_0P_5
2	0	(0, 1]				
	1	(0, 2]				
3	0	(0, 1]	*			
	1	(0, 1]	(0, 1]			
	2	(0, 1]	[1, 2]			
4	0	(0, 1]	*	*		
	1	(0, 1]	(0, 1]	(0, 1]		
	2	(0, 1]	(0, 1]	(0, 2]		
	3	(0, 1]	*	[1, 2]		
5	0	(0, 1]	(0.5, 1]	*	*	
	1	(0, 1]	(0, 1]	(0, 1]	*	
	2	(0, 1]	(0, 1]	(0, 1]	(0, 1]	
	3	(0, 1]	(0, 1]	(0, 1]	(0, 1]	[1, 2]
	4	(0, 1]	*	*	*	
6	0	(0, 1]	(0, 1]	*	*	*
	1	(0, 1]	(0, 1]	(0, 1]	*	*
	2	(0, 1]	(0, 1]	(0, 1]	(0, 1]	(0, 1]
	3	(0, 1]	(0, 1]	(0, 1]	(0, 1]	(0, 1]
	4	(0, 1]	(0, 1]	(0, 1]	(0, 1]	*
	5	(0, 1]	(0, 1]	*	[1, 2]	*

Table 1. The maximal Courant numbers for some $P_N P_M$ schemes, for $N = 0$ and $M = 1, 2, 3, 4, 5$.

cases for which we have only one value $\lambda = 1$ that gives a stable solution. The fact that $\lambda = 1$ is stable is an artifact due to the equation $v_t + v_x = 0$, because for $\lambda = 1$ the numerical solution of these schemes is the exact solution. Moreover, the range $(0, 1]$ mostly appears, but there are some semi-stable cases which are written in boldface. Also, there are two cases with higher stability $\lambda \in (0, 2]$ that are also highlighted in boldface.

6.1.2. Another experimental procedure. The von Neumann analysis for higher order $P_N P_M$ schemes with $N > 0$ is not possible without the use of computer algebra and numerical computation, see Dumbser [5], p. 8221. Therefore, we consider another numerical procedure. We continue in the study of the advection equation $v_t + v_x = 0$ with the initial solution $v_0(x) = \sin x$ for $x \in [0, 2\pi]$ and periodicity as the boundary condition. So we have $v_0(x) = \sin x$ for all $x \in \mathbb{R}$. It is well-known that the exact solution is $v_e(t, x) = v_0(x - t)$ on $[0, T] \times [0, 2\pi]$.

We found experimentally appropriate limits of the Courant numbers which guarantee the stability without resorting to von Neumann analysis. We checked the stability

of the numerical solutions at the final time $T = 100\pi$ for a mesh with $Z = 50$. Let us set $\lambda_C := 1/(2N + 1)$. Cockburn [3] considered Runge-Kutta DG schemes and took for the linear equations λ somewhat smaller than λ_C as a limit in order to avoid unstable solutions. We start with this inequality and define variables λ_s in an interval around λ_C as $\lambda_s = \alpha_N \lambda_C + 0.001s$ for $s = 0, 1, 2, \dots$, where $0 < \alpha_N < 1$ is a constant associated to the order N and determines the starting point of the search algorithm. We use the values $\alpha_0 = 0.99$, $\alpha_1 = 0.9$, $\alpha_2 = 0.8$, $\alpha_3 = 0.7$, $\alpha_4 = 0.6$, and $\alpha_5 = 0.5$.

Increasing the index s , the variable λ_s comes closer to the ratio λ_C and then larger than λ_C . For each value λ_s , we associate the value $L_s^1 = \int_I |v_e(T, x) - w(T, x)| dx$, which is the L^1 error of the reconstructed polynomial (solution) w computed using the $P_N P_M$ scheme at the last time T with time step $\Delta t = \lambda_s h$. We compute the errors L_s^1 numerically using Gaussian rules of orders large enough. As well, we compute the differences $d_{1,s} = |L_s^1 - L_{s-1}^1|$ for $s > 0$, defining $d_{1,0} = 0$. We also set a condition to stop this algorithm which is $d_{1,s} > \text{TOL}$, where TOL is a tolerance that we choose large enough, e.g. $\text{TOL} = 10$, to guarantee that the L^1 error is large, and this means that the solution is unstable. For example, we consider the $P_2 P_2$ scheme. Then we have $\lambda_C = 0.2$, $\alpha_2 = 0.8$ and $\lambda_s = 0.16 + 0.001s$. We arrange the errors starting from $s = 4$ in the following table, which leads to the conclusion that $\lambda_{\max} \approx 0.171$.

s	4	5	6	7	8	9	10	11	12
λ_s	0.164	0.165	0.166	0.167	0.168	0.169	0.170	0.171	0.172
L_s^1	0.009	0.040	0.064	0.079	0.005	0.009	0.011	0.018	4×10^{74}

Note that in the solution plots the solution for $\lambda_s = 0.170$ looks smooth, whereas for $\lambda_s = 0.171$ small oscillations occur that become stronger for larger λ_s . In the following, we give the approximate values of λ_{\max} for all $P_N P_M$ schemes for

$$M = 0, \dots, 5, \quad N = 0, \dots, M, \quad n_e = \left\lceil \frac{M+1}{N+1} \right\rceil, \dots, 6, \quad L = 0, \dots, n_e - 1.$$

Case $M = 1$. For the $P_1 P_1$ scheme we obtain $\lambda_C = 0.333$ and $\lambda_{\max} \approx 1/3$. For the $P_0 P_1$, see Table 2.

$n_e \backslash L$	0	1	2	3	4	5
2	1.003	2.006				
3	1.005	1.006	1.008			
4	1.005	1.006	1.007	1.009		
5	1.005	1.005	1.006	1.008	1.007	
6	1.006	1.006	1.005	1.008	1.007	1.007

Table 2. The maximal Courant numbers for $P_0 P_1$ scheme with $\lambda_C = 1$.

Case $M = 2$. For the P_2P_2 scheme we obtain $\lambda_C = 0.2$ and $\lambda_{\max} \approx 0.17$. Also, Table 3 gives the approximations of λ_{\max} for the P_0P_2 and P_1P_2 schemes.

The P_0P_2 schemes with $\lambda_C = 1$						
$n_e \backslash L$	0	1	2	3	4	5
3	1.002	1.01	[1,2]			
4	1.013	1.012	1.005	1.011		
5	1.003	1.01	1.006	1.005	1.011	
6	1.004	1.007	1.006	1.006	1.006	1.01
The P_1P_2 schemes with $\lambda_C = 0.333$						
$n_e \backslash L$	0	1	2	3	4	5
2	1/3	*				
3	1/3	1/3	*			
4	1/3	1/3	*	*		
5	1/3	1/3	1/3	*	*	
6	1/3	1/3	1/3	*	*	0.305

Table 3. The maximal Courant numbers for P_0P_2 and P_1P_2 schemes.

Case $M = 3$. For the P_2P_3 schemes where $\lambda_C = 0.2$ and with all stencils considered above we obtain $\lambda_{\max} \approx 0.17$ and for the P_3P_3 scheme where $\lambda_C = 0.143$ we find $\lambda_{\max} \approx 0.103$. For the P_0P_3 and P_1P_3 schemes, see Table 4.

The P_0P_3 schemes with $\lambda_C = 1$						
$n_e \backslash L$	0	1	2	3	4	5
4	1.001	1.005	2.009	[1,2]		
5	1.002	1.01	1.006	1.005	1.02	
6	1.002	1.011	1.006	1.006	1.006	1.017
The P_1P_3 schemes with $\lambda_C = 0.333$						
$n_e \backslash L$	0	1	2	3	4	5
2	0.318	*				
3	0.328	0.34	*			
4	0.331	0.33	0.338	*		
5	0.332	1/3	0.332	*	*	
6	0.332	0.332	0.316	0.335	*	*

Table 4. The maximal Courant numbers for P_0P_3 and P_1P_3 schemes.

Case $M = 4$. For the P_3P_4 schemes with $\lambda_C = 0.143$ we obtain $\lambda_{\max} \approx 0.103$ and for the P_4P_4 scheme where $\lambda_C = 0.111$ we find $\lambda_{\max} \approx 0.069$. For the P_0P_4 , P_1P_4 , and P_2P_4 schemes, see Table 5.

The P_0P_4 schemes with $\lambda_C = 1$						
$n_e \backslash L$	0	1	2	3	4	5
5	1	1.002	1.012	2.02	[1,1.5]	
6	1	1.006	1.012	1	1	[1,2]
The P_1P_4 schemes with $\lambda_C = 0.333$						
$n_e \backslash L$	0	1	2	3	4	5
3	0.316	0.346	*			
4	0.325	0.347	*	*		
5	0.328	0.338	0.337	*	*	
6	0.330	0.338	0.337	*	0.312	*
The P_2P_4 schemes with $\lambda_C = 0.2$						
$n_e \backslash L$	0	1	2	3	4	5
2	0.166	*				
3	0.169	0.176	*			
4	0.170	0.173	0.173	*		
5	0.170	0.170	0.172	0.170	0.170	
6	0.170	0.170	0.172	0.172	0.170	0.170

Table 5. The maximal Courant numbers for P_0P_4 , P_1P_4 , and P_2P_4 schemes.

Case $M = 5$. For the P_4P_5 schemes where $\lambda_C = 0.111$ we obtain $\lambda_{\max} \approx 0.069$ and for the P_5P_5 scheme where $\lambda_C = 0.091$ we find $\lambda_{\max} \approx 0.05$. For the P_0P_5 , P_1P_5 , P_2P_5 , and P_3P_5 schemes, see Table 6.

6.2. Experimental order of convergence (EOC). We investigate the orders of the accuracy numerically by calculating the EOC. Let generally X be a linear space with some norm $\|\cdot\|_X$ and let $v_h \in X$ be a numerical approximation of a given function $v \in X$ which depends on a parameter h of the discretization. The convergence of v_h towards v as h tends to zero can be quantified by $\|v_h - v\|_X \leq Ch^\kappa$, with the order of convergence κ . This gives a possibility to quantify the quality of a numerical scheme. If we can compute two numerical solutions v_h and $v_{h'}$, then the order κ can be estimated experimentally by

$$\kappa \simeq \text{EOC}(h, h') = \frac{\log(\|v_{h'} - v\|_X / \|v_h - v\|_X)}{\log(h'/h)}.$$

The maximum Courant numbers computed above are quite sharp, since oscillations occur with slightly larger time steps. We observe that the stability limits depend strongly on N and not really on M . Therefore, for our further tests, we used the restrictive bounds on the Courant number given in Table 7.

The P_0P_5 schemes with $\lambda_C = 1$						
$n_e \backslash L$	0	1	2	3	4	5
6	1	1.001	1.005	[1,2]	[1,3]	1
The P_1P_5 schemes with $\lambda_C = 0.333$						
$n_e \backslash L$	0	1	2	3	4	5
3	*	0.402	*			
4	*	0.346	*	*		
5	*	0.345	0.335	*	*	
6	0.324	0.344	0.327	0.34	*	*
The P_2P_5 schemes with $\lambda_C = 0.2$						
$n_e \backslash L$	0	1	2	3	4	5
2	*	*				
3	0.165	0.176	*			
4	0.167	0.176	0.175	*		
5	0.168	0.175	0.172	0.174	*	
6	0.169	0.172	0.172	0.172	0.172	*
The P_3P_5 schemes with $\lambda_C = 0.143$						
$n_e \backslash L$	0	1	2	3	4	5
2	0.1	*				
3	0.103	0.106	0.102			
4	0.103	0.105	0.104	0.103		
5	0.103	0.103	0.104	0.103	0.103	
6	0.103	0.103	0.104	0.104	0.103	0.103

Table 6. The maximal Courant numbers for P_0P_5 , P_1P_5 , P_2P_5 , and P_3P_5 schemes.

The order N	0	1	2	3	4	5
The Courant number λ_{used}	1	0.25	0.16	0.08	0.05	0.05

Table 7. The Courant numbers λ_{used} for $N = 0, \dots, 5$.

Now we consider the advection equation $v_t + v_x = 0$ with $v_0(x) = \sin x$ defined on $I = [0, 2\pi]$ and its solution at time $T = 2\pi$. We apply some P_NP_M schemes. The CFL numbers λ are taken from Table 7. The L^1 errors are listed in Table 8, where we always used the stencil $S_{I_j,5,2}$. The numbers for the EOC were truncated after the first decimal. Note that we always get the expected order of convergence close to $M + 1$. Some of the schemes produce a wrong experimental order on the coarsest meshes. This is not a problem, since the order is an asymptotic property for $h \rightarrow 0$.

6.3. The study of the efficiency. We again consider the advection equation $v_t + v_x = 0$ with the initial function $v_0(x) = \sin x$ defined on $I = [0, 2\pi]$ and its

Z	L^1	EOC	L^1	EOC	L^1	EOC	L^1	EOC	L^1	EOC
P_0P_0										
10	6.2e-1									
20	3.1e-1	0.9								
40	1.5e-1	1.0								
P_0P_1			P_1P_1							
10	1.5e-1		1.6e-1							
20	2.3e-2	2.6	4.1e-2	2.0						
40	4.1e-3	2.4	1.0e-2	2.0						
P_0P_2			P_1P_2		P_2P_2					
10	1.3e-1		2.6e-2		2.0e-1					
20	1.8e-2	2.8	2.1e-3	3.6	9.8e-4	7.6				
40	2.2e-3	2.9	2.1e-4	3.3	1.2e-4	3.0				
P_0P_3			P_1P_3		P_2P_3		P_3P_3			
10	7.8e-3		2.1e-2		2.0e-1		1.9e-4			
20	4.5e-4	4.1	1.3e-3	3.9	2.6e-5	12.8	1.2e-5	3.9		
40	2.7e-5	4.0	8.6e-5	3.9	1.2e-6	4.4	7.8e-7	3.9		
P_0P_4		P_1P_4		P_2P_4		P_3P_4		P_4P_4		
10	3.5e-3		1.2e-3		2.0e-1		1.2e-5		1.2e-1	
20	1.1e-4	4.8	3.5e-5	5.1	2.1e-5	13.2	2.3e-7	5.7	6.2e-2	1.0
40	3.7e-6	4.9	1.0e-6	5.0	6.7e-7	4.9	5.8e-9	5.3	5.6e-9	23.4
80									1.7e-10	4.9

Table 8. The L^1 errors and EOC of some P_NP_M schemes applied to the advection equation.

solution at time $T = 2\pi$. The CFL numbers λ were taken from Table 7. We study the efficiency of the P_NP_M schemes by setting the bound for the L^1 errors at time $T = 2\pi$ to be 0.01. We measure the speed of the schemes by the computational time and the number of time steps Z_1 . Since we consider the linear advection equation, the time step Δt is constant and then it is equal to $\Delta t = T/Z_1 = 2\pi/Z_1$. Also the time step is computed using the Courant number λ by $\Delta t = \lambda h/a$. Here we have $a = 1$ and $h = 2\pi/Z$, then $\Delta t = \lambda 2\pi/Z$. Thus we obtain $2\pi/Z_1 = \lambda 2\pi/Z$ which implies that $Z_1 = Z/\lambda$. A further indicator of the cost of the discretization is the mesh size Z . To explain how we perform these computations we take as an example the P_0P_1 scheme using the stencil $S_{I_j,2,1}$ and take the mesh size Z changing from $Z = 2$ to $Z = 35$. We ended the computation when the L^1 error became lower than 0.01. For brevity, we give only some of these results for $Z = 28, \dots, 35$. Table 9 shows that, when $Z = 33$, it is the first case where the L^1 error is less than 0.01. In this case we need 33 iterations and a computational time of 0.049 seconds.

Now we will only give the data for the solution that satisfies the error bound on the coarsest mesh, which we obtain from a sequence of finer and finer meshes as explained. The errors will be rounded to 4 decimals.

L^1	0.0132681	0.0123711	0.0115621	0.0108299	0.0101650	0.0095595
time	0.0376	0.0364	0.0399	0.0460	0.0458	0.0490
Z_1	28	29	30	31	32	33
Z	28	29	30	31	32	33

Table 9. The computational time and the mesh size for the P_0P_1 scheme.

6.3.1. The influence of the size n_e . We recall that the reconstruction stencil is given by $S_{I_j, n_e, L} = \bigcup_{c=-L}^R I_{j+c}$ and consists of the interval I_j with L and R elements to the left and right of I_j , respectively, and its size is given by $n_e = 1 + L + R$ with $L \in \{0, \dots, n_e - 1\}$ and $R \geq 0$. We used various stencils with different sizes n_e and fixed the index L at the values $L = 0$ and $L = n_e - 1$, see Table 10.

n_e	L^1	time	Z_1	Z	n_e	L^1	time	Z_1	Z
P_0P_1					P_0P_1				
2	0.00955	0.00898	33	33	2	0.00955	0.01030	33	33
3	0.00904	0.01328	45	45	3	0.00904	0.01288	45	45
4	0.00961	0.01471	52	52	4	0.00961	0.01511	52	52
5	0.00903	0.01872	61	61	5	0.00903	0.01741	61	61
6	0.00974	0.01890	65	65	6	0.00974	0.01905	65	65
P_1P_1					P_1P_1				
1	0.00879	0.01887	116	29	1	0.00879	0.01991	116	29
P_0P_2					P_0P_2				
3	0.00626	0.00586	19	19	3	0.00626	0.00630	19	19
4	0.00991	0.00657	21	21	4	0.00991	0.00638	21	21
5	0.00802	0.00768	27	27	5	0.00802	0.00826	27	27
6	0.00982	0.00915	29	29	6	0.00982	0.00919	29	29
P_1P_2					P_1P_2				
2	0.00851	0.01446	56	14	2	unstable			
3	0.00537	0.01459	76	19	3	unstable			
4	0.00721	0.01534	76	19	4	unstable			
5	0.00895	0.01483	76	19	5	unstable			
6	0.00914	0.01647	80	20	6	0.00897	0.01508	76	19
P_2P_2					P_2P_2				
1	0.00203	0.01236	75	12	1	0.00203	0.01313	75	12

Table 10. Numerical computations for some P_NP_M schemes with $M = 1$ and $M = 2$ for two values of L , $L = 0$ (left) and $L = n_e - 1$ (right).

For $M = 1$, we have two schemes, the P_0P_1 scheme with various stencils and the P_1P_1 scheme with the unique stencil $S_{I_j, 1, 0} = I_j$. In all cases $N = M$ we have $n_e = 1$, since there is no reconstruction needed. Table 10 shows that the P_0P_1

scheme is faster than the P_1P_1 scheme. This is expected, since the piecewise constant solution P_0P_1 scheme has only one unknown degree of freedom. But the P_1P_1 scheme has a higher accuracy on the same mesh. Also, we find that the computational time grows when the size of stencil becomes larger; again this is expected, since the information comes from more cells. Thus the size of the stencil has negative influence on the efficiency of the scheme, as expected. An important point is that larger stencils need more grid points to achieve the same accuracy.

For $M = 2$, we have the P_0P_2 , P_1P_2 , and P_2P_2 schemes. Table 10 shows that the P_0P_2 scheme is faster than the others. Comparing tables we see that whereas in Table 8 on the same spatial mesh the error decreases from P_0P_2 to P_1P_2 to P_2P_2 schemes, on the other hand, in terms of the actual efficiency using the smallest possible stencil in Table 10 the order in terms of computational time is reversed. This is despite the fact that the other schemes need fewer mesh points to achieve the same accuracy. However, they need more time steps due to their stability restrictions. Note also that there is no real difference between choosing the larger stencils in an upwind $L = n_e - 1$ or a downwind $L = 0$ manner.

In Table 11 we now compare the case $L = 0$ taking the smallest stencil for the different mesh sizes. The computational time of the P_0P_2 scheme is the smallest using different meshes comparing with the P_1P_2 and P_2P_2 schemes. Again we see that the stability is crucial for the comparison since severer stability limits lead to a larger number of time steps.

$L = 0$									
Z	P_0P_2 with $n_e = 3$		Z_1	P_1P_2 with $n_e = 2$		Z_1	P_2P_2		
	L^1	time		L^1	time		L^1	time	Z_1
10	0.04172	0.01471	10	0.02313	0.01881	40	0.08998	0.02117	63
11	0.99828	0.00517	12	0.25449	0.00995	45	0.04117	0.01186	69
12	0.02444	0.00462	12	0.23295	0.01008	49	0.00203	0.01282	75
13	0.84816	0.00503	14	0.01062	0.01122	52	0.10307	0.01438	82
14	0.01550	0.00806	14	0.00851	0.01421	56	0.06386	0.01830	88
15	0.73693	0.00530	16	0.18595	0.01169	61	0.02988	0.01558	94

Table 11. Numerical computations for some P_NP_M schemes with $M = 2$ using the smallest possible stencil.

6.3.2. The influence of varying L . We now use stencils of the same size n_e but with different type for the values $L = 0, \dots, n_e - 1$. We choose $M = 3$ and $n_e = 5$.

We note in Table 12 that the symmetric stencil with $L = 2$ is the best choice as concerns to the computational time and the spatial discretization. On the other hand, the one side stencils with $L = 0$ and $L = 4$ require slightly longer computational time. The difference in the choice of the stencil is not very pronounced. Moreover,

for the finite volume scheme P_0P_3 , the number of iterations relates to the type of the stencil, whereas with $N > 0$ this number seems to be constant. This is seen also in Table 10.

L	L^1	time	Z_1	Z	L	L^1	time	Z_1	Z
P_0P_3					P_2P_3				
0	0.00676	0.00714	16	16	0	0.00491	0.01139	50	8
1	0.00784	0.00540	12	12	1	0.00405	0.01156	50	8
2	0.00940	0.00421	8	8	2	0.00097	0.01165	50	8
3	0.00784	0.00485	12	12	3	0.00350	0.01178	50	8
4	0.00676	0.00632	16	16	4	0.00440	0.01220	50	8
P_1P_3					P_3P_3				
0	0.00772	0.01277	56	14	0	0.00009	0.02086	125	10
1	0.00661	0.01146	40	10					
2	0.00931	0.01017	40	10					
3		unstable							
4		unstable							

Table 12. The computational time and the mesh size of some P_NP_M schemes for $M = 3$ with different types of the stencils of the size $n_e = 5$.

Furthermore, when $N = 0$, the number of iterations Z_1 is equal to the mesh size Z , whereas for $N > 0$, this number is larger than Z by a factor due to the stability restriction. This indeed means that with larger N the cost of the computations is larger, but this improves the accuracy. This agrees with the results in Table 8 where we find for example that for $Z = 40$ the P_3P_3 scheme is more accurate than the P_2P_3 scheme which is in turn more accurate than the P_1P_3 and P_0P_3 schemes.

In Table 13 we again compare the computational time for some P_NP_M schemes with $M = 3$ using the smallest stencil for different mesh sizes. The computational

Z	P_0P_3		P_0P_3		P_1P_3		P_2P_3	
	L^1	time	L^1	time	L^1	time	L^1	time
10	0.00705	0.00375	0.00705	0.00360	0.00405	0.00986	0.08842	0.01253
11	0.99794	0.00443	0.99940	0.00436	0.25015	0.00963	0.04004	0.01684
12	0.00345	0.00458	0.00345	0.00428	0.22991	0.01044	0.00038	0.01435
13	0.84801	0.00522	0.84866	0.00480	0.00141	0.01166	0.10254	0.01487
14	0.00188	0.00502	0.00188	0.00478	0.00104	0.01117	0.06345	0.01571
15	0.73685	0.00664	0.73718	0.00604	0.18474	0.01339	0.02956	0.01866

Table 13. Numerical computations for some P_NP_M schemes with $M = 3$ using the smallest possible stencil.

time of the P_0P_3 scheme is the smallest using the different meshes comparing with the P_1P_2 and P_2P_2 schemes. Again we see that the stability is crucial to the comparison since severer stability limits lead to a larger number of time steps.

7. CONCLUSIONS

We have considered the P_NP_M DG schemes with $N \leq M$ introduced by Dumbser et al. [5] for $N, M = 0, 1, 2, 3, 4, 5$, and the simplest advection equation. Depending on the quotient $(M + 1)/(N + 1)$ we took into account where ever possible different stencils for the reconstruction.

All the allowed combinations $0 \leq N \leq M \leq 5$ had some stencils with stability for all CFL numbers between 0 and a maximal CFL number. We found a wide range of maximal stability limits being CFL numbers between 0.103 and 2. Some stencils have a strange semi-stability behaviour since they are stable for CFL numbers in an interval bounded away from 0. Also some stencils lead to unstable schemes. The stability limits depend on the parameter N and not on M .

Using the stability limits that we obtained, we checked the experimental order of convergence (EOC). We report only the cases $0 \leq N \leq M \leq 4$ for the stencil $S_{I_j,5,2}$. We always obtain an expected EOC close to $M + 1$, also in other cases we did not put into the paper.

Based on the stability limits of the various schemes we also studied the efficiency of the schemes. We found that for a given M the P_0P_M schemes are faster than the others with $M \geq N > 0$. Also, we found that the computational time grows when the size of the stencil becomes larger and there was no real difference between choosing the larger stencils in an upwind $L = n_e - 1$ or a downwind $L = 0$ manner. We noted that the symmetric stencil, i.e. with $n_e > 1$ odd and $L = n_e/2 - 1$, achieves the required accuracy on a coarser mesh leading to a faster computation in comparison to the asymmetric stencils of the same size.

APPENDIX A: EXAMPLES OF COMPUTING THE COEFFICIENTS $\widehat{w}_{i,j}^n$

7.1. Example with $n_e(N + 1) = (M + 1)$. Let $N = 0, M = 2, L = 1, R = 1$. Then we have $n_e = 1 + L + R = 3$. The stencil is $S_{I_j,3,1}$ of three elements. The system of normal equations is

$$\begin{aligned} c = -1 &\rightarrow h\widehat{w}_{0,j}^n - 2h\widehat{w}_{1,j}^n + 6h\widehat{w}_{2,j}^n = h\widehat{u}_{0,j-1}^n, \\ c = 0 &\rightarrow h\widehat{w}_{0,j}^n = h\widehat{u}_{0,j}^n, \end{aligned}$$

$$\begin{aligned}
c = 1 &\rightarrow h\widehat{w}_{0,j}^n + 2h\widehat{w}_{1,j}^n + 6h\widehat{w}_{2,j}^n = h\widehat{u}_{0,j+1}^n \\
&\Rightarrow \begin{cases} \widehat{w}_{0,j}^n = \widehat{u}_{0,j}^n, \\ \begin{pmatrix} -2 & 6 \\ 2 & 6 \end{pmatrix} \begin{pmatrix} \widehat{w}_{1,j}^n \\ \widehat{w}_{2,j}^n \end{pmatrix} = \begin{pmatrix} \widehat{u}_{0,j-1}^n - \widehat{u}_{0,j}^n \\ \widehat{u}_{0,j+1}^n - \widehat{u}_{0,j}^n \end{pmatrix}. \end{cases}
\end{aligned}$$

The coefficients matrix is invertible. The solution is unique:

$$\widehat{w}_{1,j} = \frac{1}{4}(\widehat{u}_{0,j+1} - \widehat{u}_{0,j-1}), \quad \widehat{w}_{2,j} = \frac{1}{12}(\widehat{u}_{0,j+1} - 2\widehat{u}_{0,j} + \widehat{u}_{0,j-1}).$$

7.2. Example with $n_e(N+1) > (M+1)$. Let $N = 1, M = 2, L = 1, R = 1$. The system of normal equations is

$$\begin{aligned}
c = -1 &\begin{cases} h\widehat{w}_{0,j}^n - 2h\widehat{w}_{1,j}^n + 6h\widehat{w}_{2,j}^n = h\widehat{u}_{0,j-1}^n, \\ \frac{1}{3}h\widehat{w}_{1,j}^n - 2h\widehat{w}_{2,j}^n = \frac{1}{3}h\widehat{u}_{1,j-1}^n, \end{cases} \\
c = 0 &\begin{cases} h\widehat{w}_{0,j}^n = h\widehat{u}_{0,j}^n, \\ \frac{1}{3}h\widehat{w}_{1,j}^n = \frac{1}{3}h\widehat{u}_{1,j}^n, \end{cases} \\
c = 1 &\begin{cases} h\widehat{w}_{0,j}^n + 2h\widehat{w}_{1,j}^n + 6h\widehat{w}_{2,j}^n = h\widehat{u}_{0,j+1}^n, \\ \frac{1}{3}h\widehat{w}_{1,j}^n + 2h\widehat{w}_{2,j}^n = \frac{1}{3}h\widehat{u}_{1,j+1}^n, \\ \widehat{w}_{0,j}^n = \widehat{u}_{0,j}^n, \quad \widehat{w}_{1,j}^n = \widehat{u}_{1,j}^n, \\ \begin{pmatrix} 6 \\ -2 \\ 6 \\ 2 \end{pmatrix} (\widehat{w}_{2,j}^n) = \begin{pmatrix} \widehat{u}_{0,j-1}^n - \widehat{u}_{0,j}^n + 2\widehat{u}_{1,j}^n, \\ \frac{1}{3}\widehat{u}_{1,j-1}^n - \frac{1}{3}\widehat{u}_{1,j}^n \\ \widehat{u}_{0,j+1}^n - \widehat{u}_{0,j}^n - 2\widehat{u}_{1,j}^n \\ \frac{1}{3}\widehat{u}_{1,j+1}^n - \frac{1}{3}\widehat{u}_{1,j}^n \end{pmatrix}. \end{cases}
\end{aligned}$$

The solution is non-unique. The least squares solution is

$$\widehat{w}_{2,j} = \frac{1}{120}(9\widehat{u}_{0,j+1} + \widehat{u}_{1,j+1} - 18\widehat{u}_{0,j} + 9\widehat{u}_{0,j-1} - \widehat{u}_{1,j-1}).$$

APPENDIX B: THE LOCAL SPACE TIME BASIS FUNCTIONS

For $M = 2$ we have $\mathcal{N} = (M+1)(M+2)/2 = 6$. The nodes are taken as

$$\begin{aligned}
\beta_1 &= (t_n, x_{j-1/2}) & \beta_2 &= (t_n, x_j) & \beta_3 &= (t_n, x_{j+1/2}) \\
\beta_4 &= (t_{n+1/2}, x_{j-1/2}) & \beta_5 &= (t_{n+1/2}, x_{j+1/2}) & \beta_6 &= (t_{n+1}, x_j).
\end{aligned}$$

For $x \in I_j$ and $t \in T_n$ and by using $\varsigma = 2(x - x_j)/h$ and $\zeta = (t - t_n)/k$, the nodal basis functions are given by

$$\begin{aligned}\theta_{1,j}(\zeta(t), \varsigma(x)) &= -\frac{1}{2}\zeta + \frac{1}{2}\zeta^2 - 2\zeta + \zeta\zeta + 2\zeta^2, & \theta_{2,j}(\zeta(t), \varsigma(x)) &= 1 - \zeta^2 + \zeta - 2\zeta^2, \\ \theta_{3,j}(\zeta(t), \varsigma(x)) &= \frac{1}{2}\zeta + \frac{1}{2}\zeta^2 - 2\zeta - \zeta\zeta + 2\zeta^2, & \theta_{4,j}(\zeta(t), \varsigma(x)) &= 2\zeta - \zeta\zeta - 2\zeta^2, \\ \theta_{5,j}(\zeta(t), \varsigma(x)) &= 2\zeta + \zeta\zeta - 2\zeta^2, & \theta_{6,j}(\zeta(t), \varsigma(x)) &= -\zeta + 2\zeta.\end{aligned}$$

They satisfy $\theta_{i,j}(\beta_k) = \delta_{ik}$ for $i, k = 1, \dots, 6$. We set the basis to be $\Theta_{2,j} = \{\theta_{1,j}, \theta_{2,j}, \theta_{3,j}, \theta_{4,j}, \theta_{5,j}, \theta_{6,j}\}$. Note that at $t = t_n$ we have $\theta_{4,j}(t_n, x) = \theta_{5,j}(t_n, x) = \theta_{6,j}(t_n, x) = 0$ for all $x \in I_j$. The other functions $\theta_{1,j}$, $\theta_{2,j}$, and $\theta_{3,j}$ depend on the spatial points at $t = t_n$.

APPENDIX C: SOME FORMULAS OF THE $P_N P_M$ DG SOLUTIONS

Let Δt be the time step and h the mesh size. We will use the notation $\lambda = a\Delta t/h$. $P_0 P_0$ DG scheme

$$\widehat{u}_{0,j}^{n+1} = \widehat{u}_{0,j}^n + \lambda(\widehat{u}_{0,j-1}^n - \widehat{u}_{0,j}^n),$$

$P_1 P_1$ DG scheme

$$\begin{aligned}\widehat{u}_{0,j}^{n+1} &= \widehat{u}_{0,j}^n + \lambda(\widehat{u}_{0,j-1}^n - \widehat{u}_{0,j}^n + \widehat{u}_{1,j-1}^n - \widehat{u}_{1,j}^n) - \lambda^2(\widehat{u}_{1,j-1}^n - \widehat{u}_{1,j}^n), \\ \widehat{u}_{1,j}^{n+1} &= \widehat{u}_{1,j}^n - 3\lambda(\widehat{u}_{0,j-1}^n - \widehat{u}_{0,j}^n + \widehat{u}_{1,j-1}^n + \widehat{u}_{1,j}^n) + 3\lambda^2(\widehat{u}_{1,j-1}^n - \widehat{u}_{1,j}^n),\end{aligned}$$

$P_0 P_1$ DG scheme with the stencil $S_{I_j,2,1} = I_{j-1} \cup I_j$

$$\widehat{u}_{0,j}^{n+1} = \widehat{u}_{0,j}^n - \frac{1}{2}\lambda(\widehat{u}_{0,j-2}^n - 4\widehat{u}_{0,j-1}^n + 3\widehat{u}_{0,j}^n) + \frac{1}{2}\lambda^2(\widehat{u}_{0,j-2}^n - 2\widehat{u}_{0,j-1}^n + \widehat{u}_{0,j}^n).$$

Acknowledgment. The authors would like to thank Matthias Kunik for helpful discussions of the reconstruction problem.

References

- [1] *R. A. Adams*: Sobolev Spaces. Pure and Applied Mathematics 65, Academic Press, New York, 1975. [zbl](#) [MR](#) [doi](#)
- [2] *A. Badenjki*: The $P_N P_M$ DG Schemes for the One Dimensional Hyperbolic Conservation Laws. Doctoral Thesis, Otto-von-Guericke University, Magdeburg, 2018.
- [3] *B. Cockburn*: An introduction to the discontinuous Galerkin method for convection-dominated problems. Advanced Numerical Approximation of Nonlinear Hyperbolic Equations (A. Quarteroni et al., eds.). Lecture Notes in Mathematics 1697, Springer, Berlin, 1998, pp. 151–268. [zbl](#) [MR](#) [doi](#)

- [4] *B. Cockburn, C.-W. Shu*: TVB Runge Kutta local projection discontinuous Galerkin finite element method for conservation laws. II: General framework. *Math. Comput.* 52 (1989), 411–435. [zbl](#) [MR](#) [doi](#)
- [5] *M. Dumbser, D. S. Balsara, E. F. Toro, C.-D. Munz*: A unified framework for the construction of one-step finite volume and discontinuous Galerkin schemes on unstructured meshes. *J. Comput. Phys.* 227 (2008), 8209–8253. [zbl](#) [MR](#) [doi](#)
- [6] *L. C. Evans*: *Partial Differential Equations*. Graduate Studies in Mathematics 19, American Mathematical Society, Providence, 1998. [zbl](#) [MR](#) [doi](#)
- [7] *C. R. Goetz, M. Dumbser*: A square entropy stable flux limiter for $P_N P_M$ schemes. Available at <https://arxiv.org/abs/1612.04793> (2016), 24 pages.
- [8] *C. Hirsch*: *Numerical Computation of Internal and External Flows. Volume 1: Fundamentals of Numerical Discretization*. Wiley Series in Numerical Methods in Engineering; Wiley, Chichester, 1988. [zbl](#)
- [9] *T. H. Koornwinder, R. Wong, R. Koekoek, R. F. Swarttouw*: Orthogonal polynomials. NIST Handbook of Mathematical Functions (F. W. J. Olver et al., eds.). Cambridge University Press, Cambridge, 2010, pp. 435–484. [zbl](#) [MR](#)
- [10] *I. A. Stegun*: Legendre functions. Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables (M. Abramowitz, I. A. Stegun, eds.). National Bureau of Standards Applied Mathematics Series 55, Government Printing Office, Washington, 1970. [zbl](#) [MR](#)
- [11] *G. Strang*: *Introduction to Linear Algebra*. Wellesley-Cambridge Press, Wellesley, 2003. [zbl](#)

Authors' address: Abdulatif Badenjki, Gerald Warnecke, Institute for Analysis and Numerics, Otto-von-Guericke University, Universitaetsplatz 2, 39106 Magdeburg, Germany, e-mail: a.badenjki@hotmail.com, warnecke@ovgu.de.