

# Linear schemes with several degrees of freedom for the transport equation and the long-time simulation accuracy

P.A. Bakhvalov<sup>a,\*</sup>, M.D. Surnachev<sup>a</sup>

<sup>a</sup>*Keldysh Institute of Applied Mathematics, 4, Miusskaya Sq., Moscow, 125047, Russia*

---

## Abstract

We consider linear schemes with several degrees of freedom for the transport equation on uniform meshes. For these schemes the solution error is  $O(h^p + th^q)$ , where  $p$  is equal to or greater by one than the order of the truncation error and  $q \geq p$ . We prove the existence of a mapping of smooth functions on the mesh space providing the  $q$ -th order of the truncation error and deviating from the standard mapping ( $L_2$ -projection for example) by  $O(h^p)$ . In 1D case this mapping can be found in the class of local mappings. In more dimensions the existence of a local mapping with such properties is guaranteed only under additional assumptions.

**Keywords:** consistency and accuracy, superconvergence, long-time simulation accuracy, Fourier analysis

---

## 1. Introduction

The simplest approach to establish an estimate for the solution error of a numerical method is the analysis of the truncation error. A stable scheme possessing the truncation error of order  $P_A$  gives the solution with the error of the same order. But this estimate can be very far from optimal. For instance, the DG( $k$ ) method (discontinuous Galerkin method based on the polynomials of order  $k$ ) on the uniform meshes has the truncation error of order  $\max\{k, 1\}$ , while the numerical solution of the Cauchy problem for the model transport equation  $\partial v / \partial t + \partial v / \partial x = 0$  obtained by this scheme possesses an error estimate of the form

$$\|\varepsilon_h(t)\|_2 \leq C_1 h^{k+1} + C_2 h^{2k+1} t. \quad (1.1)$$

The comparison of the solution error by the 5-th order finite-difference scheme and the DG(4) method of the same order of accuracy is presented in Fig. 1. The estimate (1.1) follows directly from the results of [1] based on negative norm estimates which essentially use the finite-element nature of the DG scheme. We will consider another two methods that can be used to obtain such estimates.

The spectral analysis is a classical approach to study linear difference schemes (both semi-discrete and fully discrete) with constant coefficients on uniform meshes under periodic boundary conditions. In the case of 1 DOF per cell a wide range of its applications can be found in [2]. This method was successfully applied also for the schemes with several DOFs. For the DG method it was used in [3, 4, 5, 6], for the spectral difference method in [7, 8], for the flux reconstruction method (the class including the first two) in [9, 10], etc. Although the main application of the spectral analysis is to ensure stability, it can be applied for the accuracy analysis also. Lowrie [3] used it to understand the enhanced accuracy of DG( $k$ ) for  $k = 1, 2, 3$ , but details were not presented. In [5] the estimate (1.1) was proved for DG(1). In [6] this was proved for DG(2) and DG(3) using symbolic computations. To find the spectrum, one needs to find the roots of the polynomial of order  $m = k + 1$  with the coefficients depending on the wave number. If  $m > 4$ , it is not generally solvable by radicals. Thus the authors of [6] could not extend their analysis to  $k > 3$ . Note also the papers [11, 12] where the spectral analysis was combined with the finite-element technique to prove  $\|\varepsilon_h(t)\|_2 \leq C_1 h^{k+1} + C_2 h^{k+3/2} t$  for DG( $k$ ),  $k \geq 1$ , on non-uniform meshes.

---

\*Corresponding author

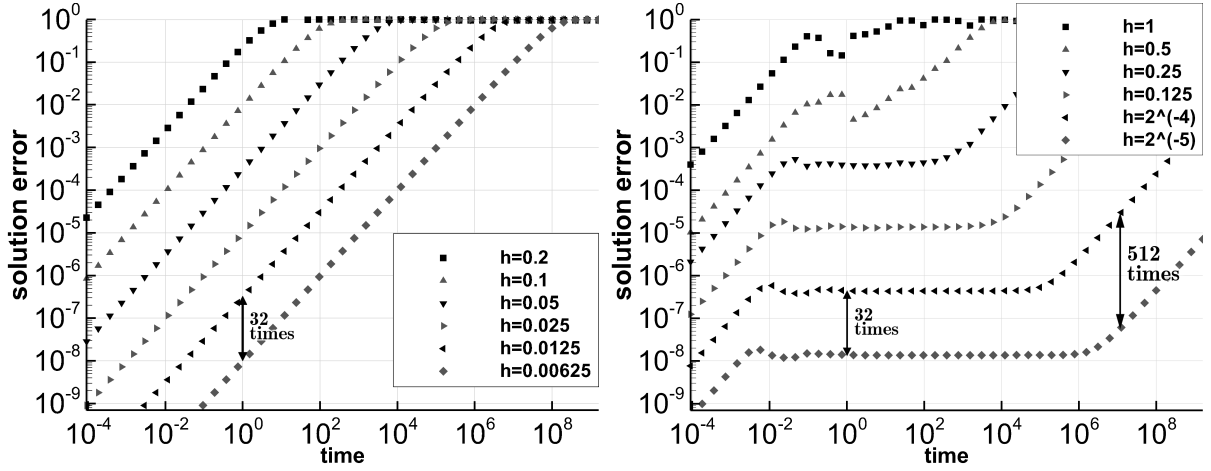


Figure 1: The  $L_2$  norm of the solution error of  $\partial v/\partial t + \partial v/\partial x = 0$ ,  $v_0 = \sin(2\pi x)$  obtained by the 5th order finite-difference scheme (left) and the DG(4) method (right)

The error estimate (1.1) for  $DG(k)$ ,  $k \in \mathbb{N}$ , was obtained by another approach which we call the method of auxiliary mapping. Let  $\Pi_h$  be an operator prescribing the initial data for the discrete (or semidiscrete) problem, i. e. taking each  $f$  to a mesh function  $f_h$ . In the case of  $DG(k)$ , a mesh function is a piecewise polynomial of order  $k$ , and the operator  $\Pi_h$  is the  $L_2$ -projection onto the space of mesh functions. The idea of the method is to introduce a new mapping  $\tilde{\Pi}_h$  such that for a sufficiently smooth function  $f$  there holds

$$(A1) \quad \|\tilde{\Pi}_h f - \Pi_h f\|_2 \leq C_0 h^P;$$

$$(A2) \quad \text{the scheme possesses the truncation error of order } Q \geq P \text{ in the sense of } \tilde{\Pi}_h.$$

By stability these properties guarantee (see Proposition 4.2) that in the sense of  $\Pi_h$  the solution error possesses an estimate of the form

$$\|\varepsilon_h(t)\|_2 \leq C_1 h^P + C_2 h^Q t, \quad Q \geq P > 0. \quad (1.2)$$

For the standard Galerkin method for elliptic and parabolic problems, the auxiliary mapping is usually chosen as the Ritz projection.

The first application of the method of auxiliary mapping to the DG scheme for the transport equation is also due to Lowrie for  $k = 1$  (see formulas (5.81) and (5.82) in [3]). In [13] the estimate (1.2) was proved for an arbitrary  $k$  and  $P = Q = k + 1$  using the Gauss–Radau projection. In [14] the estimate  $\|\varepsilon_h(t)\|_2 \leq C_1 h^{k+1} + C_2 h^{k+2} t^2$  was proved. Finally, the optimal result (1.1) was obtained in [15]. On a uniform mesh the auxiliary mapping  $\Pi_h^{(P,Q)}$  introduced in [15] takes each  $f \in W_{2,loc}^{Q+1}(\mathbb{R})$  to the mesh function  $\Pi_h^{(P,Q)} f$  defined on each cell  $\eta \in \mathbb{Z}$  by

$$\left(\Pi_h^{(P,Q)} f\right)_\eta = (\Pi_h f)_\eta + \sum_{m=P}^Q h^m \left( \mathcal{P}_h^{(m)} \left( \frac{d^m f}{dx^m} \right) \right)_\eta, \quad (1.3)$$

where  $(\Pi_h f)_\eta$  and  $(\Pi_h^{(P,Q)} f)_\eta$  are the polynomials at cell  $\eta$ , and  $\mathcal{P}_h^{(m)}$  are some linear mappings to the mesh space. The difference between the modified mapping and the original mapping (the last term in (1.3)) is called a correction function (in [15]) or a corrector (in [16]). Note that estimate (1.1) is valid for arbitrary non-uniform meshes (assuming  $h = h_{\max}$ ), but this is a specific feature of the DG method. Further results in the accuracy of the DG method were obtained in [17] (2D), [18, 19, 20] (variable coefficients), [21] (energy-conserving DG), Runge–Kutta time integration [22] etc.

Estimates of the form (1.2) with some  $P$  and  $Q$ , where  $Q$  is greater than the order of the truncation error, are valid not only for the DG method but also for some other schemes with several degrees of freedom per cell. A scheme on

a periodic mesh shows the same behavior if the mesh is refined by scaling (see Fig. 3). In this case we assume that if the mesh is scaled by a factor then  $h$  is multiplied by the same factor.

Consider an  $L_2$ -stable scheme with a local mapping  $\Pi_h$  (for the definition of “local” see Section 2.5). The existence of a mapping  $\tilde{\Pi}_h$  satisfying (A1) and (A2) for some  $P$  and  $Q$  proves the estimate (1.2). We are concerned with the following related problems.

1. How to get the optimal values of  $P$  and  $Q$  in the estimate (1.2)?
2. Does a mapping  $\tilde{\Pi}_h$  satisfying (A1) and (A2) for these  $P$  and  $Q$  exist?
3. How to construct it?

The main results of this paper are presented in three theorems. Theorem 1 gives the optimal value of  $P$ . Theorem 2 claims that the auxiliary mapping  $\tilde{\Pi}_h$  does exist. Theorem 3 states that in 1D and some multidimensional cases it can be found in a form similar to (1.3), that is the original mapping is modified using a local corrector. This leads to an algorithm which reduces the problem of finding the optimal values of  $P$  and  $Q$  in the estimate (1.2) to solving a linear system. However, in the general multidimensional case the situation is more tricky: generally the required corrector is nonlocal and the optimal value of  $Q$  can be non-integer.

The rest of the paper is organized as follows. In Section 2 we present a mathematical formulation of the problem. In Section 3 we state the main results. In Section 4 we recall the Lax – Ryabenkii theorem and present the basics of the method of auxiliary mapping. Section 5 is auxiliary and contains technical results. In Section 6 we use the spectral analysis to reduce our problem to a finite-dimensional formulation. In Section 7 we consider 1D case. While the results of this section can be treated as partial cases of the results from the next section, the 1D structure allows us to use simpler methods and understand the phenomenon of the enhanced accuracy in the long-time simulation more clearly. In Section 8 for a scheme possessing an estimate of the form (1.2) we prove the existence of an auxiliary mapping such that (A1) and (A2) hold. In Section 9 we consider two cases when this mapping is local. In Section 10 for these cases we present algorithms to get the optimal values  $P$  and  $Q$  in the estimate (1.2). In Section 11 we demonstrate our method on 1D examples. In Sections 12 and 13 we construct counterexamples for the 2D transport equation. In particular, we show that a local mapping  $\Pi_h^{(P,Q)}$  such that (A1) and (A2) hold generally does not exist. In Conclusion we summarize primary results of the paper.

## 2. Problem formulation

### 2.1. Solution spaces

Let  $d \in \mathbb{N}$  be the space dimension. Let  $\mathbf{a}_j, j = 1, \dots, d$ , form a basis in  $\mathbb{R}^d$ . We consider  $h\mathbf{a}_j$  as lattice vectors of the mesh. Let  $T$  be the map of  $\mathbb{R}^d$  to  $\mathbb{R}^d$  that takes each  $\boldsymbol{\eta}$  to

$$T\boldsymbol{\eta} = \sum_{j=1}^d \eta_j \mathbf{a}_j. \quad (2.1)$$

If  $\boldsymbol{\eta} \in \mathbb{Z}^d$  then  $T\boldsymbol{\eta}$  is the offset of the block  $\boldsymbol{\eta}$  from the zero block. We have  $(T^*\mathbf{x})_j = \mathbf{x} \cdot \mathbf{a}_j$ . If  $\mathbf{a}_j$  coincide with the vectors of the standard basis,  $T$  is the identity operator.

We say that  $f \in L_{2,loc}(\mathbb{R}^d)$  is periodic with period  $N_0 \in \mathbb{N}$  if for all  $j = 1, \dots, d$  there holds  $f(\mathbf{r} + N_0\mathbf{a}_j) = f(\mathbf{r})$  for almost all  $\mathbf{r} \in \mathbb{R}^d$ . Unless specifically stated, we will consider complex-valued functions. Denote by  $L_{2,per}(\mathbb{R}^d)$  the linear space of periodic  $f \in L_{2,loc}(\mathbb{R}^d)$  equipped with the norm

$$\|f\|^2 = \frac{1}{|\square|} \int_{\square} |f(\mathbf{r})|^2 dV,$$

where  $\square$  is the parallelepiped generated by the vectors  $N_0\mathbf{a}_1, \dots, N_0\mathbf{a}_d$  and  $N_0$  is a period of  $f$ .

Denote  $H_{per}^q(\mathbb{R}^d) = L_{2,per}(\mathbb{R}^d) \cap W_{2,loc}^q(\mathbb{R}^d)$  for  $q \geq 0$  and  $C_{per}^q(\mathbb{R}^d) = L_{2,per}(\mathbb{R}^d) \cap C^q(\mathbb{R}^d)$  for  $q \in \mathbb{N} \cup \{0\}$ . For  $w \in H_{per}^q(\mathbb{R}^d)$  (or  $w \in C_{per}^q(\mathbb{R}^d)$ ),  $q \in \mathbb{N} \cup \{0\}$ ,  $r = 0, \dots, q$ , denote

$$\|\nabla^r w\|^2 = \sum_{|\mathbf{m}|=r} \frac{r!}{\mathbf{m}!} \|D^{\mathbf{m}} w\|^2, \quad \|\nabla^r w\|_{\infty}^2 = \sum_{|\mathbf{m}|=r} \frac{r!}{\mathbf{m}!} \max_{\mathbb{R}^d} |D^{\mathbf{m}} w|^2, \quad D^{\mathbf{m}} = \frac{\partial^{|\mathbf{m}|}}{\partial x_1^{m_1} \dots \partial x_d^{m_d}}. \quad (2.2)$$

Here  $\mathbf{m} = (m_1, \dots, m_d)$  is a multiindex:  $m_j \in \mathbb{N} \cup \{0\}$ ,  $|\mathbf{m}| = m_1 + \dots + m_d$ ,  $\mathbf{m}! = m_1! \dots m_d!$ . Obviously,  $\|\nabla^{\mathbf{r}} f\| \leq \|\nabla^{\mathbf{r}} f\|_{\infty}$ . Denote  $\mathbf{r}^{\mathbf{m}} = x_1^{m_1} \dots x_d^{m_d}$  for  $\mathbf{r} = (x_1, \dots, x_d)$ . The notation  $\mathbf{l} \leq \mathbf{m}$  means that for each  $j = 1, \dots, d$  there holds  $l_j \leq m_j$ .

Each  $w \in L_{2,per}(\mathbb{R}^d)$  with a period  $N_0$  can be represented by the Fourier series

$$w = \sum_{\mathbf{k} \in \mathbb{Z}^d} \tilde{w}_{\mathbf{k}} \exp \left( i \frac{2\pi}{N_0} \mathbf{k} \cdot T^{-1} \mathbf{r} \right), \quad (2.3)$$

converging in  $L_{2,per}(\mathbb{R}^d)$ . If we consider  $w$  as a function with period  $2N_0$ , the coefficients  $\tilde{w}_{\mathbf{k}}$  will change. To avoid this, we introduce

$$\mathcal{A} = \{\boldsymbol{\alpha} \in \mathbb{R}^d : T^* \boldsymbol{\alpha} / \pi \in \mathbb{Q}^d\}.$$

Denoting  $\boldsymbol{\alpha} = 2\pi(T^*)^{-1} \mathbf{k} / N_0$  we rewrite (2.3) as

$$w = \sum_{\boldsymbol{\alpha} \in \mathcal{A}} w_{\boldsymbol{\alpha}} \exp(i \boldsymbol{\alpha} \cdot \mathbf{r}). \quad (2.4)$$

Here and below in such sums we imply that there exists a common multiple of denominators of all  $T^* \boldsymbol{\alpha} / (2\pi)$  such that  $w_{\boldsymbol{\alpha}} \neq 0$ . For  $w \in L_{2,per}(\mathbb{R}^d)$ , Parseval's identity  $\|w\|^2 = \sum |w_{\boldsymbol{\alpha}}|^2$  is valid. The correspondence between  $w \in L_{2,per}(\mathbb{R}^d)$  and its Fourier coefficients  $w_{\boldsymbol{\alpha}}$  in (2.4) establishes a bijective isometry between  $L_{2,per}(\mathbb{R}^d)$  and a dense subspace of  $l_2(\mathcal{A}, \mathbb{C})$ . The set  $\mathcal{A}$  is the frequency domain of this Fourier transform.

For  $w \in H_{per}^q(\mathbb{R}^d)$ ,  $q \in \mathbb{R}$ ,  $q \geq 0$ , denote

$$\|\nabla^q w\|^2 = \sum_{\boldsymbol{\alpha} \in \mathcal{A}} |w_{\boldsymbol{\alpha}}|^2 |\boldsymbol{\alpha}|^{2q}. \quad (2.5)$$

For  $q \in \mathbb{N} \cup \{0\}$  the definitions (2.2) and (2.5) coincide (see Lemma A.1 in Appendix). We equip  $H_{per}^q(\mathbb{R}^d)$ ,  $q \in \mathbb{R}$ ,  $q \geq 0$  and  $C_{per}^q(\mathbb{R}^d)$ ,  $q \in \mathbb{N} \cup \{0\}$  with the families of norms

$$\|f\|_{(q,h)}^2 = \|f\|^2 + h^{2q} \|\nabla^q f\|^2, \quad \|f\|_{(q,h,\infty)}^2 = \|f\|_{\infty}^2 + h^{2q} \|\nabla^q f\|_{\infty}^2.$$

Let  $\Omega$  be the unit sphere in  $\mathbb{R}^d$ . Denote by  $\mathring{\Omega}$  a set of vectors  $\mathbf{e} \in \Omega$  such that  $\lambda T^* \mathbf{e} \in \mathbb{Q}^d$  for some  $\lambda > 0$ . We say that  $f \in L_{2,loc}(\mathbb{R}^d)$  has a direction  $\mathbf{e} \in \Omega$  and write  $f \in L_{2,loc,\mathbf{e}}(\mathbb{R}^d)$  if there exists a function  $g \in L_{2,loc}(\mathbb{R})$  such that  $f(\mathbf{r}) = g(\mathbf{r} \cdot \mathbf{e})$  for almost all  $\mathbf{r} \in \mathbb{R}^d$ . For  $\mathbf{e} \in \mathring{\Omega}$  we denote  $H_{per,\mathbf{e}}^q(\mathbb{R}^d) = L_{2,loc,\mathbf{e}}(\mathbb{R}^d) \cap H_{per}^q(\mathbb{R}^d)$ ,  $C_{per,\mathbf{e}}^q(\mathbb{R}^d) = L_{2,loc,\mathbf{e}}(\mathbb{R}^d) \cap C_{per}^q(\mathbb{R}^d)$ . The spaces  $H_{per,\mathbf{e}}^q(\mathbb{R}^d)$  and  $C_{per,\mathbf{e}}^q(\mathbb{R}^d)$  contain nonconstant functions (see Lemma A.2 in Appendix).

## 2.2. Mesh function spaces

Let  $M^0$  be a finite set of the degrees of freedom (DOFs) at a mesh block and  $M = \mathbb{Z}^d \times M^0$  be the general set of DOFs. For  $f \in \mathbb{C}^M$  by  $f_{\boldsymbol{\eta}} \in \mathbb{C}^{M^0}$  we denote a part of the vector  $f$  in the block  $\boldsymbol{\eta} \in \mathbb{Z}^d$ . We shall also call  $f_{\boldsymbol{\eta}}$  the block component of  $f$  at  $\boldsymbol{\eta}$ . The set of sequences with period  $N$  is

$$V_{per}^N = \{f \in \mathbb{C}^M : \forall \boldsymbol{\eta}, \boldsymbol{\zeta} \in \mathbb{Z}^d \quad f_{\boldsymbol{\eta}+N\boldsymbol{\zeta}} = f_{\boldsymbol{\eta}}\}.$$

The set of periodic sequences  $V_{per} = \bigcup_{N \in \mathbb{N}} V_{per}^N$  is equipped with the scalar product defined for  $f \in V_{per}^{N(f)}$ ,  $g \in V_{per}^{N(g)}$  as

$$(f, g) = \frac{1}{N^d} \sum_{\boldsymbol{\eta}=(0,\dots,N-1)^d} (f_{\boldsymbol{\eta}}, g_{\boldsymbol{\eta}}), \quad N = N(f)N(g).$$

Here  $(f_{\boldsymbol{\eta}}, g_{\boldsymbol{\eta}})$  is any scalar product on  $\mathbb{C}^{M^0}$ . If the function  $f$  has a period  $N(f)$  then it has the period  $nN(f)$  for all  $n \in \mathbb{N}$ , but the substitution  $nN(f)$  for  $N(f)$  does not change the value of the scalar product. We equip  $\mathbb{C}^{M^0}$  and  $V_{per}$  with the norms induced by these scalar products:  $\|f_{\boldsymbol{\eta}}\|^2 = (f_{\boldsymbol{\eta}}, f_{\boldsymbol{\eta}})$ ;  $\|f\|^2 = (f, f)$ . The space  $V_{per}$  is incomplete (see Corollary 6.3).

### 2.3. The equation and the schemes

In this paper we consider the initial value problem for the model linear transport equation

$$\frac{\partial v}{\partial t} + \boldsymbol{\omega} \cdot \nabla v = 0, \quad \mathbf{r} \in \mathbb{R}^d, \quad (2.6)$$

$$v(0, \mathbf{r}) = v_0(\mathbf{r}) \in L_{2,per}(\mathbb{R}^d). \quad (2.7)$$

The transport velocity  $\boldsymbol{\omega}$  is constant in space and time.

To approximate (2.6) consider semi-discrete schemes of the form

$$\sum_{\zeta \in \mathcal{S}} Z_{\zeta} \frac{du_{\boldsymbol{\eta}+\zeta}}{dt}(t) + \frac{1}{h} \sum_{\zeta \in \mathcal{S}} L_{\zeta} u_{\boldsymbol{\eta}+\zeta}(t) = 0, \quad \boldsymbol{\eta} \in \mathbb{Z}^d, \quad u_{\boldsymbol{\eta}} \in \mathbb{C}^{M^0}, \quad (2.8)$$

where  $\mathcal{S} \subset \mathbb{Z}^d$  is a finite set (the stencil of the scheme),  $Z_{\zeta}$  and  $L_{\zeta}$  are real-valued matrices. For  $\zeta \notin \mathcal{S}$  put  $Z_{\zeta} = L_{\zeta} = 0$ . Let  $Z : \mathbb{C}^M \rightarrow \mathbb{C}^M$  and  $L : \mathbb{C}^M \rightarrow \mathbb{C}^M$  take each  $u$  to  $Zu$  and  $Lu$  such that

$$(Zu)_{\boldsymbol{\eta}} = \sum_{\zeta \in \mathcal{S}} Z_{\zeta} u_{\boldsymbol{\eta}+\zeta}, \quad (Lu)_{\boldsymbol{\eta}} = \sum_{\zeta \in \mathcal{S}} L_{\zeta} u_{\boldsymbol{\eta}+\zeta}. \quad (2.9)$$

With this notation, (2.8) is equivalent to

$$Z \frac{du}{dt} + \frac{1}{h} Lu = 0. \quad (2.10)$$

Obviously,  $Z V_{per}^N \subseteq V_{per}^N$ ,  $L V_{per}^N \subseteq V_{per}^N$ , hence,  $Z V_{per} \subseteq V_{per}$ ,  $L V_{per} \subseteq V_{per}$ . The operators  $L$  and  $Z$  are bounded on  $V_{per}$  with

$$\|Z\| \leq \sum_{\zeta \in \mathcal{S}} \|Z_{\zeta}\|, \quad \|L\| \leq \sum_{\zeta \in \mathcal{S}} \|L_{\zeta}\|,$$

where the operator norms on the LHS are induced by the norm on  $V_{per}$  and the operator norms on the RHS are induced by the norm on  $\mathbb{C}^{M^0}$ .

**Lemma 2.1.** *The operator  $Z : V_{per} \rightarrow V_{per}$  has a bounded inverse if and only if for each  $N$  there exists an inverse for the restriction of  $Z$  to  $V_{per}^N$ , with its norm being limited by a constant independent of  $N$ .*

*Proof.* Let  $Z$  be invertible and  $\|Z^{-1}\| < \infty$ . Then  $\text{Ker } Z = \{0\}$ . In particular, for each  $N$  the restriction of  $Z$  to  $V_{per}^N$  has zero kernel, and is invertible. Then  $Z^{-1} V_{per}^N = V_{per}^N$ . Clearly, the norm of the restriction of  $Z^{-1}$  to  $V_{per}^N$  does not exceed the norm of  $Z^{-1}$  on  $V_{per}$ . To prove the reverse implication note that the norm of  $Z^{-1}$  on  $V_{per}$  is the supremum of the norms of its restrictions to  $V_{per}^N$ .  $\square$

Further in this paper we assume that for  $Z : V_{per} \rightarrow V_{per}$  there exists a bounded inverse.

**Lemma 2.2.** *For any initial data  $u_0 \in V_{per}$  the ODE system (2.10) has a unique solution  $u \in C^\infty([0, \infty), V_{per})$  satisfying  $u(0) = u_0$ . Besides, if  $u_0 \in V_{per}^N$ , then for any  $t > 0$  we have  $u(t) \in V_{per}^N$ .*

*Proof.* Let  $u_0 \in V_{per}^N$ . Since  $Z$  is invertible on  $V_{per}^N$  and  $Z^{-1} L V_{per}^N \subset V_{per}^N$ , (2.10) has a solution  $u \in C^\infty([0, \infty), V_{per}^N)$ . The uniqueness on  $V_{per}$  follows from the boundedness of  $Z^{-1} L$  on  $V_{per}$ .  $\square$

Since the operator  $Z^{-1} L$  is bounded on  $V_{per}$ , and  $V_{per}^N$  are its invariant subspaces, it generates the uniformly continuous operator group  $\exp(-z Z^{-1} L) = \sum_{k=0}^{\infty} (-z Z^{-1} L)^k / k!$ ,  $z \in \mathbb{C}$ , which is analytical on the whole complex plane and also has  $V_{per}^N$  as its invariant subspaces. Since  $V_{per}^N$  are finite dimensional, on these subspaces  $\exp(-Z^{-1} L z)$  can be represented as the standard matrix exponential. For  $u_0 \in V_{per}$  the function  $u(t) = \exp(-Z^{-1} L t / h) u_0$  solves (2.10) with the initial data  $u(0) = u_0$ .

A scheme (2.8) is called *stable* on  $V_{per}$  with the stability constant  $K$ , if for all  $h$  and each  $u \in C^\infty([0, \infty), V_{per})$  satisfying (2.8) and each  $t \geq 0$  there holds  $\|u(t)\| \leq K \|u(0)\|$ . In other words, a scheme (2.8) is stable with stability

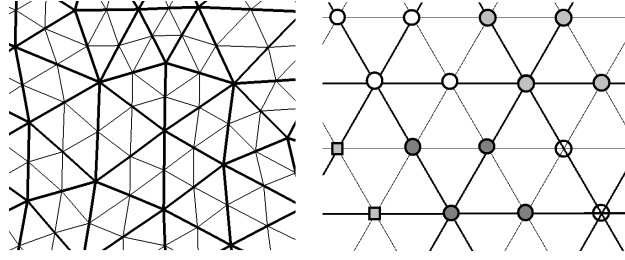


Figure 2: Fragments of a mesh with spectral elements. Left: unstructured mesh. Right: translationally-invariant mesh. In the right figure, nodes from different mesh blocks are marked by different symbols

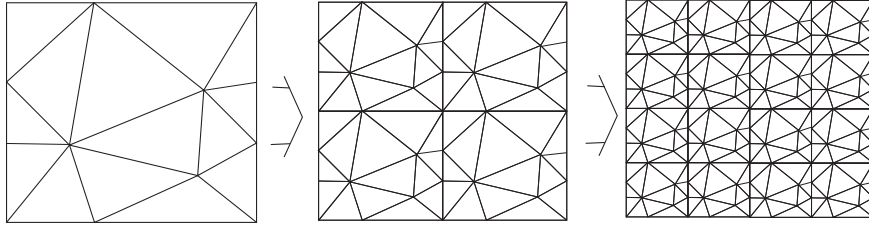


Figure 3: Block refinement of a mesh

105 constant  $K$  iff  $\sup_{t,h>0} \|\exp(-Z^{-1}Lt/h)\| = \sup_{\nu \geq 0} \|\exp(-Z^{-1}L\nu)\| \leq K$ .

The most simple schemes of the form (2.8) are the finite-difference schemes on the uniform meshes. In this case  $|M^0| = 1$ , the order of accuracy coincides with the order of the truncation error and there is no enhanced accuracy in the long-time simulation (see Proposition 9.9). The schemes of the form (2.8) such that  $|M^0| > 1$  arise in the following situations.

- Schemes with several DOFs per cell on the uniform meshes: discontinuous Galerkin, spectral difference, active flux etc. In this case  $\mathbf{a}_j$  are the linearly independent vectors of the mesh edges and  $M^0$  is the set of the DOFs in a cell. We assume that the ratio of the mesh steps along different directions remains constant under mesh refinement.
- 115 • The use of spectral elements. For example, in the flux correction method [24] one needs to compute the gradients at nodes of an unstructured mesh with the second order of the truncation error. It is convenient to use an unstructured mesh obtained by refining a base mesh (see Fig. 2, left). Elements of the base mesh are called spectral elements. On each spectral element one can construct the 2-nd order interpolation polynomial based on the nodal values of the mesh function. A gradient of the mesh function at a node can be defined as a weighted average of the gradients of interpolation polynomials at all the spectral elements containing this node. The resulting scheme on the translationally-invariant meshes has the form (2.8) with  $|M^0| = 4$  (see Fig. 2, right) since the gradients are defined differently in the nodes on the base mesh and in the nodes at edges of the spectral elements. We say that the mesh is *translationally-invariant* if it is invariant with respect to the translation by the vector of any mesh edge.
- 125 • Finite-difference or finite-volume schemes if the block refinement is used (see Fig. 3). In this case the whole mesh is uniform-block, i. e. space is tessellated by blocks of an unstructured mesh. Blocks can be indexed by  $\eta_1, \dots, \eta_d$ , as if they were cells of the uniform mesh. Vectors  $\mathbf{a}_j, j = 1, \dots, d$  are the offsets of blocks adjacent to a reference block and  $|M^0|$  is the number of DOFs per block.
- 130 • Combination of the cases mentioned above. For example, the DG method on a simplicial translationally-invariant mesh.

#### 2.4. Homogeneous mappings to the mesh space

Throughout this paper we assume that  $1/h \in \mathbb{N}$ . In order to transform data to the mesh space (prescribe the initial data for the difference problem etc.) one needs a mapping. Let  $\mathcal{H} = H_{per}^q(\mathbb{R}^d)$  or  $\mathcal{H} = C_{per}^q(\mathbb{R}^d)$ . Consider a family of linear mappings  $\Pi_h : \mathcal{H} \rightarrow V_{per}$ ,  $1/h \in \mathbb{N}$ , with the following properties:

- (P1) for each  $f, g \in \mathcal{H}$  such that almost everywhere  $g(\mathbf{r}) = f(\mathbf{r} + hT\zeta)$  there holds  $(\Pi_h f)_{\eta+\zeta, \xi} = (\Pi_h g)_{\eta, \xi}$  for each  $\eta, \zeta \in \mathbb{Z}^d, \xi \in M^0$ ;
- (P2) for each  $f, g \in \mathcal{H}$  such that almost everywhere  $g(\mathbf{r}) = f(\mathbf{r}/h)$  there holds  $(\Pi_h g)_{\eta, \xi} = (\Pi_1 f)_{\eta, \xi}$ ;
- (P3) if  $f \in \mathcal{H}$  is real-valued then so is  $\Pi_h f$ .

Each  $\Pi_h$  from this family we call a *homogeneous mapping*. By (P2), a homogeneous mapping  $\Pi_h$  uniquely defines the family it belongs to.

A homogeneous mapping  $\Pi_h$  maps a function  $f \in \mathcal{H}$  with a period  $N_0$  to a sequence  $\Pi_h f \in V_{per}$  with the period  $N_0/h$ . This follows directly from the definition. A homogeneous mapping  $\Pi_h$  of  $H_{per}^r(\mathbb{R}^d)$  to  $V_{per}$  is *bounded* if  $\|\Pi_h\|_{(r, h)} := \sup \|\Pi_h f\| / \|f\|_{(r, h)} < \infty$ . A homogeneous mapping  $\Pi_h$  of  $C_{per}^r(\mathbb{R}^d)$  to  $V_{per}$  is *bounded* if  $\|\Pi_h\|_{(r, h, \infty)} := \sup \|\Pi_h f\| / \|f\|_{(r, h, \infty)} < \infty$ . By construction,  $\|\Pi_h\|_{(r, h)}$  and  $\|\Pi_h\|_{(r, h, \infty)}$  do not depend on  $h$ .

#### 2.5. Local mappings to the mesh space

Consider a family of mappings  $\Pi_h$ ,  $1/h \in \mathbb{N}$ , from  $L_{2, loc}(\mathbb{R}^d)$  to  $\mathbb{C}^M$  of the form

$$(\Pi_h f)_{\eta, \xi} = \langle \mu_\xi, f(h(\cdot + T\eta)) \rangle \quad (2.11)$$

where  $\eta \in \mathbb{Z}^d$  is the index of a mesh block,  $\xi \in M^0$  is the index of a variable inside the block,  $\mu_\xi \in (W_2^q(G))^*$  or  $(C^q(G))^*$  for some  $q \in \mathbb{N} \cup \{0\}$  and a bounded domain  $G$ ,  $\Pi_h f$  is real-valued for real-valued  $f$ , and  $\langle \mu_\xi, 1 \rangle \neq 0$  for at least one  $\xi \in M^0$ . Each  $\Pi_h$  from this family we call a *local mapping* and  $\mu_\xi$  we call its kernel. Particularly, if  $\mu_\xi$  is a measure then

$$(\Pi_h f)_{\eta, \xi} = \int_G f(h(\mathbf{r} + T\eta)) d\mu_\xi = \int_G f\left(h\mathbf{r} + h \sum_{j=1}^d \eta_j \mathbf{a}_j\right) d\mu_\xi.$$

If  $\mu_\xi \in L_2(G)$ , then  $\Pi_h$  maps  $L_{2, loc}(\mathbb{R}^d)$  to  $\mathbb{C}^M$ . For example, the standard  $L_2$ -projection onto a space of functions, which are polynomials on each mesh cell, can be represented as a local mapping. If  $\mu_\xi \in (W_2^q(G))^*$  for some  $q \in \mathbb{N} \cup \{0\}$ , then  $\Pi_h$  maps  $H_{loc}^q(\mathbb{R}^d)$  to  $\mathbb{C}^M$ . If  $\mu_\xi \in (C^q(G))^*$  for some  $q \in \mathbb{N} \cup \{0\}$ , then  $\Pi_h$  maps  $C^q(\mathbb{R}^d)$  to  $\mathbb{C}^M$ . An example is the pointwise mapping used in finite-difference schemes, namely,  $(\Pi_h f)_{\eta, \xi} = f(h(\rho_\xi + T\eta))$  (i. e.  $\mu_\xi = \delta(\mathbf{r} - \rho_\xi)$ ) where  $\rho_\xi$  are collocation points.

**Lemma 2.3.** *A local mapping  $\Pi_h$  with  $\mu_\xi \in (W_2^q(G))^*$  is a bounded homogeneous mapping of  $H_{per}^q(\mathbb{R}^d)$  to  $V_{per}$ . A local mapping  $\Pi_h$  with  $\mu_\xi \in (C^q(G))^*$  is a bounded homogeneous mapping of  $C_{per}^q(\mathbb{R}^d)$  to  $V_{per}$ .*

*Proof.* Let  $N_0$  be a period of  $f$ . Then for each  $\eta, \zeta \in \mathbb{Z}^d, \xi \in M^0$  there holds

$$(\Pi_h f)_{\eta+\zeta N_0/h, \xi} = \langle \mu_\xi, f(h(\cdot + T\eta) + N_0 T\zeta) \rangle = \langle \mu_\xi, f(h(\cdot + T\eta)) \rangle = (\Pi_h f)_{\eta, \xi},$$

and the same for  $\Pi_h'$ . Thus  $\Pi_h H_{per}^q(\mathbb{R}^d) \subseteq V_{per}$ ,  $\Pi_h' C_{per}^q(\mathbb{R}^d) \subseteq V_{per}$ . The properties (P1) and (P2) are obvious. The boundedness of  $\Pi_h'$  is obvious; the proof of the boundedness of  $\Pi_h$  is in Appendix (see Proposition A.4).  $\square$

In this paper we will use mappings  $\Pi_h^{(p, q)}$  and  $\Pi_{h, \mathbf{e}}^{(p, q)}$ ,  $p, q \in \mathbb{N}$ , given by

$$\left(\Pi_h^{(p, q)} f\right)_\eta = (\Pi_h f)_\eta + \sum_{p \leq |\mathbf{m}| \leq q} h^{|\mathbf{m}|} \mathfrak{C}^{(\mathbf{m})} (\mathcal{P}_h D^{\mathbf{m}} f)_\eta, \quad (2.12)$$

$$\left(\Pi_{h, \mathbf{e}}^{(p, q)} f\right)_\eta = (\Pi_h f)_\eta + \sum_{m=p}^q h^m \mathfrak{C}_e^{(m)} \left(\mathcal{P}_h \frac{\partial^m f}{\partial \mathbf{e}^m}\right)_\eta, \quad (2.13)$$

where  $\Pi_h$  and  $\mathcal{P}_h$  are some local mappings, the kernel  $\hat{\mu}_\xi$  of  $\mathcal{P}_h$  satisfies  $\langle \hat{\mu}_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ ,  $\mathfrak{C}_e^{(m)}$  and  $\mathfrak{C}_e^{(m)}$  are diagonal square real-valued matrices of size  $|M^0|$ , and  $e$  is a unit vector. It is easy to see that  $\Pi_h^{(p,q)}$  and  $\Pi_{h,e}^{(p,q)}$  are local mappings. If the kernels of  $\Pi_h$  and  $\mathcal{P}_h$  belong to  $L_2(G)$  then  $\Pi_h^{(p,q)}$  maps  $W_{2,loc}^q(\mathbb{R}^d)$  to  $\mathbb{C}^M$  and  $H_{per}^q(\mathbb{R}^d)$  to  $V_{per}$ . If the kernels of  $\Pi_h$  and  $\mathcal{P}_h$  belong to  $(C(G))^*$  then  $\Pi_h^{(p,q)}$  maps  $C^q(\mathbb{R}^d)$  to  $\mathbb{C}^M$  and  $C_{per}^q(\mathbb{R}^d)$  to  $V_{per}$ .

## 155 2.6. The truncation error and the solution error

Suppose  $\Pi$  maps a subspace of  $H_{loc}^1(\mathbb{R}^d)$  to  $\mathbb{C}^M$ . The *truncation error* on the function  $f$  in the sense of  $\Pi$  is the vector  $\epsilon_h(f, \Pi) \in \mathbb{C}^M$  defined as

$$\epsilon_h(f, \Pi) = -Z\Pi(\omega \cdot \nabla f) + \frac{1}{h}L\Pi f. \quad (2.14)$$

Let  $\Pi_h$  be a homogeneous mapping of  $f \in H_{per}^q(\mathbb{R}^d)$  (or  $C_{per}^q(\mathbb{R}^d)$ ) to  $V_{per}$ . Then for each  $f$  with a period  $N_0$  and each  $h$  there holds  $\epsilon_h(f, \Pi_h) \in V_{per}^{(N_0/h)}$ .

**Definition 1.** We say that the scheme possesses the *truncation error of order  $P_A$*  in the sense of  $\Pi_h$  if for each  $\alpha \in \mathcal{A}$  there holds  $\|\epsilon_h(e^{i\alpha \cdot r}, \Pi_h)\| \leq C(\alpha)h^{P_A}$ . We say that the scheme possesses the truncation error of order  $P_A$  in the sense of  $\Pi_h$  on the direction  $e \in \Omega$  if for each  $\alpha \in \mathbb{R}$  such that  $\alpha e \in \mathcal{A}$  there holds  $\|\epsilon_h(e^{i\alpha e \cdot r}, \Pi_h)\| \leq C(\alpha)h^{P_A}$ . We say that the optimal order of the truncation error is  $P_A$  if the scheme possesses the truncation error of order  $P_A$  and does not possess the truncation error of order  $P_A + \delta$  for each  $\delta > 0$ .

Later (see Corollary 6.21) we will show that for a local mapping  $\Pi_h$  with kernel  $\mu \in (W_2^q(G))^*$  the scheme possesses the truncation error of order  $P_A \in \mathbb{N}$  iff for some  $c_1, c_2 > 0$  there holds

$$\|\epsilon_h(v_0, \Pi_h)\| \leq c_1 \|\nabla^{P_A+1} v_0\| h^{P_A} + c_2 \|\nabla^{\max\{P_A, q\}+1} v_0\| h^{\max\{P_A, q\}},$$

for each  $v_0 \in H_{per}^{\max\{P_A, q\}+1}(\mathbb{R}^d)$ . A similar property holds if  $\mu \in (C^q(G))^*$ .

**Definition 2.** Suppose  $\Pi$  maps  $H_{per}^r(\mathbb{R}^d)$  or  $C_{per}^r(\mathbb{R}^d)$  to  $V_{per}$  for some  $r$ . The *solution error* in the sense of  $\Pi$  with the initial data  $v_0$  is the vector

$$\varepsilon_h(t, v_0, \Pi) = u(t) - \Pi v(t, \cdot) \in V_{per}, \quad (2.15)$$

where  $u(t) = \exp(-Z^{-1}Lt/h)\Pi v_0$  is the solution of (2.10) with the initial data  $u(0) = \Pi v_0$ , and  $v(t, r) = v_0(r - \omega t)$ .

If  $\Pi_h$  is a local mapping with  $\mu_\xi \in (W_2^q(G))^*$  (or  $(C^q(G))^*$ ), then for each  $f \in H_{per}^q(\mathbb{R}^d)$  (or  $C_{per}^q(\mathbb{R}^d)$ ) with a period  $N_0$  and each  $h$  and  $t \geq 0$  there holds  $\varepsilon_h(t, f, \Pi_h) \in V_{per}^{(N_0/h)}$ . Note that by definition we have  $\varepsilon_h(0, v_0, \Pi) \equiv 0$ , so in contrast to the concepts conventional in the finite-element analysis we do not consider the error of the initial data mapping to be a part of the solution error.

## 170 2.7. Long-time simulation accuracy

We introduce two different definitions of the scheme order in the long-time simulation.

**Definition 3.** Consider a scheme of the form (2.8) and a homogeneous mapping  $\Pi_h$  of  $\mathcal{H} \subseteq L_{2,per}(\mathbb{R}^d)$  to  $V_{per}$ . Let  $P, Q$  satisfy  $0 < P \leq Q \leq \infty$ . Suppose for each  $v_0 \in \mathcal{H}$  there exist non-negative constants  $c_1(v_0), c_2(v_0)$  such that for each  $h$ , and each  $t \geq 0$  the scheme possesses the error estimate

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq c_1(v_0)h^P + c_2(v_0)th^Q \quad (2.16)$$

(here we assume  $h^\infty = 0$ ). Then we say that in the sense of  $\Pi_h$  on  $\mathcal{H}$  the scheme possesses the *formal order of accuracy  $P$*  and the *long-time simulation order  $Q$  in the weak sense*.



**Definition 4.** Let  $\mathcal{H}$  be one of the following spaces:  $H_{per}^R(\mathbb{R}^d)$  for some  $R \geq 0$ ,  $C_{per}^R(\mathbb{R}^d)$  for some  $R \in \mathbb{N} \cup \{0\}$ ,  $H_{per,e}^R(\mathbb{R}^d)$  for some  $R \geq 0$  and  $e \in \hat{\Omega}$ ,  $C_{per,e}^R(\mathbb{R}^d)$  for some  $R \in \mathbb{N} \cup \{0\}$  and  $e \in \hat{\Omega}$ . Consider a scheme of the form (2.8) and a homogeneous mapping  $\Pi_h$  with  $\mathcal{H}$  in the domain of  $\Pi_h$ .

**i.** Let  $P$  and  $Q$  be real numbers satisfying  $0 < P \leq Q < \infty$  and  $Q + 1 \leq R$ . We say that the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  in the sense of  $\Pi_h$  on  $\mathcal{H}$  if for each initial data  $v_0 \in \mathcal{H}$ , each  $h$ , and each  $t \geq 0$  the scheme possesses the error estimate

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq C_1 h^P \|\nabla^P v_0\| + C_2 t h^Q \|\nabla^{Q+1} v_0\| + C_3 h^R \|\nabla^R v_0\|_* \quad (2.17)$$

where  $C_1, C_2, C_3$  are non-negative constants independent of  $v_0, h$ , and  $t$ . Here  $*$   $\equiv \infty$  for  $\mathcal{H} = C_{per}^R(\mathbb{R}^d)$  or  $C_{per,e}^R(\mathbb{R}^d)$  and dropped otherwise<sup>1</sup>.

**ii.** Let  $P$  be a real number satisfying  $0 < P \leq R$ . We say that the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q = \infty$  in the sense of  $\Pi_h$  on  $\mathcal{H}$  if for each initial data  $v_0 \in \mathcal{H}$ , each  $h$ , and each  $t \geq 0$  the scheme possesses the error estimate

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq C_1 h^P \|\nabla^P v_0\| + C_3 h^R \|\nabla^R v_0\|_* \quad (2.18)$$

with nonnegative constants  $C_1$  and  $C_3$  independent of  $v_0, h$ , and  $t$ , and  $*$  treated as in (2.17).

**iii.** If for each initial data  $v_0 \in \mathcal{H}$  there holds  $\varepsilon_h(t, v_0, \Pi_h) \equiv 0$  then we say that the scheme possesses the formal order of accuracy  $P = \infty$  and the long-time simulation order  $Q = \infty$  on  $\mathcal{H}$ .

Obviously, if the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  then it possesses the same orders in the weak sense. Unless specifically stated, speaking about the formal order of accuracy and the long-time simulation order we imply the use of Definition 4.

For  $y \in \mathbb{R}$  denote by  $\lfloor y \rfloor$  the largest integer that is less than or equal to  $y$  and by  $\lceil y \rceil$  the lowest integer that is greater than or equal to  $y$ . The formal order of accuracy and the long-time simulation order given by Definition 4 have the following properties, which are intuitively clear.

**Lemma 2.4.** Let  $\Pi_h$  be a local mapping with  $\mu_\xi \in (W_2^s(G))^*$  (or  $(C^s(G))^*$ ). Let the scheme (2.8) be stable and possess the formal order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per}^r(\mathbb{R}^d)$  (or  $C_{per}^{\lceil r \rceil}(\mathbb{R}^d)$ ),  $r \geq \max\{Q + 1, s\}$  (or  $r \geq \max\{P, s\}$  if  $Q = \infty$ ). Then for each  $p$  and  $q$  such that  $0 < p \leq P$  and  $p \leq q \leq Q$  the scheme possesses the formal order  $p$  and the long-time simulation order  $q$  on  $H_{per}^r(\mathbb{R}^d)$  (or  $C_{per}^{\lceil r \rceil}(\mathbb{R}^d)$ ). If  $q < Q = \infty$  we additionally assume  $q \leq r - 1$ . The same holds for the definition in the weak sense.

In the sense of Definition 3 this is obvious; in the sense of Definition 4 this will be proved below (see Corollary 6.18). So we can say that the scheme possesses the order of accuracy  $P$ ,  $0 < P \leq \infty$ , if the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $P$  on  $H_{per}^r(\mathbb{R}^d)$  (or  $C_{per}^{\lceil r \rceil}(\mathbb{R}^d)$ ) with  $r \geq P$ .

**Lemma 2.5.** Let  $P > 0$ ,  $Q \geq P$ ,  $R \geq Q + 1$ . Let  $\Pi_h$  be a homogeneous mapping of  $H_{per}^R(\mathbb{R}^d)$ . Let  $r \geq R$ . The scheme (2.8) possesses the formal order  $P$  and the long-time simulation order  $Q$  on  $H_{per}^r(\mathbb{R}^d)$  iff it possesses these orders on  $H_{per}^r(\mathbb{R}^d)$ .

**Lemma 2.6.** Let  $P > 0$ ,  $Q \geq P$ ,  $R \geq Q + 1$ ,  $R \in \mathbb{N}$ . Let  $\Pi_h$  be a local mapping with  $\mu_\xi \in (W_2^R(G))^*$ . The scheme (2.8) possesses the formal order  $P$  and the long-time simulation order  $Q$  on  $H_{per}^R(\mathbb{R}^d)$  iff it possesses the formal order  $P$  and the long-time simulation order  $Q$  on  $C_{per}^R(\mathbb{R}^d)$ .

These lemmas follow from Theorem 6.17. Thus we will drop the class of functions on which the scheme possesses these orders, unless this leads to ambiguity.

Definitions 3 and 4 introduce the concepts of the formal order and the long-time simulation order simultaneously. Now we introduce the concept of the optimal values in such a way that one should first define the optimal value of the formal order of accuracy and then find the optimal value of the long-time simulation order.

<sup>1</sup> If  $P, Q \in \mathbb{N}$ , then  $\|\nabla^P v_0\|$  and  $\|\nabla^{Q+1} v_0\|$  can be replaced by  $\|\nabla^P v_0\|_*$  and  $\|\nabla^{Q+1} v_0\|_*$ , which leads to an equivalent definition.

**Definition 5.** Let  $\Pi_h$  be a homogeneous mapping. Let the scheme (2.8) possess the formal order  $P$  and the long-time simulation order  $Q$  in the sense of  $\Pi_h$  on  $H_{per}^R(\mathbb{R}^d)$  (or  $H_{per,e}^R(\mathbb{R}^d)$ ). If for each  $P' > P$ ,  $Q' \geq P'$  and for each  $P' = P$ ,  $Q' > Q$  the scheme does not possess the formal order  $P'$  and the long-time simulation order  $Q'$  in the sense of  $\Pi_h$  on any  $H_{per}^r(\mathbb{R}^d)$  (or  $H_{per,e}^r(\mathbb{R}^d)$ ), then we call the values  $P$  and  $Q$  *optimal*. The optimal values for homogeneous mappings of  $C_{per}^q(\mathbb{R}^d)$  and the optimal values in the weak sense are defined in the same way.

### 3. Main results

For  $n \in \mathbb{N}$ , let  $\mathring{\Omega}_n$  be a set of vectors  $\{e_k \in \mathring{\Omega}, k = 1, \dots, C_{n+d-1}^{d-1}\}$  such that  $\{(e_k \cdot r)^n\}$  form a basis in the space of homogeneous polynomials of order  $n$ . For the existence of this set see Lemma A.6. Denote  $L(0) = \sum_{\zeta} L_{\zeta}$ .

**Theorem 1.** Consider a scheme of the form (2.8), stable with a constant  $K$ , and local mappings  $\Pi_h, \mathcal{P}_h$  with kernels  $\mu_{\xi}, \hat{\mu}_{\xi} \in (W_2^q(G))^*$  (or  $(C^q(G))^*$ ). Let  $P_A$  and  $P$  be the optimal orders of the truncation error and accuracy, correspondingly, in the sense of  $\Pi_h$ . Let  $\hat{\kappa} = \min_{\xi \in M^0} |\langle \hat{\mu}_{\xi}, 1 \rangle| > 0$  and  $R = \max\{q, P\} + 1$ . Then the following holds.

1.  $P_A$  and  $P$  are integers.
2. Either  $P = P_A$  or  $P = P_A + 1$ .
3. If  $P = P_A + 1$  then there exist real-valued diagonal matrices  $\mathfrak{C}^{(m)}, |\mathbf{m}| = P$ , such that the scheme possesses the truncation error of order  $P$  in the sense of  $\Pi_h^{(P,P)}$  given by (2.12). Moreover,  $\|\mathfrak{C}^{(m)}\| \leq \tilde{\delta} C_1$ , where  $C_1$  is the constant in estimate (2.17) and  $\tilde{\delta}$  depends only on  $K, \hat{\kappa}, P, |M^0|$ , and the norm on  $\mathbb{C}^{M^0}$ .
4. If  $P = P_A$  then there exists no set of matrices  $\{\mathfrak{C}^{(m)}, |\mathbf{m}| = P_A + 1\}$ , such that the scheme possesses the truncation error of order  $P_A + 1$  in the sense of  $\Pi_h^{(P_A+1, P_A+1)}$  given by (2.12).
5. If  $L(0) = 0$  then  $P = P_A$ .
6. If the scheme possesses the formal order of accuracy  $P_A + 1$  on  $H_{per,e}^R(\mathbb{R}^d)$  (or  $C_{per,e}^R(\mathbb{R}^d)$ ) for each  $e \in \mathring{\Omega}_{P_A+1}$ , then it possesses the formal order of accuracy  $P = P_A + 1$  on  $H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^R(\mathbb{R}^d)$ ).
7.  $P$  coincides with the optimal order of accuracy in the weak sense.

**Theorem 2.** Let  $\Pi_h$  be a local mapping with  $\mu_{\xi} \in (W_2^q(G))^*$  or  $(C^q(G))^*$ . Let  $P, Q > 0$ ,  $r \geq \max\{P, q\}$ , and  $R = \max\{Q + 1, r\}$ . Let the scheme (2.8) be stable and possess the error estimate (2.17) on  $H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^{[R]}(\mathbb{R}^d)$ ). Then there exists a homogeneous mapping  $\tilde{\Pi}_h : L_{2,per}(\mathbb{R}^d) \rightarrow V_{per}$  such that

$$\|\tilde{\Pi}_h f - \Pi_h f\| \leq C(h^P \|\nabla^P f\| + h^r \|\nabla^r f\|_*), \quad \|\epsilon_h(f, \tilde{\Pi}_h)\| \leq Ch^Q \|\nabla^{Q+1} f\|$$

for each  $h$  and  $f \in H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^{[R]}(\mathbb{R}^d)$ ), where  $C$  does not depend on  $h$  and  $f$ , and  $\|\cdot\|_*$  means either  $\|\cdot\|$  or  $\|\cdot\|_{\infty}$  depending on the case.

Note that the mapping  $\tilde{\Pi}_h$  given by this theorem is generally *not* local.

In order to formulate the third result, we need two more definitions.

**Definition 6.** We say that a scheme of the form (2.8) is *quasi-1D* if the stencil  $\mathcal{S}$  of the scheme belongs to a 1D subset of  $\mathbb{Z}^d$ , i. e. there exists  $\eta \in \mathbb{Z}^d$  such that  $\mathcal{S} \subset \{m\eta, m \in \mathbb{Z}\}$ .

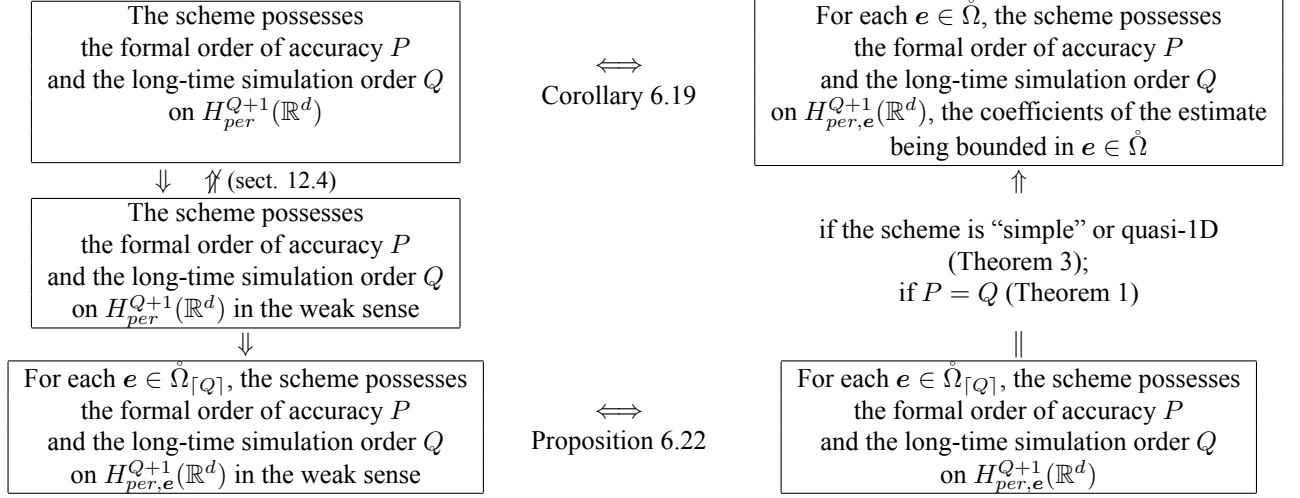
In 1D case, each scheme is quasi-1D.

**Definition 7.** We say that a scheme of the form (2.8) is *simple* if the matrix  $L(0) \equiv \sum_{\zeta} L_{\zeta}$  has rank  $|M^0| - 1$ .

Note that if  $L(0)$  has rank  $|M^0|$  then the scheme does not preserve a constant solution, thus its order of the truncation error is equal to  $P_A = -1$ , then  $P \leq 0$  by Theorem 1, i. e. there is no solution convergence.

**Theorem 3.** Let  $\Pi_h$  and  $\mathcal{P}_h$  be local mappings with kernels  $\mu_{\xi}, \hat{\mu}_{\xi} \in (W_2^r(G))^*$  (or  $(C^r(G))^*$ ). Let  $\hat{\kappa} = \min_{\xi \in M^0} |\langle \hat{\mu}_{\xi}, 1 \rangle| > 0$ . Let the scheme (2.8) be 1) stable and 2) either quasi-1D or simple. Let  $P_A, P$ , and

Table 1: Relations between the definitions



$Q$  be the optimal values of the order of the truncation error, the formal order of accuracy, and the long-time simulation order, correspondingly. Then the following holds.

1.  $Q \in \mathbb{N} \cup \{0\}$ .
2. There exist real-valued diagonal matrices  $\mathfrak{C}^{(\mathbf{m})}$ ,  $P \leq |\mathbf{m}| \leq Q$ , such that the scheme possesses the truncation error of order  $Q$  in the sense of  $\Pi_h^{(P,Q)}$  given by (2.12).
3. Let  $p, q > 0$ . If for each  $e \in \check{\Omega}_{[q]}$  the scheme possesses the formal order of accuracy  $p$  and the long-time simulation order  $q$  on  $H_{per,e}^R(\mathbb{R}^d)$  (or  $C_{per,e}^{[R]}(\mathbb{R}^d)$ ), where  $R \geq \max\{r, p, q + 1\}$ , then it possesses the formal order of accuracy  $p$  and the long-time simulation order  $q$  on  $H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^{[R]}(\mathbb{R}^d)$ ). Particularly, if  $p = P$ , then  $Q \geq q$ .
4.  $Q$  coincides with the optimal value of the long-time simulation order in the weak sense.
5. If the scheme is simple, then either  $Q \geq P = P_A + 1$  or  $Q = P = P_A$ .

If the scheme is neither quasi-1D nor simple, then the statements 1–4 of Theorem 3 may be wrong, see counter-examples to the statements 1 and 2 in Section 12.4, to the statement 3 in Sections 12.2 and 12.3, to the statement 4 in Section 12.1. If the scheme is not simple, then the statement 5 of Theorem 3 may also be wrong even in 1D case, see Section 11.3.

The optimal value of the formal order of accuracy is the same for Definitions 3 and 4, but this generally does not hold for the long-time simulation order. Relations between the definitions are shown in Table 1. For each implication, a reference to a corresponding statement is given unless it is obvious. For an implication that does not hold, a reference to a counter-example is given.

Basing on Theorems 1 and 3 we construct algorithms giving the optimal values of the formal order of accuracy (for each stable scheme) and of the long-time simulation order (for simple and quasi-1D schemes), see Section 10. Theorem 2 is used below to establish Theorem 8.5 which states that a scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  if and only if a specific function is bounded at zero.

#### 4. The method of auxiliary mapping

**Theorem 4.1** (Lax – Ryabenkii). *Let  $\Pi_h$  be a map of  $H_{per}^{Q+1}(\mathbb{R}^d)$  to  $V_{per}$ . Let  $K$  be the stability constant of the scheme (2.8). Let for each  $f \in H_{per}^{Q+1}(\mathbb{R}^d)$  there hold*

$$\|\epsilon(f, \Pi_h)\| \leq Ch^Q \|\nabla^{Q+1} f\|. \quad (4.1)$$

Then for each  $v_0$  there holds

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq CKth^Q \|Z^{-1}\| \|\nabla^{Q+1} v_0\|. \quad (4.2)$$

*Proof.* Put  $v(t, \mathbf{r}) = v_0(\mathbf{r} - \boldsymbol{\omega}t)$ . Taking the time derivative of (2.15), multiplying by  $Z$  and using (2.10) we get

$$\begin{aligned} Z \frac{d}{dt} \varepsilon_h(t, v_0, \Pi_h) &= -\frac{1}{h} Lu(t) + Z \Pi_h(\boldsymbol{\omega} \cdot \nabla) v(t, \cdot) = \\ &= -\frac{1}{h} L(\varepsilon_h(t, v_0, \Pi_h) + \Pi_h v(t, \cdot)) + Z \Pi_h(\boldsymbol{\omega} \cdot \nabla) v(t, \cdot). \end{aligned}$$

By (2.14) this is equivalent to

$$Z \frac{d}{dt} \varepsilon_h(t, v_0, \Pi_h) + \frac{1}{h} L \varepsilon_h(t, v_0, \Pi_h) = -\varepsilon_h(v(t, \cdot), \Pi_h).$$

Since  $\varepsilon_h(0, v_0, \Pi_h) = 0$  (see (2.15)),

$$\varepsilon_h(t, v_0, \Pi_h) = - \int_0^t \exp\left(-Z^{-1} L \frac{t-\tau}{h}\right) Z^{-1} \varepsilon_h(v(\tau, \cdot), \Pi_h) d\tau.$$

By stability  $\|\exp(-\nu Z^{-1} L)\| \leq K$  for all  $\nu \geq 0$ , so using (4.1) we have

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq Kt \max_{0 \leq \tau \leq t} \|Z^{-1} \varepsilon_h(v(\tau, \cdot), \Pi_h)\| \leq CKth^Q \|Z^{-1}\| \cdot \max_{0 \leq \tau \leq t} \|\nabla^{Q+1} v(\tau, \cdot)\|.$$

It remains to note that  $\|\nabla^{Q+1} v(\tau, \cdot)\| = \|\nabla^{Q+1} v_0\|$  for all  $\tau$ . □

The following proposition describes the method of auxiliary mapping, which is the cornerstone of this paper.

**Proposition 4.2.** Consider a stable scheme (2.8) and homogeneous mappings  $\Pi_h$  and  $\tilde{\Pi}_h$  from  $L_{2,per}(\mathbb{R}^d)$  to  $V_{per}$ . Suppose for some  $P$  and  $Q$  for each  $f \in H_{per}^{\max\{P, Q+1\}}$  there holds

$$\|(\Pi_h - \tilde{\Pi}_h)f\| \leq C_1 h^P \|\nabla^P f\| \quad \text{and} \quad \|\epsilon(f, \tilde{\Pi}_h)\| \leq C_2 h^Q \|\nabla^{Q+1} f\|.$$

Then the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  in the sense of  $\Pi_h$ .

*Proof.* Let  $u(t)$  and  $\tilde{u}(t)$  be the solutions of (2.8) with the initial data  $u(0) = \Pi_h v_0$  and  $\tilde{u}(0) = \tilde{\Pi}_h v_0$ , correspondingly. Denote by  $v$  the exact solution  $v(t, \mathbf{r}) = v_0(\mathbf{r} - \boldsymbol{\omega}t)$ . Then

$$\begin{aligned} \|\varepsilon_h(t, v_0, \Pi_h)\| &= \|u(t) - \Pi_h v(t, \cdot)\| \leq \\ &\leq \|u(t) - \tilde{u}(t)\| + \|\tilde{u}(t) - \tilde{\Pi}_h v(t, \cdot)\| + \|\tilde{\Pi}_h v(t, \cdot) - \Pi_h v(t, \cdot)\|. \end{aligned}$$

By stability there holds

$$\|u(t) - \tilde{u}(t)\| \leq K \|u(0) - \tilde{u}(0)\| = K \|\tilde{\Pi}_h v_0 - \Pi_h v_0\|,$$

where  $K$  is the stability constant of the scheme. Note that  $\|v(t, \cdot)\| \equiv \|v_0\|$ . Then we get

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq \|\varepsilon_h(t, v_0, \tilde{\Pi}_h)\| + (K+1)C_1 h^P \|\nabla^P v_0\|. \quad (4.3)$$

270 By Theorem 4.1 we get the desired estimate (2.17) or (2.18). □

## 5. Preliminaries

### 5.1. Some properties of the functional spaces

Recall that for  $n \in \mathbb{N}$ , by  $\mathring{\Omega}_n$  we denote a set of vectors  $\{e_k \in \mathring{\Omega}, k = 1, \dots, C_{n+d-1}^{d-1}\}$  such that the polynomials  $(e_k \cdot r)^n$  form a basis in the space of homogeneous polynomials of order  $n$  (see Lemma A.6).

**Lemma 5.1.** *Let  $m$  be a multiindex,  $\mathring{\Omega}_{|m|} = \{e_k\}$ . Then there exists  $\gamma_k^{(m)} \in \mathbb{R}, k = 1, \dots, C_{|m|+d-1}^{d-1}$ , such that*

$$D^m = \sum_{k=1}^{C_{|m|+d-1}^{d-1}} \gamma_k^{(m)} (e_k \cdot \nabla)^{|m|}. \quad (5.1)$$

*Proof.* Since the polynomials  $(e_k \cdot r)^n$  form a basis in the space of homogeneous polynomials of order  $n$ , there exists a set of coefficients  $\gamma_k^{(m)}$  such that

$$r^m = \sum_{k=1}^{C_{|m|+d-1}^{d-1}} \gamma_k^{(m)} (e_k \cdot r)^{|m|}.$$

Since this is a polynomial identity, it holds with the substitution of  $\nabla$  for  $r$ . Thus we get (5.1).  $\square$

**Lemma 5.2.** *Let  $G$  be a neighborhood of zero in  $\mathbb{R}^d$ ,  $f \in C^p(G)$ ,  $p \in \mathbb{N}$  and  $|f(x)| \leq c_f |x|^p$  in  $G$ . Then  $|D^m f(0)| \leq c_f c_p$  for each  $|m| = p$ , where  $c_p$  does not depend on  $f$ .*

*Proof.* For the 1D case, this follows from the Taylor expansion with remainder term in Peano form with  $c_p = p!$ . Thus for each direction  $e \in \Omega$  we have  $|(e \cdot \nabla)^p f| \leq c_f p!$ , and it remains to use Lemma 5.1.  $\square$

**Lemma 5.3.** *Let  $f(\phi)$  be a holomorphic function at  $\phi = 0$ , and  $n, m \in \mathbb{N} \cup \{0\}$ ,  $m \geq n$ . Suppose for each  $e \in \mathring{\Omega}_m$  there exists  $c_e \geq 0$  such that  $|f(e\psi)| \leq c_e |\psi|^{n+1}$  in a neighborhood of  $\psi = 0$ . Then there exists  $\tilde{c}$  such that  $|f(\phi)| \leq \tilde{c} |\phi|^{n+1}$  in a neighborhood of  $\phi = 0$ .*

*Proof.* By Lemma A.5, the system  $\{(e \cdot r)^n, e \in \mathring{\Omega}_m\}$  is complete in the space of homogeneous polynomials of order  $n$ . Let  $g(e, \phi) = \partial^n f / \partial e^n(\phi)$ . As a function of  $e$ ,  $g$  is a homogeneous polynomial of order  $n$ . By assumption  $g(e, 0) = 0$  for  $e \in \mathring{\Omega}_m$ . So we have  $g(e, 0) = 0$  for each  $e \in \Omega$ . By Lemma 5.1 for each  $m$ ,  $|m| = n$ , we have  $D^m f(\phi) = 0$ . Then the statement of the Lemma is obvious.  $\square$

**Lemma 5.4.** *For each  $p \leq m \leq q$  and each  $w \in H_{per}^q(\mathbb{R}^d)$  there holds*

$$h^{2m} \|\nabla^m w\|^2 \leq h^{2p} \|\nabla^p w\|^2 + h^{2q} \|\nabla^{2q} w\|^2. \quad (5.2)$$

*Proof.* For  $p = m = q$  this is obvious so assume  $q > p$ . Let  $w_\alpha$  be Fourier coefficients of  $w$ . By (2.5),

$$\begin{aligned} h^{2m} \|\nabla^m w\|^2 &= \sum_{\alpha \in \mathcal{A}} |w_\alpha|^2 |h\alpha|^{2m} \leq \\ &\leq \sum_{\alpha \in \mathcal{A}} |w_\alpha|^2 \left( \frac{q-m}{q-p} |h\alpha|^{2p} + \frac{m-p}{q-p} |h\alpha|^{2q} \right) \leq h^{2p} \|\nabla^p w\|^2 + h^{2q} \|\nabla^q w\|^2, \end{aligned}$$

where we have used the Young inequality for products.  $\square$

**Lemma 5.5.** *For each  $p, q \in \mathbb{N} \cup \{0\}$  there exists  $C(p, q)$  such that for each  $m = p, \dots, q$  and  $f \in C_{per}^q(\mathbb{R}^d)$  there holds*

$$h^m \|\nabla^m f\|_\infty \leq C(p, q) (h^p \|\nabla^p f\|_\infty + h^q \|\nabla^q f\|_\infty).$$

*Proof.* By the coordinate transformation we can assume without loss that  $h = 1$ . For 1D case this follows from the Landau – Kolmogorov inequality [25]. For each  $e \in \Omega$  this yields

$$\left\| \frac{\partial^m f}{\partial e^m} \right\|_\infty \leq C(p, q) \left( \left\| \frac{\partial^p f}{\partial e^p} \right\|_\infty + \left\| \frac{\partial^q f}{\partial e^q} \right\|_\infty \right).$$

Then by Lemma 5.1 we have

$$\begin{aligned} \|\nabla^m f\|_\infty^2 &= \sum_{|\mathbf{m}|=m} \frac{m!}{\mathbf{m}!} \max_{\mathbb{R}^d} |D^{\mathbf{m}} f|^2 \leq \sum_{|\mathbf{m}|=m} \frac{m!}{\mathbf{m}!} \gamma^2 \sup_{e \in \Omega} \left\| \frac{\partial^m f}{\partial e^m} \right\|_\infty^2 \leq \\ &\leq \sum_{|\mathbf{m}|=m} \frac{m!}{\mathbf{m}!} \gamma^2 C^2(p, q) \sup_{e \in \Omega} \left( \left\| \frac{\partial^p f}{\partial e^p} \right\|_\infty + \left\| \frac{\partial^q f}{\partial e^q} \right\|_\infty \right)^2 \leq \tilde{C}(p, q) (\|\nabla^p f\|_\infty + \|\nabla^q f\|_\infty)^2, \end{aligned}$$

290 where  $\gamma = \max_{|\mathbf{m}|=m} \sum_k |\gamma_k^{(\mathbf{m})}|$  and  $\gamma_k^{(\mathbf{m})}$  are given by Lemma 5.1. □

**Lemma 5.6.** *Let  $0 \leq r \leq q < \infty$ . Then for each  $f \in H_{per}^q(\mathbb{R}^d)$  there holds  $\|f\|_{(r,h)} \leq 2\|f\|_{(q,h)}$ . For  $r, q \in \mathbb{N} \cup \{0\}$ ,  $r \leq q$ , there exists  $C(p, q)$  such that for each  $f \in C_{per}^q(\mathbb{R}^d)$  there holds  $\|f\|_{(r,h,\infty)} \leq C(p, q)\|f\|_{(q,h,\infty)}$ .*

The proof is obvious by Lemma 5.4 and Lemma 5.5.

**Lemma 5.7.** *For each  $v_0 \in H_{per}^{r+1}(\mathbb{R}^d)$  there holds*

$$\|\nabla^r(\omega \cdot \nabla v_0)\| \leq |\omega| \|\nabla^{r+1} v_0\|, \quad h\|\omega \cdot \nabla v_0\|_{(r,h)} \leq \sqrt{2}|\omega| \|v_0\|_{(r+1,h)}.$$

*There exists  $c_r$  such that for each  $v_0 \in C_{per}^{r+1}(\mathbb{R}^d)$  there holds*

$$\|\nabla^r(\omega \cdot \nabla v_0)\|_\infty \leq |\omega| \|\nabla^{r+1} v_0\|_\infty, \quad h\|\omega \cdot \nabla v_0\|_{(r,h,\infty)} \leq c_r |\omega| \|v_0\|_{(r+1,h,\infty)}.$$

*Proof.* For  $v_0 \in H_{per}^{r+1}(\mathbb{R}^d)$ ,  $r \geq 0$ , we have

$$\|\nabla^r(\omega \cdot \nabla v_0)\|^2 = \sum_{\alpha \in \mathcal{A}} |f_\alpha|^2 (\omega \cdot \alpha)^2 |\alpha|^{2r} \leq \sum_{\alpha \in \mathcal{A}} |f_\alpha|^2 |\omega|^2 |\alpha|^{2r+2} = |\omega|^2 \|\nabla^{r+1} v_0\|^2.$$

Then, by Lemma 5.4,

$$\begin{aligned} h^2 \|\omega \cdot \nabla v_0\|_{(r,h)}^2 &= h^2 \|\omega \cdot \nabla v_0\|^2 + h^{2r+2} \|\nabla^r(\omega \cdot \nabla v_0)\|^2 \leq h^2 |\omega|^2 \|\nabla v_0\|^2 + h^{2r+2} |\omega|^2 \|\nabla^{r+1} v_0\|^2 \leq \\ &\leq |\omega|^2 (\|v_0\|^2 + 2h^{2r+2} \|\nabla^{r+1} v_0\|^2) \leq 2|\omega|^2 \|v_0\|_{(r+1,h)}^2. \end{aligned}$$

For  $v_0 \in C_{per}^{r+1}(\mathbb{R}^d)$ ,  $r \in \mathbb{N} \cup \{0\}$ , we have

$$\begin{aligned} \|\nabla^r(\omega \cdot \nabla v_0)\|_\infty^2 &= \sum_{|\mathbf{m}|=r} \frac{r!}{\mathbf{m}!} \|D^{\mathbf{m}}(\omega \cdot \nabla v_0)\|_\infty^2 = \\ &= \sum_{|\mathbf{m}|=r} \frac{r!}{\mathbf{m}!} \left\| \sum_{|\mathbf{l}|=1} \omega^{\mathbf{l}} D^{\mathbf{m}+\mathbf{l}} v_0 \right\|_\infty^2 \leq \sum_{|\mathbf{m}|=r} \frac{r!}{\mathbf{m}!} \left( \sum_{|\mathbf{l}|=1} |\omega^{\mathbf{l}}| \|D^{\mathbf{m}+\mathbf{l}} v_0\|_\infty \right)^2 \leq |\omega|^2 \sum_{|\mathbf{m}|=r} \frac{r!}{\mathbf{m}!} \sum_{|\mathbf{l}|=1} \|D^{\mathbf{m}+\mathbf{l}} v_0\|_\infty^2. \end{aligned}$$

Denote  $\mathbf{n} = \mathbf{m} + \mathbf{l}$ , then we continue the chain:

$$\|\nabla^r(\omega \cdot \nabla v_0)\|_\infty^2 \leq |\omega|^2 \sum_{|\mathbf{n}|=r+1} \sum_{|\mathbf{l}|=1, \mathbf{l} \leq \mathbf{n}} \frac{(r+1)!}{(\mathbf{n}-\mathbf{l})! (r+1)} \|D^{\mathbf{n}} v_0\|_\infty^2 =$$

$$= |\omega|^2 \sum_{|n|=r+1} \frac{(r+1)!}{n!} \|D^n v_0\|_\infty^2 \sum_{|l|=1, l \leq n} \frac{n^l}{|n|} = |\omega|^2 \sum_{|n|=r+1} \frac{(r+1)!}{n!} \|D^n v_0\|_\infty^2 = |\omega|^2 \|\nabla^{r+1} v_0\|_\infty^2.$$

Then, by Lemma 5.5,

$$\begin{aligned} h^2 \|\omega \cdot \nabla v_0\|_{(r,h,\infty)}^2 &= h^2 \|\omega \cdot \nabla v_0\|_\infty^2 + h^{2r+2} \|\nabla^r(\omega \cdot \nabla v_0)\|_\infty^2 \leq h^2 |\omega|^2 \|\nabla v_0\|_\infty^2 + h^{2r+2} |\omega|^2 \|\nabla^{r+1} v_0\|_\infty^2 \leq \\ &\leq |\omega|^2 (C(1, r+1) \|v_0\|^2 + (1 + C(1, r+1)) h^{2r+2} \|\nabla^{r+1} v_0\|^2) \leq (1 + C(1, r+1)) |\omega|^2 \|v_0\|_{(r+1,h,\infty)}^2. \end{aligned}$$

□

### 5.2. Approximation of local mappings

In this subsection we prove that any local mapping can be approximated by a local mapping with a continuous kernel.

**Lemma 5.8.** *Let  $s \in \mathbb{N}$ . Then there exists  $\varkappa^{(s)} \in C^s(\{|x| \leq 1\})$  such that the following conditions hold.*

300

- (i)  $\varkappa^{(s)}$  and its derivatives up to the order  $s$  are equal to zero on  $|x| = 1$ ;
- (ii) there holds

$$\int_{|x| < 1} \varkappa^{(s)}(x) dx = 1; \quad (5.3)$$

- (iii) for each  $\sigma = 0, \dots, s+1$ , each  $f \in C_{per}^\sigma(\mathbb{R}^d)$ , and each  $y \in \mathbb{R}^d$  there holds

$$\left| \int_{|x| < 1} \varkappa^{(s)}(x) f(y + hx) dx - f(y) \right| \leq C(d, s) h^\sigma \|\nabla^\sigma f\|_\infty. \quad (5.4)$$

*Proof.* Put  $\varkappa^{(s)}(x) = |x|^{2s} (1 - |x|)^s g(|x|)$ , where  $g(r)$  is a smooth function on  $[0, 1]$  such that

$$|\Omega| \int_0^1 r^{2s+d-1} (1-r)^s g(r) dr = 1, \quad (5.5)$$

$$\int_0^1 r^{2s+d-1+m} (1-r)^s g(r) dr = 0, \quad m = 1, \dots, s. \quad (5.6)$$

The factor  $|x|^{2s}$  guarantees the existence and continuity of the derivatives of  $\varkappa^{(s)}(x)$  up to the order  $s$  at  $r = 0$ , and the factor  $(1 - |x|)^s$  makes the boundary conditions satisfied. The condition (5.5) yields (5.3). Now we show (5.4). For  $\sigma = 0$  inequality (5.4) is obvious, so assume  $\sigma > 0$ . For each  $y$  and  $x$  there holds

$$\left| f(y + hx) - f(y) - \sum_{k=1}^{\sigma-1} \frac{1}{k!} \frac{\partial^k f(y)}{\partial e^k} h^k |x|^k \right| \leq \|\nabla^\sigma f\|_\infty \frac{h^\sigma |x|^\sigma}{\sigma!},$$

where  $e$  is the unit vector along  $x$ . Hence, using (5.3) we get

$$\begin{aligned} \left| \int \varkappa^{(s)}(x) f(y + hx) dx - f(y) \right| &= \left| \int \varkappa^{(s)}(x) (f(y + hx) - f(y)) dx \right| \leq \\ &\leq \left| \int \varkappa^{(s)}(x) \left[ \sum_{k=1}^{\sigma-1} \frac{1}{k!} \frac{\partial^k f(y)}{\partial e^k} h^k |x|^k \right] dx \right| + \|\nabla^\sigma f\|_\infty \frac{h^\sigma}{\sigma!} \int |\varkappa^{(s)}(x)| dx \end{aligned}$$

(all integrals are over  $|x| < 1$ ). We claim that the first term in the right-hand size is equal to zero. Indeed, put  $x = r\gamma$ , where  $0 < r < 1$  and  $\gamma$  varies over  $\Omega$ . Then for  $k = 1, \dots, \sigma-1$  we have

$$\int_{|x| < 1} \varkappa^{(s)}(x) \frac{\partial^k f(y)}{\partial e^k} |x|^k dx = \int_0^1 r^{2s} (1-r)^s r^{d-1} r^k g(r) dr \int_\Omega \frac{\partial^k f(y)}{\partial e_\gamma^k} d\gamma,$$

that is equal to zero by (5.6). Thus the inequality (5.4) is proved. □

**Lemma 5.9.** Let  $\tilde{\Pi}_h$  be a local mapping with kernel  $\tilde{\mu}_\xi \in (C^q(G))^*$ ,  $q \in \mathbb{N} \cup \{0\}$ . Let  $s \geq q$ ,  $s \in \mathbb{N}$ , and  $\varkappa^{(s)}$  be given by Lemma 5.8. Then the convolution  $\mu_\xi = \tilde{\mu}_\xi * \varkappa^{(s)}$  is a continuous function with support in  $G + B_1(0)$ , and the local mapping  $\Pi_h$  with kernel  $\mu_\xi$  satisfies

$$\|\tilde{\Pi}_h f - \Pi_h f\| \leq Ch^r \|\nabla^r f\|_\infty$$

for each  $r = q, \dots, s$  and all  $f \in C_{per}^r(\mathbb{R}^d)$ , where  $C$  depends only on  $d, s, \tilde{\mu}_\xi$ , and the norm on  $\mathbb{C}^{M^0}$ .

*Proof.* For each  $\xi \in M^0$ , each  $h$ , and  $\eta \in \mathbb{Z}^d$  we have

$$(\tilde{\Pi}_h f)_{\eta, \xi} - (\Pi_h f)_{\eta, \xi} = \langle \tilde{\mu}_\xi - \mu_\xi, f(h(\cdot + T\eta)) \rangle = \langle \tilde{\mu}_\xi, \Phi[f] \rangle,$$

where  $\Phi[f]$  is given by

$$\Phi[f](\mathbf{r}) = \int_{|\mathbf{x}| < 1} \varkappa^{(s)}(\mathbf{x}) f(h(\mathbf{r} + T\eta + \mathbf{x})) d\mathbf{x} - f(h(\mathbf{r} + T\eta)).$$

For  $|\mathbf{m}| \leq q$  we have  $D^{\mathbf{m}} \Phi[f] = h^{|\mathbf{m}|} \Phi[D^{\mathbf{m}} f]$ . Then by (5.4) with  $\sigma = r - |\mathbf{m}|$  we get

$$\sup_{\mathbf{r} \in \mathbb{R}^d} |D^{\mathbf{m}} \Phi[f](\mathbf{r})| \leq h^{|\mathbf{m}|} C(d, s) h^{r-|\mathbf{m}|} \|\nabla^{r-|\mathbf{m}|}(D^{\mathbf{m}} f)\|_\infty \leq h^r C(d, s) \|\nabla^r f\|_\infty.$$

Equip  $C^q(G)$  with the norm

$$\|g\|_{C^q(G)} = \max_{|\mathbf{m}| \leq q} \max_{\mathbf{r} \in G} |D^{\mathbf{m}} g|.$$

Then  $\|\Phi[f]\|_{C^q(G)} \leq h^r C(d, s) \|\nabla^r f\|_\infty$  and

$$\left| (\tilde{\Pi}_h f)_{\eta, \xi} - (\Pi_h f)_{\eta, \xi} \right| \leq \|\tilde{\mu}_\xi\| h^r C(d, s) \|\nabla^r f\|_\infty.$$

Then the inequality to prove is by the definition of the norm on  $V_{per}$ . □

### 5.3. Properties of the matrix exponential

305

In this subsection we consider a matrix norm induced by a vector norm.

**Lemma 5.10.** If  $A(\phi)$  is a holomorphic matrix function at  $\phi = \phi_0$ , then so is  $\exp(A(\phi))$ .

This follows from the representation

$$\exp(A(\phi)) = \frac{1}{2\pi i} \int_{\gamma} \frac{e^z}{zI - A(\phi)} dz,$$

where  $I$  is the identity matrix and  $\gamma$  is any closed contour such that all eigenvalues of  $A(\phi_0)$  lay inside  $\gamma$ .

**Lemma 5.11.** For each square matrix  $Y$  there holds

$$\|e^Y - I\| \leq (\|e^Y\| + e - 1) \min\{1, \|Y\|\}.$$

*Proof.* If  $\|Y\| \geq 1$  then

$$\|e^Y - I\| \leq \|e^Y\| + 1 \leq \|e^Y\| + e - 1.$$

If  $\|Y\| \leq 1$  then

$$\|e^Y - I\| = \left\| Y \sum_{k=1}^{\infty} \frac{Y^{k-1}}{k!} \right\| \leq \|Y\| \sum_{k=1}^{\infty} \frac{1}{k!} = (e - 1)\|Y\| \leq (\|e^Y\| + e - 1)\|Y\|.$$

□



**Lemma 5.12.** Suppose  $\|A\| \leq 1$  and

$$f(A) = \sum_{k=1}^{\infty} \frac{A^{k-1}}{k!}. \quad (5.7)$$

Then  $\|(f(A))^{-1}\| \leq 4$ .

*Proof.* First,

$$\|f(A) - I\| = \left\| \sum_{k=2}^{\infty} \frac{A^{k-1}}{k!} \right\| \leq \sum_{k=2}^{\infty} \frac{\|A\|^{k-1}}{k!} \leq \sum_{k=2}^{\infty} \frac{1}{k!} = e - 2.$$

Hence

$$\|(f(A))^{-1}\| = \|(I + (f(A) - I))^{-1}\| \leq \frac{1}{1 - \|f(A) - I\|} \leq \frac{1}{1 - (e - 2)} \leq 4.$$

310

□

**Lemma 5.13.** Let  $A$  be non-degenerate and  $\|A\| \leq 1$ . Then  $\|(e^A - I)^{-1}\| \leq 4\|A^{-1}\|$ .

*Proof.* For  $f(A)$  given by (5.7) we have  $e^A - I = Af(A)$ . Then

$$\|(e^A - I)^{-1}\| \leq \|(f(A))^{-1}\| \|A^{-1}\| \leq 4\|A^{-1}\|.$$

□

#### 5.4. Transformation of a matrix to a block-diagonal form

Let  $\mathbb{C}^{n \times n}$  be the space of complex matrices of size  $n$ . Denote by  $\mathcal{A}(\mathbb{C}^d, \mathbb{C}^{n \times n})$  the set of functions from  $\mathbb{C}^d$  to  $\mathbb{C}^{n \times n}$  holomorphic at  $\phi = 0$ . Let  $\sigma(Y)$  be the spectrum of a matrix  $Y$  and  $\varkappa(Y)$  be its condition number, i. e.  $\varkappa(Y) = \|Y\| \|Y^{-1}\|$ .

To proceed, we need to recall some facts from the perturbation theory.

**Lemma 5.14** ([26], §1.5.3). Let  $A \in \mathbb{C}^{n \times n}$ ,  $\lambda \in \sigma(A)$ . Let  $0 < \rho < \text{dist}(\lambda, \sigma(A) \setminus \{\lambda\})$  (the distance to the empty set is assumed to be  $+\infty$ ). Then the matrix

$$P_\lambda(A) = \frac{1}{2\pi i} \oint_{|z-\lambda|=\rho} (zI - A)^{-1} dz \quad (5.8)$$

is a projection onto the algebraic eigenspace corresponding to  $\lambda$ . Besides, for  $\lambda, \mu \in \sigma(A)$  there holds  $P_\lambda(A)P_\mu(A) = \delta_{\lambda\mu}P_\lambda(A)$ .

315 **Lemma 5.15** ([26], §2.5.1). Let  $A(\phi) \in \mathcal{A}(\mathbb{C}^d, \mathbb{C}^{n \times n})$ . Then the eigenvalues  $\lambda_j(\phi)$ ,  $j = 1, \dots, n$ , can be reordered so that  $\lambda_j(\phi) \rightarrow \lambda_j(0)$  as  $\phi \rightarrow 0$ .

For  $A(\phi) \in \mathcal{A}(\mathbb{C}^d, \mathbb{C}^{n \times n})$  and  $\lambda \in \sigma(A(0))$  the set of eigenvalues  $\lambda_j(\phi)$  of the matrix  $A(\phi)$  such that  $\lambda_j(\phi) \rightarrow \lambda$  is referred to as  $\lambda$ -group, and the sum of the corresponding algebraic eigenspaces is the *total eigenspace* of the  $\lambda$ -group. The operator  $\mathcal{P}_\lambda(\phi) = \sum P_{\lambda_j(\phi)}(A(\phi))$ , where the sum is over the  $\lambda$ -group, is the *total eigenprojection* of the  $\lambda$ -group. It is a projection onto the total eigenspace of the  $\lambda$ -group (see [26], §2.2.1). The sum of all total eigenprojections is the identity operator and  $\mathcal{P}_\lambda(\phi)\mathcal{P}_\mu(\phi) = \delta_{\lambda\mu}\mathcal{P}_\lambda(\phi)$ .

325

**Lemma 5.16.** Let  $A(\phi) \in \mathcal{A}(\mathbb{C}^d, \mathbb{C}^{n \times n})$  and  $\lambda \in \sigma(A(0))$ . Then  $\mathcal{P}_\lambda(\phi) \in \mathcal{A}(\mathbb{C}^d, \mathbb{C}^{n \times n})$ .

*Proof.* Let  $\rho = \frac{1}{2}\text{dist}(\lambda, \sigma(A(0)) \setminus \{\lambda\})$ . By continuity in a neighborhood of  $\phi = 0$  there holds  $|\lambda_j(\phi) - \lambda| < \rho$  for each  $\lambda_j(\phi)$  belonging to the  $\lambda$ -group corresponding to  $\lambda$  and  $|\lambda_j(\phi) - \lambda| > \rho$  otherwise. Thus the total eigenprojection of this  $\lambda$ -group is

$$\mathcal{P}_\lambda(\phi) = \frac{1}{2\pi i} \oint_{|z-\lambda|=\rho} (zI - A(\phi))^{-1} dz. \quad (5.9)$$

Since the integrand is analytical with respect to  $\phi$  and  $z$  over the integral path, we have  $\mathcal{P}_\lambda(\phi) \in \mathcal{A}(\mathbb{C}^d, \mathbb{C}^{n \times n})$ . □

**Lemma 5.17.** For any  $n \in \mathbb{N}$ ,  $K \geq 1$  there exists  $\hat{C} = \hat{C}(n, K)$  such that for any  $B \in \mathbb{C}^{n \times n}$  satisfying  $\sup_{\nu \geq 0} \|e^{\nu B}\| \leq K$  there holds

$$\frac{|\lambda(B)|_{\max}}{|\lambda(B)|_{\min}} \leq \kappa(B) \leq \hat{C} \frac{|\lambda(B)|_{\max}}{|\lambda(B)|_{\min}},$$

where  $|\lambda(B)|_{\max}$  and  $|\lambda(B)|_{\min}$  are the maximal and minimal module of the eigenvalues of  $B$ .

*Proof.* The first inequality follows from the fact that any norm of the matrix induced by a vector norm is greater than or equal to its spectral radius. The proof of the second inequality is based on the Kreiss matrix theorem [27], see Lemma 6 and Lemma 8 in [28].  $\square$

**Theorem 5.18** ([28]). For any  $n \in \mathbb{N}$ ,  $K \geq 1$  there exists  $C = C(n, K)$  such that for any  $B \in \mathbb{C}^{n \times n}$  satisfying  $\sup_{\nu \geq 0} \|e^{\nu B}\| \leq K$  there exists a matrix  $U$  such that  $\kappa(U) \leq C$ , the matrix  $M = U^{-1}BU$  is block-diagonal, and each block  $M_j$  of  $M$  satisfies either  $M_j = 0$  or  $\kappa(M_j) \leq C$ .

Without loss we assume the spectra of  $M_j$  do not intersect. Indeed, if two blocks  $M_j$  have common eigenvalues then by Lemma 5.17 they can be united into one block with the condition number not greater than the product of the conditional numbers of the original blocks multiplied by  $\hat{C}$ .

**Theorem 5.19.** Let  $A \in \mathcal{A}(\mathbb{C}^d, \mathbb{C}^{n \times n})$  and  $\sup_{\nu \geq 0} \|\exp(\nu A(0))\| \leq K$ . Let  $A(0)$  have zero eigenvalue of algebraic multiplicity  $n_0$ . Then in a neighborhood of  $\phi = 0$  there holds  $A(\phi) = S(\phi)M(\phi)S^{-1}(\phi)$  where  $S, M, S^{-1} \in \mathcal{A}(\mathbb{C}^d, \mathbb{C}^{n \times n})$ ,

$$M(\phi) = \begin{pmatrix} M^{(1)}(\phi) & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & M^{(m)}(\phi) & 0 \\ 0 & \dots & 0 & M^{(0)}(\phi) \end{pmatrix}, \quad (5.10)$$

$M^{(0)}(\phi)$  is a  $n_0 \times n_0$ -matrix,  $M^{(0)}(0) = 0$ ,  $\kappa(M^{(j)}(0)) \leq \delta$  for  $j \neq 0$ ,  $\|S(0)\| \leq \delta$ ,  $\|S^{-1}(0)\| \leq \delta$ , and  $\delta$  depends on  $n$  and  $K$  only.

*Proof.* By Theorem 5.18 there exists a matrix  $Y$  such that  $A(0) = YBY^{-1}$ , where  $B$  is a block-diagonal matrix with blocks  $B_j$  of size  $n_j$ ,  $\kappa(Y) \leq C(n, K)$ , and  $\kappa(B_j) \leq C(n, K)$ . Denote  $B(\phi) = Y^{-1}A(\phi)Y$ .

Denote by  $P_j(\phi)$  the sum of the total eigenprojections of the  $\lambda$ -groups for the matrix function  $B(\phi)$  corresponding to  $\sigma(B_j)$ . By construction

$$P_j(0) = \begin{pmatrix} 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & I & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 \end{pmatrix},$$

where  $I$  is the identity matrix of size  $n_j$  and its position corresponds to the position of  $B_j$  in the matrix  $B$ . By continuity the submatrix of  $P_j(\phi)$  taking the same position as  $B_j$  is non-degenerate in a neighborhood of  $\phi = 0$ . Therefore, for a given  $j$  the columns of  $P_j(\phi)$  corresponding to the block  $B_j$  form a basis in the sum of algebraic eigenspaces of  $B(\phi)$  corresponding to  $\sigma(B_j)$ .

Introduce the matrix function  $\tilde{S}(\phi)$  composed from these basis columns of all  $P_j(\phi)$ . By construction the matrix  $M(\phi) = \tilde{S}^{-1}(\phi)B(\phi)\tilde{S}(\phi)$  has the block-diagonal form (5.10), where  $M^{(j)}(0) = B_j$  and  $\tilde{S}(0) = I$ . The matrix  $\tilde{S}(\phi)$  is analytical in a neighborhood of  $\phi = 0$  since all of its columns are analytical by Lemma 5.16.

By continuity  $\det \tilde{S}(\phi) \neq 0$  in a neighborhood of  $\phi = 0$ . The elements of  $\tilde{S}^{-1}(\phi)$  are products of elements of  $\tilde{S}(\phi)$  divided by  $\det \tilde{S}(\phi)$ , so  $\tilde{S}^{-1}(\phi)$  is also holomorphic. Hence so is  $M(\phi)$ . Thus we get representation (5.10) with  $S(\phi) = Y\tilde{S}(\phi)$ .  $\square$

## 6. Spectral analysis

The purpose of this section is to show that the error estimate for an arbitrary smooth solution is equivalent to the error estimate for a single complex exponent. This allows us to reduce our further analysis to a single wave and thus to a finite-dimensional analysis.

### 6.1. The basics

Recall that  $\mathcal{A} = \{\alpha \in \mathbb{R}^d : T^* \alpha / \pi \in \mathbb{Q}^d\}$ . For  $\beta > 0$  denote  $\mathcal{A}_\beta = \{\phi \in \mathcal{A} : T^* \phi \in [-\beta, \beta]^d\}$ . Let  $l_2(\mathcal{A}_\pi, \mathbb{C}^{M^0})$  be the space of maps  $w$  of  $\mathcal{A}_\pi$  to  $\mathbb{C}^{M^0}$  such that  $\sum_\phi \|w_\phi\|^2 < \infty$  with the scalar product  $(w, w') = \sum_\phi (w_\phi, w'_\phi)$ . Let  $V_F \subset l_2(\mathcal{A}_\pi, \mathbb{C}^{M^0})$  be the space of maps  $w$  such that  $w(\phi) = 0$  for all but a finite number  $\phi \in \mathcal{A}_\pi$ . Clearly  $l_2(\mathcal{A}_\pi, \mathbb{C}^{M^0})$  is a Hilbert space, and  $V_F$  is its dense subspace. Also denote  $U = (T^*)^{-1}$ .

Let  $F : V_{per} \rightarrow V_F$  be the mapping taking each  $f \in V_{per}^N$  to

$$F[f](\phi) = \frac{1}{N^d} \sum_{\eta \in \{0, \dots, N-1\}^d} f_\eta \exp(-i\phi \cdot T\eta) \quad \text{for } \phi = 2\pi U \frac{\mathbf{k}}{N}, \quad \mathbf{k} \in \mathbb{Z}^d \cap [-N/2, N/2]^d, \quad (6.1)$$

and  $F[f](\phi) = 0$  otherwise. Note that  $f \in V_{per}^N$  also belongs to  $f \in V_{per}^{nN}$  for each  $n \in \mathbb{N}$ , but this does not affect the value  $F[f]$ . It is easy to see that  $F$  is a linear mapping. This mapping is invertible. Indeed, for  $w \in V_F$  let  $F^{-1}[w]$  be the sequence with components

$$(F^{-1}[w])_\eta = \sum_{\phi \in \mathcal{A}_\pi} w(\phi) \exp(i\phi \cdot T\eta). \quad (6.2)$$

Obviously,  $F^{-1}[w] \in V_{per}^N$ , where  $N$  is the product of all denominators of  $\phi$  such that  $w(\phi) \neq 0$  (by definition, the number of such  $\phi$  is finite). It is easy to see that  $F^{-1}[F[f]] = f$  and  $F[F^{-1}[w]] = w$ .

Denote

$$\delta_{\phi_1, \phi_2} = \begin{cases} 1, & \phi_1 = \phi_2; \\ 0, & \text{otherwise}; \end{cases} \quad \delta_{\phi_1, \phi_2}^{\text{mod}} = \sum_{\mathbf{c} \in \mathbb{Z}^d} \delta_{\phi_1, \phi_2 + 2\pi U \mathbf{c}} = \begin{cases} 1, & \exists \mathbf{c} \in \mathbb{Z}^d : \phi_1 = \phi_2 + 2\pi U \mathbf{c}; \\ 0, & \text{otherwise}. \end{cases}$$

**Lemma 6.1.** *If  $f \in V_{per}$  has components  $f_\eta = y \exp(i\phi \cdot T\eta)$  with some  $\phi \in \mathcal{A}$ ,  $y \in \mathbb{C}^{M^0}$ , then*

$$F[f](\phi') = y \delta_{\phi, \phi'}^{\text{mod}} \quad (6.3)$$

for each  $\phi' \in \mathcal{A}_\pi$ .

The proof is obvious.

**Lemma 6.2.** *For  $f, g \in V_{per}$  there holds  $(F[f], F[g]) = (f, g)$  and thus  $F$  is an isometry:  $\|F[f]\| = \|f\|$ .*

*Proof.* Let  $f \in V_{per}^{N_1}$ ,  $g \in V_{per}^{N_2}$ , then  $f, g \in V_{per}^N$  with  $N = N_1 N_2$ . We have

$$\begin{aligned} (F[f], F[g]) &= \sum_{\mathbf{k} \in \mathbb{Z}^d \cap [-N/2, N/2]^d} (F[f](2\pi U \mathbf{k}/N), F[g](2\pi U \mathbf{k}/N)) = \\ &= \frac{1}{N^{2d}} \sum_{\eta \in \{0, \dots, N-1\}^d} \sum_{\zeta \in \{0, \dots, N-1\}^d} (f_\eta, g_\zeta) \sum_{\mathbf{k} \in \mathbb{Z}^d \cap [-N/2, N/2]^d} \exp\left(2\pi i \frac{\mathbf{k} \cdot (\eta - \zeta)}{N}\right). \end{aligned}$$

The sum over  $\mathbf{k}$  is equal to zero for  $\eta \neq \zeta$  and to  $N^d$  for  $\eta = \zeta$ , thus

$$(F[f], F[g]) = \frac{1}{N^d} \sum_{\eta \in \{0, \dots, N-1\}^d} (f_\eta, g_\eta) = (f, g).$$

□

**Corollary 6.3.** *The space  $V_{per}$  is incomplete.*

*Proof.* Since  $F$  is a bijective isometry from  $V_{per}$  to  $V_F$ , and the latter is incomplete, so is the former.  $\square$

For each  $\phi \in \mathbb{C}^d$  put by definition

$$Z(\phi) = \sum_{\eta \in S \subset \mathbb{Z}^d} Z_\eta \exp(i\phi \cdot T\eta), \quad L(\phi) = \sum_{\eta \in S \subset \mathbb{Z}^d} L_\eta \exp(i\phi \cdot T\eta). \quad (6.4)$$

370 Functions  $Z(\phi)$  and  $L(\phi)$  are periodic with the periodic cell  $(T^*)^{-1}[-\pi, \pi]^d$ .

**Lemma 6.4.** *Let  $Z$  and  $L$  be given by (2.9). Then for  $f \in V_{per}$  there holds*

$$F[Zf](\phi) = Z(\phi)F[f](\phi), \quad F[Lf](\phi) = L(\phi)F[f](\phi), \quad (6.5)$$

$$F[\exp(-Z^{-1}Lt/h)f](\phi) = \exp(-Z^{-1}(\phi)L(\phi)t/h)F[f](\phi). \quad (6.6)$$

*Proof.* By linearity it is enough to check this for  $f_\eta = \exp(i\phi \cdot T\eta)y$ , with some  $\phi \in \mathcal{A}_\pi$ ,  $y \in \mathbb{C}^{M^0}$ . Then

$$(Zf)_\eta = \sum_{\zeta \in S} Z_\zeta \exp(i\phi \cdot T(\eta + \zeta))y = \exp(i\phi \cdot T\eta)Z(\phi)y,$$

so  $F[Zf](\phi') = Z(\phi)y\delta_{\phi, \phi'} = Z(\phi)F[f](\phi')$ . The proof of the second identity is similar. To prove (6.6), note that  $u(t) = \exp(-Z^{-1}Lt/h)f$  is the solution of (2.10) with the initial data  $f$ . Applying  $F$  to (2.10) and using (6.5) we get

$$Z(\phi) \frac{d}{dt} F[u(t)](\phi) + \frac{1}{h} L(\phi) F[u(t)](\phi) = 0, \quad F[u(0)](\phi) = F[f](\phi).$$

Solving this equation we get (6.6).  $\square$

## 6.2. Mappings to the mesh space

For  $\phi \in \mathbb{R}^d$  there holds  $e^{i\phi \cdot r} \in L_{2,per}(\mathbb{R}^d)$  iff  $\phi \in \mathcal{A}$ . For a homogeneous mapping  $\Pi_h$  and  $\phi \in \mathcal{A}$  put by definition

$$v(\phi, \Pi_h) = (\Pi_1 e^{i\phi \cdot r})_0, \quad (6.7)$$

where the subscript 0 means the block component at  $\eta = 0$ . If  $\Pi_h$  is a local mapping, then (6.7) defines  $v(\phi, \Pi_h)$  for each  $\phi \in \mathbb{C}^d$ , and  $v(\phi, \Pi_h)$  is a holomorphic function.

**Lemma 6.5.** *Let  $\phi, \psi \in \mathcal{A}$ ,  $\phi' \in \mathcal{A}_\pi$ , and  $\Pi_h$  be a homogeneous mapping. Then for  $\alpha = \phi/h$ ,  $\beta = \psi/h$  there holds  $\Pi_h e^{i\alpha \cdot r} \in V_{per}$  and*

$$F[\Pi_h e^{i\alpha \cdot r}](\phi') = \delta_{\phi, \phi'}^{\text{mod}} v(\phi, \Pi_h), \quad (6.8)$$

$$(\Pi_h e^{i\alpha \cdot r}, \Pi_h e^{i\beta \cdot r}) = \delta_{\phi, \psi}^{\text{mod}} (v(\phi, \Pi_h), v(\psi, \Pi_h)). \quad (6.9)$$

*Proof.* By definition,

$$(\Pi_h e^{i\alpha \cdot r})_\eta = (\Pi_h e^{i\alpha \cdot (r+hT\eta)})_0 = e^{i\alpha \cdot hT\eta} (\Pi_h e^{i\alpha \cdot r})_0 = e^{i\phi \cdot T\eta} v(\phi, \Pi_h).$$

Applying Lemma 6.1 we get (6.8). By Lemma 6.2

$$\begin{aligned} (\Pi_h e^{i\alpha \cdot r}, \Pi_h e^{i\beta \cdot r}) &= \sum_{\psi' \in \mathcal{A}_\pi} (F[\Pi_h e^{i\alpha \cdot r}](\psi'), F[\Pi_h e^{i\beta \cdot r}](\psi')) = \\ &= \sum_{\psi' \in \mathcal{A}_\pi} \delta_{\phi, \psi'}^{\text{mod}} \delta_{\psi, \psi'}^{\text{mod}} (v(\phi, \Pi_h), v(\psi, \Pi_h)) = \delta_{\phi, \psi}^{\text{mod}} (v(\phi, \Pi_h), v(\psi, \Pi_h)). \end{aligned}$$

**Lemma 6.6.** *Let  $G \subset \mathcal{A}_\pi$  be symmetric with respect to the origin. Let  $v$  be a map of  $G$  to  $\mathbb{C}^{M^0}$  such that  $v(-\phi) = \overline{v(\phi)}$ . Then there exists a homogeneous mapping  $\Pi_h$  of  $L_{2,per}(\mathbb{R}^d)$  to  $V_{per}$  such that  $v(\phi) = v(\phi, \Pi_h)$  for  $\phi \in G$  and  $\|\Pi_h\| \leq \sup_G \|v(\phi)\|$ .*

*Proof.* Extend  $v$  to  $\mathcal{A}_\pi \setminus G$  by zero. Let  $f \in L_{2,per}(\mathbb{R}^d)$  have Fourier series  $\sum_{\alpha \in \mathcal{A}} f_\alpha \exp(i\alpha \cdot \mathbf{r})$ . Note that there is only finitely many nonzero  $f_\alpha$  in any bounded part of the frequency domain. Define  $w_f \in V_F$  by  $w_f(\phi) = v(\phi)f_{\phi/h}$ . Let  $\Pi_h$  is a linear mapping of  $L_{2,per}(\mathbb{R}^d)$  to  $V_{per}$  defined by  $\Pi_h f = F^{-1}[w_f]$ . Since

$$\|\Pi_h f\|^2 = \|F^{-1}[w_f]\|^2 = \|w_f\|^2 = \sum_{\phi \in G} \|v(\phi)f_{\phi/h}\|^2 \leq \sup_G \|v(\phi)\|^2 \sum_{\phi \in \mathcal{A}} |f_\phi|^2 = \sup_G \|v(\phi)\|^2 \|f\|^2,$$

there holds  $\|\Pi_h\| \leq \sup_G \|v(\phi)\|$ .

Now we check that  $\Pi_h$  is a homogeneous mapping. To prove (P3), consider a real-valued function  $f \in L_{2,per}(\mathbb{R}^d)$ . Then  $f_{-\alpha} = \overline{f_\alpha}$  and  $w_f(-\phi) = \overline{w_f(\phi)}$ . From (6.2),  $F^{-1}[w_f]$  is real-valued, i. e. (P3) is proved. To prove (P1) and (P2), it is enough to consider a single exponent  $f(\mathbf{r}) = \exp(i\alpha \cdot \mathbf{r})$ ,  $\alpha \in \mathcal{A}$ . Clearly,

$$(\Pi_h e^{i\alpha \cdot \mathbf{r}})_\eta = \begin{cases} v(\alpha h) \exp(i\alpha h \cdot T\eta), & \alpha h \in G, \\ 0, & \alpha h \notin G. \end{cases} \quad (6.10)$$

For

$$g(\mathbf{r}) = f(\mathbf{r} + hT\zeta) = \exp(i\alpha \cdot hT\zeta)f(\mathbf{r}),$$

we have

$$(\Pi_h g)_\eta = \exp(i\alpha \cdot hT\zeta) \exp(ih\alpha \cdot T\eta) v(h\alpha) = \exp(ih\alpha \cdot T(\eta + \zeta)) v(h\alpha) = (\Pi_h f)_{\eta+\zeta}.$$

The property (P1) is proved. For  $g(\mathbf{r}) = f(\mathbf{r}/h) = \exp(i\alpha \cdot \mathbf{r}/h)$  we have

$$(\Pi_h g)_\eta = \exp(ih\alpha \cdot T\eta/h) v(h\alpha/h) = (\Pi_1 f)_\eta.$$

380 The property (P2) is proved.

The property  $v(\phi, \Pi_h) = v(\phi)$  follows directly from (6.10). □

**Lemma 6.7.** *Let  $\Pi_h$  and  $\mathcal{P}_h$  be local mappings with the kernels  $\mu_\xi$  and  $\hat{\mu}_\xi$ , correspondingly, and  $\langle \hat{\mu}_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ . Suppose  $W(\phi)$  maps a neighborhood of zero in  $\mathbb{R}^d$  to  $\mathbb{C}^{M^0}$ , has  $q+1$  continuous derivatives at  $\phi = 0$ , satisfies  $W(-\phi) = \overline{W(\phi)}$ , and  $\|v(\phi, \Pi_h) - W(\phi)\| = O(|\phi|^p)$  as  $\phi \rightarrow 0$ , where  $p, q \in \mathbb{N}$ . Then there exist real-valued diagonal matrices  $\mathfrak{C}^{(m)}$ ,  $p \leq |\mathbf{m}| \leq q$ , such that for the mapping  $\Pi_h^{(p,q)}$  given by (2.12) there holds*

$$v(\phi, \Pi_h^{(p,q)}) = W(\phi) + O(|\phi|^{q+1}). \quad (6.11)$$

*Proof.* Let the diagonal matrix  $\mathfrak{C}(\phi)$  be defined in a neighborhood of  $\phi = 0$  by

$$W(\phi) - v(\phi, \Pi_h) = \mathfrak{C}(\phi)v(\phi, \mathcal{P}_h). \quad (6.12)$$

Since  $(v(\phi, \mathcal{P}_h))_\xi \rightarrow \langle \hat{\mu}_\xi, 1 \rangle \neq 0$  as  $\phi \rightarrow 0$ , then in a neighborhood of  $\phi = 0$  the matrix function  $\mathfrak{C}(\phi)$  is well-defined and  $\|\mathfrak{C}(\phi)\| = O(|\phi|^p)$ . Put

$$\mathfrak{C}^{(m)} = \frac{1}{\mathbf{m}!} \frac{1}{i^{|\mathbf{m}|}} D^{\mathbf{m}} \mathfrak{C}(\phi)|_{\phi=0}. \quad (6.13)$$

Taking complex conjugation from (6.12) we get

$$W(-\phi) - v(-\phi, \Pi_h) = \overline{\mathfrak{C}(\phi)} v(-\phi, \mathcal{P}_h),$$

thus  $\overline{\mathfrak{C}(\phi)} = \mathfrak{C}(-\phi)$  and

$$\overline{\mathfrak{C}^{(\mathbf{m})}} = \frac{1}{\mathbf{m}!} \frac{1}{(-i)^{|\mathbf{m}|}} D^{\mathbf{m}}(\overline{\mathfrak{C}(\phi)}) \Big|_{\phi=0} = \frac{1}{\mathbf{m}!} \frac{1}{(-i)^{|\mathbf{m}|}} D^{\mathbf{m}}(\mathfrak{C}(-\phi)) \Big|_{\phi=0} = \mathfrak{C}^{(\mathbf{m})},$$

so  $\mathfrak{C}^{(\mathbf{m})}$  is real-valued. By definition

$$v(\phi, \Pi_h^{(p,q)}) = v(\phi, \Pi_h) + \left[ \sum_{p \leq |\mathbf{m}| \leq q} \mathfrak{C}^{(\mathbf{m})}(i\phi)^{\mathbf{m}} \right] v(\phi, \mathcal{P}_h).$$

Since  $\|\mathfrak{C}(\phi)\| = O(|\phi|^p)$ , there holds  $\mathfrak{C}^{(\mathbf{m})} = 0$  whenever  $|\mathbf{m}| < p$ . Thus the sum inside the brackets is the Taylor polynomial of the function  $\mathfrak{C}(\phi)$  of order  $q$  and

$$v(\phi, \Pi_h^{(p,q)}) = v(\phi, \Pi_h) + (\mathfrak{C}(\phi) + O(|\phi|^{q+1}))v(\phi, \mathcal{P}_h).$$

Combining this with (6.12) we get (6.11).  $\square$

**Lemma 6.8.** *Let  $\Pi_h$  and  $\mathcal{P}_h$  be local mappings with the kernels  $\mu_\xi$  and  $\hat{\mu}_\xi$ , correspondingly, and  $\langle \hat{\mu}_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ . Let  $\mathbf{e} \in \hat{\Omega}$ . Let  $W(\psi)$  have  $q+1$  continuous derivatives at  $\psi = 0$ ,  $W(-\psi) = \overline{W(\psi)}$ , and  $\|v(\psi\mathbf{e}, \Pi_h) - W(\psi)\| = O(|\psi|^p)$  as  $\psi \rightarrow 0$ , where  $p, q \in \mathbb{N}$ . Then there exist real-valued diagonal matrices  $\mathfrak{C}_e^{(\mathbf{m})}$ ,  $\mathbf{m} = p, \dots, q$ , such that for the mapping  $\Pi_{h,\mathbf{e}}^{(p,q)}$  given by (2.13) there holds*

$$v(\psi\mathbf{e}, \Pi_{h,\mathbf{e}}^{(p,q)}) = w(\psi) + O(|\psi|^{q+1}).$$

The proof is similar to the previous Lemma.

### 6.3. Truncation error, stability, and accuracy

We need the following auxiliary result.

**Lemma 6.9.** *Let  $\Pi_h$  be a homogeneous mapping of  $H_{per}^q(\mathbb{R}^d)$  (or  $C_{per}^q(\mathbb{R}^d)$ ) to  $V_{per}$ . Then the operator  $F_h^1$  defined as  $F_h^1 f = h\epsilon_h(f, \Pi_h)$  is a homogeneous mapping of  $H_{per}^{q+1}(\mathbb{R}^d)$  (or  $C_{per}^{q+1}(\mathbb{R}^d)$ ). For each  $t \geq 0$  the operator  $F_h^2$  defined as  $F_h^2 f = \varepsilon_h(t, f, \Pi_h)$  is a homogeneous mapping of  $H_{per}^q(\mathbb{R}^d)$  (or  $C_{per}^q(\mathbb{R}^d)$ ) to  $V_{per}$ . Besides,*

$$\|F_h^1\|_{(q+1,h,*)} \leq \|\Pi_h\|_{(q,h,*)} (\|Z\|\|\omega\| + \|L\|), \quad \|F_h^2\|_{(q,h,*)} \leq (K+1)\|\Pi_h\|_{(q,h,*)},$$

where  $K$  is the stability constant of the scheme, and  $*$  means either nothing or infinity depending on the case.

*Proof.* By the definitions of  $\epsilon_h(f, \Pi_h)$  (see (2.14)) and  $\varepsilon_h(t, f, \Pi_h)$  (see (2.15)) it immediately follows that the mappings  $F_h^1$  and  $F_h^2$  are homogeneous.

For each  $v_0 \in H_{per}^{q+1}(\mathbb{R}^d)$  (or  $C_{per}^{q+1}(\mathbb{R}^d)$ ) by Lemma 5.7 there holds

$$\|F_h^1 v_0\| = \|h\epsilon_h(v_0, \Pi_h)\| = \|-hZ\Pi_h(\omega \cdot \nabla v_0) + L\Pi_h v_0\| \leq (\|Z\|\|\omega\| + \|L\|)\|\Pi_h\|_{(q,h,*)}\|v_0\|_{(q+1,h,*)}.$$

Now let  $v_0 \in H_{per}^q(\mathbb{R}^d)$  (or  $C_{per}^q(\mathbb{R}^d)$ ),  $v(t, \mathbf{r}) = v_0(\mathbf{r} - \omega t)$ , and  $u(t)$  be the solution of (2.8) with the initial data  $u(0) = \Pi_h v_0$ . Using  $\|v_0\| = \|v(t, \cdot)\|$  we get

$$\|F_h^2 v_0\| = \|\varepsilon_h(t, v_0, \Pi_h)\| = \|u(t) - \Pi_h v(t, \cdot)\| \leq K\|\Pi_h v_0\| + \|\Pi_h v(t, \cdot)\| \leq (K+1)\|\Pi_h\|_{(q,h,*)}\|v_0\|_{(q,h,*)},$$

where we have used the stability condition.  $\square$

By definition, put

$$A(\phi) = -Z^{-1}(\phi)L(\phi) + i\omega \cdot \phi I, \quad (6.14)$$

$$\hat{\epsilon}(\phi, \Pi_h) = A(\phi)v(\phi, \Pi_h) \quad (6.15)$$

$$\hat{\epsilon}(\phi, \nu, \Pi_h) = \left( e^{\nu A(\phi)} - I \right) v(\phi, \Pi_h). \quad (6.16)$$

390 where  $I$  is the identity matrix of size  $|M^0|$ ,  $Z(\phi)$  and  $L(\phi)$  are given by (6.4), and  $v(\phi, \Pi_h)$  is given by (6.7).

**Lemma 6.10.** *Let  $\Pi_h$  be a homogeneous mapping,  $\alpha \in \mathcal{A}$ ,  $\phi = \alpha h$ . Then there holds*

$$\|\epsilon_h(e^{i\alpha \cdot r}, \Pi_h)\| = \frac{1}{h} \|Z(\phi) \hat{\epsilon}(\phi, \Pi_h)\|, \quad (6.17)$$

$$\|\epsilon_h(t, e^{i\alpha \cdot r}, \Pi_h)\| = \|\hat{\epsilon}(\phi, t/h, \Pi_h)\|. \quad (6.18)$$

*Proof.* First we obtain the spectral representation of the truncation error. By (6.5) and (6.8), for  $\phi' \in \mathcal{A}_\pi$ ,

$$\begin{aligned} F[\epsilon_h(e^{i\alpha \cdot r}, \Pi_h)](\phi') &= \left( -(\omega \cdot i\alpha)Z(\phi') + \frac{1}{h}L(\phi') \right) F[\Pi_h e^{i\alpha \cdot r}](\phi') = \\ &= -\frac{1}{h}Z(\phi') \left( (\omega \cdot i\alpha h)I - Z^{-1}(\phi')L(\phi') \right) \delta_{\phi, \phi'}^{\text{mod}} v(\phi, \Pi_h) = -\frac{1}{h}Z(\phi) \hat{\epsilon}(\phi, \Pi_h) \delta_{\phi, \phi'}^{\text{mod}}. \end{aligned} \quad (6.19)$$

The last identity is due to the  $2\pi U\mathbb{Z}^d$ -periodicity of  $L$  and  $Z$ . Using Lemma 6.2 we get (6.17).

Now we move to the solution error. By definition,

$$\epsilon_h(t, e^{i\alpha \cdot r}, \Pi_h) = u(t) - e^{-i\alpha \cdot \omega t} \Pi_h e^{i\alpha \cdot r},$$

where  $u(t) = \exp(-Z^{-1}Lt/h) \Pi_h e^{i\alpha \cdot r}$  (see (2.15)). By (6.6), for  $\phi' \in \mathcal{A}_\pi$  there holds

$$F[u(t)](\phi') = \exp\left(-Z^{-1}(\phi')L(\phi')\frac{t}{h}\right) F[\Pi_h e^{i\alpha \cdot r}](\phi').$$

Thus by (6.8) we get

$$F[\epsilon_h(t, e^{i\alpha \cdot r}, \Pi_h)](\phi') = \left[ \exp\left(-Z^{-1}(\phi')L(\phi')\frac{t}{h}\right) - \exp(-i\alpha \cdot \omega t) \right] \delta_{\phi, \phi'}^{\text{mod}} v(\phi, \Pi_h).$$

By the  $2\pi U\mathbb{Z}^d$ -periodicity of  $Z(\phi)$  and  $L(\phi)$  we replace  $\phi'$  by  $\phi$  in their arguments and finally obtain

$$F[\epsilon_h(t, e^{i\alpha \cdot r}, \Pi_h)](\phi') = \exp(-i\alpha \cdot \omega t) \hat{\epsilon}(\phi, t/h, \Pi_h) \delta_{\phi, \phi'}^{\text{mod}}.$$

Now Lemma 6.2 yields (6.18). □

**Lemma 6.11.** *Consider a scheme of the form (2.8) and a homogeneous mapping  $\Pi_h$ . Let the error estimate (2.17) hold. Then for  $\phi \in \mathcal{A}$  there holds*

$$\|\hat{\epsilon}(\phi, \nu, \Pi_h)\| \leq C_1 |\phi|^P + C_2 \nu |\phi|^{Q+1} + C_3 |\phi|^R, \quad (6.20)$$

where  $C_1, C_2, C_3$  are the same as in (2.17). If  $\Pi_h$  is a local mapping, then (6.20) holds for each  $\phi \in \mathbb{R}^d$ .

*Proof.* Using (2.17) for  $v_0 = e^{i\alpha \cdot r}$ ,  $\alpha = \phi/h$ , we get

$$\|\epsilon_h(t, e^{i\alpha \cdot r}, \Pi_h)\| \leq C_1 h^P |\alpha|^P + C_2 t h^Q |\alpha|^{Q+1} + C_3 h^R |\alpha|^R,$$

and it remains to use Lemma 6.10. If  $\Pi_h$  is a local mapping, then (6.20) extends to  $\phi \in \mathbb{R}^d$  by continuity. □

**Lemma 6.12.** *The optimal value of the order of the truncation error in the sense of a local mapping is an integer.*

*Proof.* By Lemma 6.10 this order is less by one then the order of smallness of  $\hat{\epsilon}(\phi, \Pi_h)$  as  $\phi$  tends to zero, which is integer since  $\hat{\epsilon}(\phi, \Pi_h)$  is a holomorphic function of  $\phi$ .  $\square$

**Lemma 6.13.** *A scheme (2.8) is stable with the constant  $K$  if and only if for each  $\phi \in \mathcal{A}$  and each  $\nu > 0$  there holds*

$$\|\exp(A(\phi)\nu)\| \leq K. \quad (6.21)$$

*Proof.* For an operator  $Y : V_F \rightarrow V_F$  such that  $(Yw)(\phi) = y(\phi)w(\phi)$  there obviously holds  $\|Y\| = \sup_{\phi \in \mathcal{A}_\pi} \|y(\phi)\|$ . Using this fact for  $y(\phi) = \exp(-Z^{-1}(\phi)L(\phi)t/h)$  and Lemmas 6.2 and 6.4, we get

$$\left\| \exp\left(-Z^{-1}L\frac{t}{h}\right) \right\| = \sup_{\phi \in \mathcal{A}_\pi} \left\| \exp\left(-Z^{-1}(\phi)L(\phi)\frac{t}{h}\right) \right\|.$$

By continuity and  $2\pi U\mathbb{Z}^d$ -periodicity of  $Z^{-1}(\phi)$  and  $L(\phi)$ , and finally by the definition of  $A(\phi)$  we get

$$\begin{aligned} \left\| \exp\left(-Z^{-1}L\frac{t}{h}\right) \right\| &= \sup_{\phi \in \mathcal{A}} \left\| \exp\left(-Z^{-1}(\phi)L(\phi)\frac{t}{h}\right) \right\| = \\ &= \sup_{\phi \in \mathcal{A}} \left\| \exp\left(A(\phi)\frac{t}{h}\right) \exp(-i\omega \cdot \phi I) \right\| = \sup_{\phi \in \mathcal{A}} \left\| \exp\left(A(\phi)\frac{t}{h}\right) \right\|. \end{aligned}$$

Taking the supremum over  $t, h > 0$ , we get

$$\sup_{t, h > 0} \left\| \exp\left(-Z^{-1}L\frac{t}{h}\right) \right\| = \sup_{\nu > 0} \sup_{\phi \in \mathcal{A}} \|\exp(A(\phi)\nu)\|.$$

It remains to recall the definition of the stability constant of a scheme.  $\square$

#### 6.4. From a sine wave to an arbitrary smooth solution

The section is a toolbox for extending various estimates from single Fourier modes to arbitrary smooth solutions.

**Lemma 6.14.** *Let  $F_h$  be a homogeneous mapping of  $H_{per}^r(\mathbb{R}^d)$  to  $V_{per}$ ,  $r \geq 0$ . Let  $p(s) = \sum_{j=1}^J c_j s^{p_j}$  with some  $J \in \mathbb{N}$  and  $c_j, p_j \geq 0$ . Suppose there exists  $0 < \beta \leq \pi$  such that for each  $h$  and each  $\alpha \in \mathcal{A}_{\beta/h}$  there holds*

$$\|F_h e^{i\alpha \cdot r}\| \leq p(|\alpha|).$$

*Then for each  $w \in H_{per}^{\max\{r, p_1, \dots, p_J\}}(\mathbb{R}^d)$  there holds*

$$\|F_h w\| \leq \sum_{j=1}^J c_j \|\nabla^{p_j} w\| + 2h^r c^r \|F_h\|_{(r, h)} \|\nabla^r w\|, \quad (6.22)$$

where  $c = 1 + \beta^{-1} \|T^*\|_2$  and  $\|F_h\|_{(r, h)} = \sup \|F_h f\| / \|f\|_{(r, h)}$ .

*Proof.* Consider a function  $w \in H_{per}^{\max\{r, p_1, \dots, p_J\}}(\mathbb{R}^d)$ . Since  $w$  is periodic it has the Fourier representation (2.4):

$$w = \sum_{\alpha \in \mathcal{A}} w_\alpha \exp(i\alpha \cdot r).$$

Recall that  $\mathcal{A}_{\beta/h} = \{\alpha \in \mathcal{A} : T^*\alpha \in [-\beta/h, \beta/h]^d\}$ . By definition, put

$$S_N(r) = \sum_{\alpha \in \mathcal{A}_{\beta/h}} w_\alpha \exp(i\alpha \cdot r).$$



Obviously,  $F_h w = F_h S_N + F_h(w - S_N)$ . By Lemma 6.5 for  $\alpha, \alpha' \in \mathcal{A}_{\beta/h}$ ,  $\alpha \neq \alpha'$  there holds  $(F_h e^{i\alpha \cdot r}, F_h e^{i\alpha' \cdot r}) = 0$ . By assumption and using the Minkowski inequality for sums (the triangle inequality for  $l_2$ ) we get

$$\begin{aligned} \|F_h S_N\|^2 &= \sum_{\alpha \in \mathcal{A}_{\beta/h}} |w_\alpha|^2 \|F_h \exp(i\alpha \cdot r)\|^2 \leq \sum_{\alpha \in \mathcal{A}_{\beta/h}} |w_\alpha|^2 (p(|\alpha|))^2 = \\ &= \sum_{\alpha \in \mathcal{A}_{\beta/h}} \left( \sum_{j=1}^J c_j |w_\alpha| |\alpha|^{p_j} \right)^2 \leq \left( \sum_{j=1}^J c_j \left( \sum_{\alpha \in \mathcal{A}_{\beta/h}} |w_\alpha|^2 |\alpha|^{2p_j} \right)^{1/2} \right)^2. \end{aligned}$$

Extending the limits of the sum to  $\alpha \in \mathcal{A}$  we obtain

$$\|F_h S_N\| \leq \sum_{j=1}^J c_j \left( \sum_{\alpha \in \mathcal{A}} |w_\alpha|^2 |\alpha|^{2p_j} \right)^{1/2} = \sum_{j=1}^J c_j \|\nabla^{p_j} w\|. \quad (6.23)$$

Now consider  $F_h(w - S_N)$ . By the definition of  $S_N$  we have

$$w - S_N = \sum_{\alpha \in \mathcal{A} \setminus \mathcal{A}_{\beta/h}} w_\alpha \exp(i\alpha \cdot r).$$

For  $\alpha \in \mathcal{A} \setminus \mathcal{A}_{\beta/h}$  we have

$$|\alpha|^{-1} \leq \|T^*\|_2 |T^* \alpha|^{-1} \leq \|T^*\|_2 \|T^* \alpha\|_\infty^{-1} \leq h\beta^{-1} \|T^*\|_2 \leq hc.$$

Thus

$$\|w - S_N\|^2 = \sum_{\alpha \in \mathcal{A} \setminus \mathcal{A}_{\beta/h}} |w_\alpha|^2 = \sum_{\alpha \in \mathcal{A} \setminus \mathcal{A}_{\beta/h}} \frac{|w_\alpha|^2 |\alpha|^{2r}}{|\alpha|^{2r}} \leq (hc)^{2r} \|\nabla^r w\|^2.$$

Taking into account that  $c \geq 1$  we obtain

$$\|w - S_N\|_{(r,h)}^2 \leq \|\nabla^r w\|^2 ((hc)^{2r} + h^{2r}) \leq 2\|\nabla^r w\|^2 (hc)^{2r}$$

and

$$\|F_h(w - S_N)\| \leq \sqrt{2} \|F\|_{(r,h)} \|\nabla^r w\| h^r c^r.$$

Combining this (6.23) we get (6.22).  $\square$

**Theorem 6.15.** Consider a scheme of the form (2.8) and a homogeneous mapping  $\Pi_h$ . [Let  $e \in \mathring{\Omega}$ .] Suppose there exist  $C \geq 0$ ,  $P_A \geq -1$ , and  $0 < \beta \leq \pi$  such that there holds

$$\|\hat{e}(\phi, \Pi_h)\| \leq C|\phi|^{P_A+1} \quad (6.24)$$

for each  $|\phi| \in \mathcal{A}_\beta$  [aligned with  $e$ ].

1) If  $\Pi_h$  is a bounded homogeneous mapping of  $H_{per}^q(\mathbb{R}^d)$  to  $V_{per}$ , then for each  $v_0 \in H_{per[\cdot, e]}^{r+1}(\mathbb{R}^d)$  and  $r \geq \max\{P_A, q\}$  there holds

$$\|\epsilon_h(v_0, \Pi_h)\| \leq Ch^{P_A} \|\nabla^{P_A+1} v_0\| + \tilde{C} h^r \|\nabla^{r+1} v_0\|, \quad (6.25)$$

where  $\tilde{C} = 4c^{r+1} \|\Pi_h\|_{(q,h)} (\|Z\| \|\omega\| + \|L\|)$ ,  $c = 1 + \beta^{-1} \|T^*\|_2$ .

2) If  $\Pi_h$  is a local mapping with  $\mu_\xi \in (C^q(G))^*$ ,  $q \in \mathbb{N} \cup \{0\}$ , then there exists  $\tilde{C}$  such that for each  $v_0 \in C_{per[\cdot, e]}^{r+1}(\mathbb{R}^d)$ ,  $r \in \mathbb{N}$ ,  $r \geq \max\{P_A, q\}$ , there holds

$$\|\epsilon_h(v_0, \Pi_h)\| \leq Ch^{P_A} \|\nabla^{P_A+1} v_0\| + \tilde{C} h^r \|\nabla^{r+1} v_0\|_\infty.$$

*Proof.* If  $e$  is specified, let  $\mathfrak{T} : L_{2,per}(\mathbb{R}^d) \rightarrow L_{2,per}(\mathbb{R}^d)$  be the projection operator taking each  $e^{i\alpha \cdot r}$  to  $e^{i\alpha \cdot r}$  if  $\alpha$  is aligned with  $e$  and to zero otherwise. If  $e$  is not specified,  $\mathfrak{T}$  is the identity operator. Clearly,  $\|\mathfrak{T}\| = 1$  and maps  $H_{per}^\sigma(\mathbb{R}^d)$  to  $H_{per}^\sigma(\mathbb{R}^d)$  for each  $\sigma$ .

To prove the first statement, put  $F_h v_0 = h\epsilon_h(v_0, \Pi_h)$ . By (6.24) and Lemma 6.10 we have  $\|F_h e^{i\alpha \cdot r}\| \leq C|h\alpha|^{P_A+1}$  for each  $\alpha \in \mathcal{A}_{\beta/h}$  [aligned with  $e$ ]. By Lemma 6.9 we have  $\|F_h\|_{(q+1,h)} \leq \|\Pi_h\|_{(q,h)}(\|Z\|\|\omega\| + \|L\|)$ . Since  $r \geq q$  then by Lemma 5.6 we have  $\|F_h\|_{(r+1,h)} \leq 2\|F_h\|_{(q+1,h)}$ . Since  $\|\mathfrak{T}\| = 1$  we have

$$\|F_h \mathfrak{T}\|_{(r+1,h)} \leq 2\|\Pi_h\|_{(q,h)}(\|Z\|\|\omega\| + \|L\|).$$

Using Theorem 6.14 for the mapping  $F_h \mathfrak{T}$  and  $p(s) = C(hs)^{P_A+1}$  we get (6.25).

Now we prove the second statement. By Lemma 5.9 there exists a local mapping  $\tilde{\Pi}_h$  with  $\tilde{\mu}_\xi \in L_2(\tilde{G})$ ,  $\tilde{G} = G + B_1(0)$ , such that for  $k \in \{r, r+1\}$  and each  $f \in C_{per}^k(\mathbb{R}^d)$  there holds  $\|\tilde{\Pi}_h f - \Pi_h f\| \leq ch^k \|\nabla^k f\|_\infty$ . By Lemma 2.3,  $\tilde{\Pi}_h$  is a bounded homogeneous mapping of  $L_{2,per}(\mathbb{R}^d)$  to  $V_{per}$ . Then for each  $f \in C_{per}^{r+1}(\mathbb{R}^d)$  there holds

$$\begin{aligned} \|\epsilon_h(f, \tilde{\Pi}_h) - \epsilon_h(f, \Pi_h)\| &\leq \|Z\| \|\tilde{\Pi}_h - \Pi_h\|(\omega \cdot \nabla) f + h^{-1} \|L\| \|(\tilde{\Pi}_h - \Pi_h) f\| \leq \\ &\leq ch^r \|Z\| \|\nabla^r(\omega \cdot \nabla) f\|_\infty + ch^r \|L\| \|\nabla^{r+1} f\|_\infty \leq ch^r (\|Z\| \|\omega\| + \|L\|) \|\nabla^{r+1} f\|_\infty. \end{aligned} \quad (6.26)$$

The last inequality in this chain is by Lemma 5.7. Particularly, this holds for  $f = e^{i\alpha \cdot r}$ , so by the triangle inequality for each  $\alpha \in \mathcal{A}_{\beta/h}$  [aligned with  $e$ ] we get

$$\|\epsilon_h(e^{i\alpha \cdot r}, \tilde{\Pi}_h)\| \leq C|\alpha|^{P_A+1} h^{P_A} + \tilde{c}|\alpha|^{r+1} h^r.$$

Put  $F_h v_0 = h\epsilon_h(v_0, \tilde{\Pi}_h)$ . Since  $\|\mathfrak{T}\| = 1$ , by Lemma 5.6 and Lemma 6.9 we have

$$\|F_h \mathfrak{T}\|_{(r+1,h)} \leq \|F_h\|_{(r+1,h)} \leq \|\tilde{\Pi}_h\|_{(q,h)}(\|Z\|\|\omega\| + \|L\|) \leq \|\tilde{\Pi}_h\|(\|Z\|\|\omega\| + \|L\|).$$

Using Lemma 6.14 for the mapping  $F_h \mathfrak{T}$  and  $p(s) = C(hs)^{P_A+1} + \tilde{c}(hs)^{r+1}$  we get

$$\|\epsilon_h(\mathfrak{T}v_0, \tilde{\Pi}_h)\| \leq Ch^{P_A} \|\nabla^{P_A+1} v_0\| + \tilde{C}h^r \|\nabla^{r+1} v_0\|.$$

From (6.26) by the triangle inequality for each  $v_0 \in C_{per[e]}^{r+1}(\mathbb{R}^d)$  we obtain

$$\|\epsilon_h(v_0, \Pi_h)\| \leq \|\epsilon_h(v_0, \tilde{\Pi}_h)\| + \hat{c}h^r \|\nabla^{r+1} v_0\|_\infty \leq Ch^{P_A} \|\nabla^{P_A+1} v_0\| + \hat{C}h^r \|\nabla^{r+1} v_0\|_\infty.$$

410 In the last inequality we used the inequality  $\|\nabla^{r+1} v_0\| \leq \|\nabla^{r+1} v_0\|_\infty$ . □

**Lemma 6.16.** Consider a scheme of the form (2.8) and a local mapping  $\Pi_h$  with  $\mu_\xi \in (W_2^q(G))^*$  (or  $(C^q(G))^*$ ),  $q \in \mathbb{N} \cup \{0\}$ . [Let  $e \in \mathring{\Omega}$ .] Suppose for  $\phi \in \mathcal{A}_\beta$  [aligned with  $e$ ] there holds  $\hat{e}(\phi, \Pi_h) = 0$ . Then for each  $v_0 \in H_{per[e]}^{q+1}(\mathbb{R}^d)$  (or  $C_{per[e]}^{q+1}(\mathbb{R}^d)$ ), there holds  $\epsilon_h(v_0, \Pi_h) = 0$ .

*Proof.* By assumption

$$y(\phi) := Z(\phi)\hat{e}(\phi, \Pi_h) = (-L(\phi) + i\omega \cdot \phi Z(\phi))v(\phi, \Pi_h) \equiv 0$$

415 for  $\phi \in \mathcal{A}_\beta$  [aligned with  $e$ ]. If  $e$  is not defined,  $y$  is an entire function of  $\phi$ , therefore  $y \equiv 0$  on  $\mathbb{C}^d$  by the uniqueness theorem for a holomorphic function of several complex variables and thus  $\hat{e}(\phi, \Pi_h) = Z^{-1}(\phi)y(\phi) = 0$  for  $\phi \in \mathbb{R}^d$ . If  $e$  is defined, the function  $y(e\psi)$  is an entire function of  $\psi$ , thus  $y(e\psi) = 0$  for each  $\psi \in \mathbb{C}$  and so  $\hat{e}(e\psi, \Pi_h) = 0$  for  $\psi \in \mathbb{R}$ . By Lemma 6.10 for each  $\alpha \in \mathcal{A}$  [aligned with  $e$ ] there holds  $\epsilon_h(e^{i\alpha \cdot r}, \Pi_h) = 0$ . Since the linear span of  $e^{i\alpha \cdot r}$ ,  $\alpha \in \mathcal{A}$ , is dense in  $H_{per}^{r+1}(\mathbb{R}^d)$  and  $C_{per}^{r+1}(\mathbb{R}^d)$ , and the linear span of  $e^{i\alpha \cdot r}$ ,  $\alpha \in \mathcal{A}$  and aligned with  $e$ , is dense in  $H_{per,e}^{r+1}(\mathbb{R}^d)$  and  $C_{per,e}^{r+1}(\mathbb{R}^d)$ , the statement to prove is by boundedness of the linear operator  $\epsilon_h(\cdot, \Pi_h)$  in the corresponding space (see Lemma 6.9). □

**Theorem 6.17.** Consider a stable scheme of the form (2.8) and a homogeneous mapping  $\Pi_h$ . Let  $P, Q, C_1, C_2 \geq 0$ . [Let  $\mathbf{e} \in \tilde{\Omega}$ .] Assume for all  $\phi \in \mathcal{A}_\beta$ ,  $\beta \in (0, \pi]$ , [aligned with  $\mathbf{e}$ ] and all  $\nu \geq 0$  there holds

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq C_1|\phi|^P + C_2\nu|\phi|^{Q+1}. \quad (6.27)$$

1) If  $\Pi_h$  is a bounded homogeneous mapping of  $H_{per}^q(\mathbb{R}^d)$  to  $V_{per}$  then there exists  $\tilde{C}$  such that for each initial data  $v_0 \in H_{per[\mathbf{e}]}^r(\mathbb{R}^d)$ ,  $r \geq \max\{Q+1, P, q\}$ , there holds

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq C_1 h^P \|\nabla^P v_0\| + C_2 t h^Q \|\nabla^{Q+1} v_0\| + \tilde{C} h^r \|\nabla^r v_0\|. \quad (6.28)$$

2) If  $\Pi_h$  is a local mapping with kernel  $\mu_\xi \in (C^q(G))^*$ ,  $q \in \mathbb{N} \cup \{0\}$ , then there exists  $\tilde{C}$  such that for each initial data  $v_0 \in C_{per[\mathbf{e}]}^r(\mathbb{R}^d)$ ,  $r \geq \max\{\lceil Q \rceil + 1, \lceil P \rceil, q\}$ , there holds

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq C_1 h^P \|\nabla^P v_0\| + C_2 t h^Q \|\nabla^{Q+1} v_0\| + \tilde{C} h^r \|\nabla^r v_0\|_\infty. \quad (6.29)$$

For  $Q = +\infty$  the terms with  $C_2$  vanish and the expression for  $r$  is replaced by  $r = \max\{P, q\}$  (or  $r = \max\{\lceil P \rceil, q\}$ ). The constants  $\tilde{C}$  in 1) and 2) depend only on  $\Pi_h, r, T, \beta, K$ , and the norm on  $\mathbb{C}^{M^0}$ .

*Proof.* By (6.18) for each  $t > 0$ , each  $h$ , and each  $\alpha \in \mathcal{A}_{\beta/h}$  [aligned with  $\mathbf{e}$ ] there holds

$$\|\varepsilon_h(t, e^{i\alpha \cdot \mathbf{r}}, \Pi_h)\| = \|\hat{\varepsilon}(h\alpha, t/h, \Pi_h)\| \leq C_1 h^P |\alpha|^P + C_2 t h^Q |\alpha|^{Q+1}.$$

Let  $\mathfrak{T}$  be the same as in Theorem 6.15.

To prove the first statement, put  $F_h v_0 = \varepsilon_h(t, \mathfrak{T} v_0, \Pi_h)$ . By Lemma 6.9

$$\|F_h\|_{(r,h)} \leq (K+1)\|\Pi_h\|_{(r,h)} \leq 2(K+1)\|\Pi_h\|_{(q,h)},$$

where  $K$  is the stability constant. Then by Lemma 6.14 with  $p(s) = C_1(hs)^P + C_2 t h^Q s^{Q+1}$  we have (6.28).

Now we prove the second statement. By Lemma 5.9 there exists a local mapping  $\tilde{\Pi}_h : L_{2,per}(\mathbb{R}^d) \rightarrow V_{per}$  with kernel  $\tilde{\mu}_\xi \in L_2(\tilde{G})$ ,  $\tilde{G} = G + B_1(0)$ , such that for each  $f \in C_{per}^r(\mathbb{R}^d)$  there holds  $\|\tilde{\Pi}_h f - \Pi_h f\| \leq c h^r \|\nabla^r f\|_\infty$  with  $c = c(d, r, \Pi_h)$ . In particular,  $\|\tilde{\Pi}_h e^{i\alpha \cdot \mathbf{r}} - \Pi_h e^{i\alpha \cdot \mathbf{r}}\| \leq c h^r |\alpha|^r$ . By construction and Lemma A.4,  $\|\tilde{\Pi}_h\|$  depends only on  $\mu_\xi$  and the norm on  $\mathbb{C}^{M^0}$ . By the triangle inequality and stability we have

$$\begin{aligned} \|\varepsilon_h(t, e^{i\alpha \cdot \mathbf{r}}, \tilde{\Pi}_h)\| &\leq \|\varepsilon_h(t, e^{i\alpha \cdot \mathbf{r}}, \Pi_h)\| + (K+1)\|(\tilde{\Pi}_h - \Pi_h)e^{i\alpha \cdot \mathbf{r}}\| \leq \\ &\leq C_1 |\alpha|^P h^P + C_2 |\alpha|^{Q+1} t h^Q + c(K+1) |\alpha|^r h^r. \end{aligned}$$

Put  $F_h v_0 = \varepsilon_h(t, \mathfrak{T} v_0, \tilde{\Pi}_h)$ . By Lemma 6.9 we have

$$\|F_h\|_{(r,h)} \leq (K+1)\|\tilde{\Pi}_h\|_{(r,h)} \leq (K+1)\|\tilde{\Pi}_h\|.$$

Then by Lemma 6.14 with  $p(s) = C_1(hs)^P + C_2 t h^Q s^{Q+1} + c(K+1)(hs)^r$  we have

$$\|\varepsilon_h(t, v_0, \tilde{\Pi}_h)\| \leq C_1 h^P \|\nabla^P v_0\| + C_2 t h^Q \|\nabla^{Q+1} v_0\| + \tilde{c} h^r \|\nabla^r v_0\|,$$

where  $\tilde{c}$  depends only on  $\Pi_h, r, T, \beta, K$ , and the norm on  $\mathbb{C}^{M^0}$ . By the triangle inequality and stability

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq \|\varepsilon_h(t, v_0, \tilde{\Pi}_h)\| + K\|\tilde{\Pi}_h v_0 - \Pi_h v_0\| + \|\tilde{\Pi}_h v(t, \cdot) - \Pi_h v(t, \cdot)\|$$

425 where  $v(t, \mathbf{r}) = v_0(t - \omega \cdot \mathbf{r})$ . From here, (6.29) easily follows.  $\square$

### 6.5. Miscellaneous facts

**Corollary 6.18.** Let  $\Pi_h$  be a bounded homogeneous mapping of  $H_{per}^s(\mathbb{R}^d)$  to  $V_{per}$  or a local mapping with  $\mu_\xi \in (C^s(G))^*$ . Let the scheme (2.8) be stable and possess the formal order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per}^r(\mathbb{R}^d)$  (or  $C_{per}^{[r]}(\mathbb{R}^d)$ ),  $r \geq \max\{Q+1, s\}$  (or  $r \geq \max\{P, s\}$  if  $Q = \infty$ ). Then for each  $p$  and  $q$  such that  $0 < p \leq P$  and  $p \leq q \leq Q$  the scheme possesses the formal order  $p$  and the long-time simulation order  $q$  on  $H_{per}^r(\mathbb{R}^d)$  (or  $C_{per}^{[r]}(\mathbb{R}^d)$ ). If  $q < Q = \infty$  we additionally assume  $q \leq r-1$ .

*Proof.* Let the scheme possess the formal order of accuracy  $P$  and the long-time simulation order  $Q$ . Then by Lemma 6.11 for all  $\phi \in \mathbb{R}^d$  and  $\nu \geq 0$  we have

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq C_1|\phi|^P + C_2\nu|\phi|^{Q+1} + C_3|\phi|^R.$$

Hence, in a neighborhood of  $\phi = 0$  we get

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq 2C_1|\phi|^P + C_2\nu|\phi|^{q+1}.$$

By Theorem 6.17 we obtain that the scheme possesses the formal order of accuracy  $p$  and the long-time simulation order  $q$  on  $H_{per}^r(\mathbb{R}^d)$  (or  $C_{per}^{[r]}(\mathbb{R}^d)$ ). For  $Q = \infty$  the proof remains valid if we drop the terms with  $C_2$ .  $\square$

**Corollary 6.19.** Consider a stable scheme of the form (2.8) and a bounded homogeneous mapping  $\Pi_h$  of  $H_{per}^q(\mathbb{R}^d)$  (or  $C_{per}^q(\mathbb{R}^d)$ ) to  $V_{per}$ . Then the following two statements are equivalent.

- (Y1) The scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$ .
- (Y2) For each  $e \in \dot{\Omega}$  the scheme possesses a error estimate of the form

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq C_1 h^P \|\nabla^P v_0\| + C_2 t h^Q \|\nabla^{Q+1} v_0\| + C_3 h^r \|\nabla^r v_0\|_* \quad (6.30)$$

for  $v_0 \in H_{per,e}^r(\mathbb{R}^d)$ ,  $r \geq \max\{Q+1, P, q\}$ , (or  $v_0 \in C_{per,e}^r(\mathbb{R}^d)$ ,  $r \geq \max\{[Q]+1, [P], q\}$ ), where  $C_1, C_2$ , and  $C_3$  are independent of  $e$ , and  $\|\cdot\|_*$  means either  $\|\cdot\|$  or  $\|\cdot\|_\infty$  depending on the case.

*Proof.* The implication (Y1) to (Y2) is obvious. Conversely, assume (Y2). Taking (6.30) for  $v_0 = \exp(i\phi \cdot \mathbf{r}/h)$ , by (6.18) we get

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq C_1|\phi|^P + C_2\nu|\phi|^{Q+1} + C_3|\phi|^r$$

for each  $\phi \in \mathcal{A}$ . Since  $r \geq P$ , for  $\phi \in \mathcal{A}_\pi$  there holds (6.27) with another  $C_1$ . It remains to use Theorem 6.17.  $\square$

**Lemma 6.20.** Let  $f(x) : \mathbb{R} \rightarrow \mathbb{C}^n$  be analytical at  $x = 0$  and satisfy  $\|f(1/k)\| \leq g(1/k)$  for each  $k \in \mathbb{N}$  where  $g(x) \in C^1(\mathbb{R})$  is a function such that  $g(0) = 0$  and  $g'(x) > 0$  for  $0 < x < 1$ . Then there exists  $\delta > 0$  such that there holds  $\|f(x)\| \leq Cg(2x)$  for  $0 < x < \delta$ , where  $C$  depends on the norm on  $\mathbb{C}^n$  only.

*Proof.* We interpret  $\mathbb{C}^n$  as  $\mathbb{R}^{2n}$  and denote by  $f_j$ ,  $j = 1, \dots, 2n$ , the components of  $f$ . Since  $f_j$  is analytical then there exists  $\delta_j$  such that  $f'_j(x) \neq 0$  for  $x \in (0, \delta_j)$ . Put  $\delta = \min\{\min\{\delta_j\}, 1\}/2$ . Then  $f_j$  is monotone on  $(0, \delta)$ ,  $f_j(0) = 0$  and so for any  $k \in \mathbb{N}$  we have

$$|f_j(x)| \leq |f_j(2^{-k})| \leq cg(2^{-k}) \leq cg(2x), \quad x \in (2^{-k-1}, 2^{-k}] \cap (0, \delta).$$

The inequality for  $\|f(x)\|$  easily follows.  $\square$

**Corollary 6.21.** Let  $\Pi_h$  be a local mapping with kernel  $\mu \in (W_2^q(G))^*$  (or  $(C^q(G))^*$ ). If for each  $\alpha \in \mathcal{A}$  and for each  $h$  there holds  $\|\epsilon_h(e^{i\alpha \cdot \mathbf{r}}, \Pi_h)\| \leq c(\alpha)h^{P_A}$ , then there holds

$$\|\epsilon_h(v_0, \Pi_h)\| \leq C_1 \|\nabla^{P_A+1} v_0\| h^{P_A} + C_2 \|\nabla^{\max\{P_A, q\}+1} v_0\|_* h^{\max\{P_A, q\}} \quad (6.31)$$

with some  $C_1, C_2 > 0$ , where  $\|\cdot\|_* \equiv \|\cdot\|$  if  $\mu \in (W_2^q(G))^*$  and  $\|\cdot\|_* \equiv \|\cdot\|_\infty$  otherwise.

*Proof.* Let  $e \in \mathring{\Omega}$ , then there exists  $\lambda > 0$  such that  $\lambda e \in \mathcal{A}$ . Let  $h$  be such that  $1/h \in \mathbb{N}$ . Then we get

$$\|\hat{e}(\lambda h e, \Pi_h)\| \leq h \|Z^{-1}\| \|\epsilon_h(e^{i\lambda e \cdot r}, \Pi_h)\| \leq c(e) \|Z^{-1}\| h^{P_A+1}.$$

Here the first inequality is by Lemma 6.10 and the second one is by assumption. For  $\phi = \lambda h e$ , where  $1/h \in \mathbb{N}$ , we have

$$\|\hat{e}(\phi, \Pi_h)\| \leq \frac{c(e)}{|\lambda e|^{P_A+1}} |\phi|^{P_A+1} := \hat{c}(e) |\phi|^{P_A+1}.$$

445 Since  $\hat{e}$  is holomorphic with respect to the first argument by Lemma 6.20 there holds  $\|\hat{e}(\phi, \Pi_h)\| \leq C \hat{c}(e) |2\phi|^{P_A+1}$  for each  $\phi$  aligned with  $e$  in a neighborhood of zero. By Lemma 5.3 there holds  $\|\hat{e}(\phi, \Pi_h)\| \leq \tilde{c} |\phi|^{P_A+1}$  in a neighborhood of  $\phi = 0$ . It remains to apply Theorem 6.15.  $\square$

**Proposition 6.22.** Consider a scheme of the form (2.8), a local mapping  $\Pi_h$  with the kernel  $\mu_\xi \in (W_2^q(G))^*$ , and  $e \in \mathring{\Omega}$ . Let  $0 < P \leq Q < \infty$  and  $R = \max\{Q+1, q\}$ . The scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per,e}^R(\mathbb{R}^d)$  iff it possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per,e}^R(\mathbb{R}^d)$  in the weak sense. 450

*Proof.* Let the scheme possess the formal order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per,e}^Q(\mathbb{R}^d)$  in the weak sense. Since  $e \in \mathring{\Omega}$ , there exists  $\lambda > 0$  such that  $\lambda e \in \mathcal{A}$ . For  $v_0 = e^{i\lambda e \cdot r}$  Definition 3 gives

$$\|\varepsilon_h(t, e^{i\lambda e \cdot r}, \Pi_h)\| \leq c_1 h^P + c_2 t h^Q.$$

Here  $c_1$  and  $c_2$  may depend on  $e$  but are independent of  $t$  and  $h$ . By Lemma 6.10 there holds

$$\|\hat{e}(\lambda h e, t/h, \Pi_h)\| = \|\varepsilon_h(t, e^{i\lambda e \cdot r}, \Pi_h)\| \leq c_1 h^P + c_2 t h^Q.$$

For  $\phi = \lambda h e$ , where  $1/h \in \mathbb{N}$ , and each  $\nu \geq 0$  we have

$$\|\hat{e}(\phi, \nu, \Pi_h)\| \leq |\phi|^P \frac{c_1}{|\lambda e|^P} + \nu |\phi|^{Q+1} \frac{c_2}{|\lambda e|^{Q+1}}.$$

Since  $\hat{e}$  is holomorphic with respect to the first argument by Lemma 6.20 there holds  $\|\hat{e}(\phi, \nu, \Pi_h)\| \leq \tilde{c}_1 |\phi|^P + \tilde{c}_2 \nu |\phi|^{Q+1}$  for each  $\phi$  aligned with  $e$  in a neighborhood of zero. By Theorem 6.17 the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per,e}^Q(\mathbb{R}^d)$ . The reverse implication is obvious.  $\square$

455 In particular, in the 1D case for local mappings Definition 3 is equivalent to Definition 4.

## 7. The 1D case

In this section we consider the one-dimensional case. The main result of this section (Theorem 7.7) can be obtained from the quasi-one-dimensional case considered in Section 9.3 (see Lemma 9.20). However, the 1D case allows us to use powerful analytical tools that provide a clearer understanding of the enhanced accuracy in the long-time simulation.

460 Throughout this section we consider a stable scheme of the form (2.8) and a local mapping  $\Pi_h$ . In 1D the set of the vectors  $\mathbf{a}_j$  contains only one vector  $\mathbf{a}_1$ , which has only one component. We assume it to be unit, i. e. the value  $h$  coincides with the mesh period, which is the natural definition for the schemes with several DOFs per cell. Hence,  $T$  and  $U$  are the identity operators on  $\mathbb{R}^1$ .

### 7.1. Matrix decomposition

465 Let  $\mathbb{R}^{n \times n}$  and  $\mathbb{C}^{n \times n}$  be the spaces of real and complex matrices of size  $n$ . Denote by  $\mathcal{A}(\mathbb{C}, \mathbb{R}^{n \times n})$  the set of functions  $A(\phi)$  from  $\mathbb{C}$  to  $\mathbb{C}^{n \times n}$  such that each element of the matrix  $A(\phi)$  is a holomorphic function at  $\phi = 0$  and for all  $i\phi \in \mathbb{R}$  in a neighborhood of  $\phi = 0$  there holds  $A(\phi) \in \mathbb{R}^{n \times n}$ .

We need the following result.

**Theorem 7.1** ([29]). Let  $A \in \mathcal{A}(\mathbb{C}, \mathbb{R}^{n \times n})$ . Suppose there exists  $K > 0$  such that for all  $\phi \in \mathbb{R}$  and  $\nu \geq 0$  there holds  $\|\exp(\nu A(\phi))\| \leq K$ . Then in a neighborhood of  $\phi = 0$  the matrix  $A(\phi)$  can be represented in the form

$$A(\phi) = S(\phi)M(\phi)S^{-1}(\phi), \quad (7.1)$$

$$M(\phi) = \begin{pmatrix} M_0(\phi) & 0 & \dots & 0 & 0 \\ 0 & \phi M_1(\phi) & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \phi^m M_m(\phi) & 0 \\ 0 & 0 & \dots & 0 & M_\infty(\phi) \equiv 0 \end{pmatrix}, \quad (7.2)$$

where  $S, M, S^{-1} \in \mathcal{A}(\mathbb{C}, \mathbb{R}^{n \times n})$ , the square matrices  $M_k(\phi)$ ,  $k \in \mathbb{N} \cup \{0, \infty\}$ , are non-degenerate for  $\phi = 0$  (except  $k = \infty$ ), and some of them are absent.

Throughout this section we denote by  $S(\phi)$ ,  $M_j(\phi)$ , and  $M(\phi)$  the matrices given by Theorem 7.1 for the matrix  $A(\phi)$  defined in (6.14).

Denote by  $\bar{\mathbb{N}}$  the set  $j \in \mathbb{N} \cup \{0, \infty\}$  such that the block  $\phi^j M_j(\phi)$  is present in the matrix  $M(\phi)$ .

The function  $v(\phi, \Pi_h)$  defined in (6.7) has the form  $v(\phi, \Pi_h) = (\Pi_1 e^{i\phi x})_0$ . Let  $v_j(\phi, \Pi_h)$ ,  $j \in \bar{\mathbb{N}}$ , be the components of  $S^{-1}(\phi)v(\phi, \Pi_h)$  corresponding to the blocks  $M_j(\phi)$ . Since  $v_j(\phi, \Pi_h)$  are holomorphic functions there exist  $p_j \in \mathbb{N} \cup \{0\}$  and  $c_j \in \mathbb{R} \setminus \{0\}$  such that in a neighbourhood of  $\phi = 0$  there holds

$$c_j |\phi|^{p_j} \leq \|v_j(\phi, \Pi_h)\| \leq 2c_j |\phi|^{p_j}. \quad (7.3)$$

If  $v_j(\phi) \equiv 0$  we put by definition  $p_j = \infty$  and  $c_j = 1$ . Put

$$\mathbb{N} = \{j \in \mathbb{N} \cup \{0\} : p_j < \infty\} \subseteq \bar{\mathbb{N}}.$$

This means that  $j \in \bar{\mathbb{N}} \setminus \mathbb{N}$  iff  $j = \infty$  or  $p_j = \infty$ .

Below we show that the values  $p_j$ ,  $j \in \bar{\mathbb{N}}$ , are responsible for the structure of the numerical error and define the order of the truncation error, the order of accuracy, and the long-time simulation order.

**Lemma 7.2.** *There exists  $j \in \bar{\mathbb{N}}$  such that  $p_j = 0$ .*

*Proof.* Assume the converse:  $p_j > 0$  for each  $j \in \bar{\mathbb{N}}$ . Then for each  $j$  there holds  $v_j(0) = 0$ , therefore  $(\Pi_1 1)_0 = v(0, \Pi_h) = 0$ . On the other hand, by definition for at least one  $\xi \in M^0$  there holds  $(\Pi_1 1)_{0,\xi} = \langle \mu_\xi, 1 \rangle \neq 0$ . This contradiction proves the lemma.  $\square$

## 7.2. The structure of the truncation error

**Lemma 7.3.** *Let  $\Pi_h$  be a local mapping with  $\mu_\xi \in (W_2^q(G))^*$  or  $(C^q(G))^*$ . If  $\mathbb{N} = \emptyset$  then the scheme is exact, i. e. for each  $v_0 \in H_{per}^{q+1}(\mathbb{R})$  (or  $C_{per}^{q+1}(\mathbb{R})$ ) there holds  $\epsilon_h(v_0, \Pi_h) = 0$ .*

*Proof.* Let  $j \in \bar{\mathbb{N}}$ . Since  $\mathbb{N} = \emptyset$ , either  $j = \infty$  or  $p_j = \infty$ . In a neighborhood of  $\phi = 0$  there holds  $M_j(\phi) = 0$  in the first case and  $v_j(\phi, \Pi_h) = 0$  in the second case. Thus for each  $j \in \bar{\mathbb{N}}$  there holds  $M_j(\phi)v_j(\phi, \Pi_h) \equiv 0$  and hence

$$\hat{\epsilon}(\phi, \Pi_h) \equiv A(\phi)v(\phi, \Pi_h) \equiv 0. \quad (7.4)$$

It remains to apply Lemma 6.16.  $\square$

**Lemma 7.4.** *Let  $\Pi_h$  be a local mapping. If  $\mathbb{N} \neq \emptyset$ , then the scheme possesses the truncation error of order*

$$P_A = \min_{j \in \mathbb{N}} \{p_j + j - 1\}. \quad (7.5)$$

*This value is optimal, i. e. the scheme does not possess the truncation error of order  $p > P_A$ .*

*Proof.* Since the matrices  $Z(\phi)$ ,  $S(\phi)$ , and  $M_j(\phi)$  for each  $j$  are holomorphic functions in a neighborhood of  $\phi = 0$ , for small  $|\phi|$  we have

$$\|\hat{\epsilon}(\phi, \Pi_h)\| \leq C \sum_{j \in \mathbb{N}} |\phi|^j \|v_j(\phi, \Pi_h)\| \leq \tilde{C} \max_{j \in \mathbb{N}} |\phi|^{j+p_j}. \quad (7.6)$$

By Theorem 6.15 the scheme possesses the truncation error of order  $P_A$ .

Conversely, let the scheme possess the truncation error of order  $p$ . Then by (6.17) we have  $\|\hat{\epsilon}(\phi, \Pi_h)\| \leq C|\phi|^{p+1}$ . Since  $S(\phi)$  is non-degenerate, for small  $|\phi|$  for each  $j \in \mathbb{N}$  there holds

$$\|\phi^j M_j(\phi) v_j(\phi, \Pi_h)\| \leq \tilde{c} |\phi|^{p+1}.$$

Since  $M_j(0)$  is non-degenerate, we have  $\|M_j^{-1}(\phi)\| \leq 2\|M_j^{-1}(0)\|$  in a neighborhood of  $\phi = 0$ . Hence,  $\|v_j(\phi, \Pi_h)\| \leq 2\tilde{c}\|M_j^{-1}(0)\| |\phi|^{p+1-j}$ . By (7.3) we get  $p_j \geq p + 1 - j$ . Since this holds for each  $j \in \mathbb{N}$ , we have  $p \leq P_A$ .  $\square$

### 490 7.3. The structure of the solution error

Now substitute the representation of  $A(\phi)$  given by Theorem 7.1 into (6.16). We have

$$\hat{\epsilon}(\phi, \nu, \Pi_h) = S(\phi) \left( e^{\nu M(\phi)} - I \right) S^{-1}(\phi) v(\phi, \Pi_h) = S(\phi) \begin{pmatrix} \vdots \\ E_j(\phi, \nu, \Pi_h) \\ \vdots \end{pmatrix},$$

where

$$E_j(\phi, \nu, \Pi_h) = \left( e^{\nu \phi^j M_j(\phi)} - I \right) v_j(\phi, \Pi_h) \quad (7.7)$$

and  $v_j(\phi, \Pi_h)$  is defined in Section 7.1. For  $j \in \mathbb{N} \setminus \mathbb{N}$  in a neighborhood of  $\phi = 0$  we have  $E_j(\phi, \nu, \Pi_h) \equiv 0$ . For  $j \in \mathbb{N}$  put  $Y = \nu \phi^j M_j(\phi)$ . By the stability condition (6.21) we have  $\|e^Y\| \leq K\kappa(\phi)$ , where  $\kappa(\phi)$  is the condition number of  $S(\phi)$ . By Lemma 5.11 from (7.3) we get

$$\|E_j(\phi, \nu, \Pi_h)\| \leq (K\kappa(\phi) + e) \min \{1, \nu |\phi|^j \|M_j(\phi)\|\} 2c_j |\phi|^{p_j}. \quad (7.8)$$

**Lemma 7.5.** *Let the scheme possess the formal order of accuracy  $P$  and the long-time simulation order  $Q \geq P$ . Then for each  $j \in \mathbb{N}$  there holds*

$$p_j \geq \min\{P, Q + 1 - j\}. \quad (7.9)$$

*Proof.* By Lemma 6.11 in a neighborhood of  $\phi = 0$  there holds

$$\|\hat{\epsilon}(\phi, \nu, \Pi_h)\| \leq \tilde{C}_1 |\phi|^P + C_2 \nu |\phi|^{Q+1}.$$

Since  $S(0)$  is non-degenerate, we get the estimate for each error component:

$$\|E_j(\phi, \nu, \Pi_h)\| \leq C'_1 |\phi|^P + C'_2 \nu |\phi|^{Q+1} \quad (7.10)$$

in a neighborhood of  $\phi = 0$ . Particularly, for  $\nu = |\phi|^{-j} \|M_j(0)\|^{-1}$  we get

$$\|E_j(\phi, |\phi|^{-j} \|M_j(0)\|^{-1}, \Pi_h)\| \leq C' |\phi|^{\min\{P, Q+1-j\}}.$$

From (7.7) we have

$$E_j(\phi, |\phi|^{-j} \|M_j(0)\|^{-1}, \Pi_h) = \left( \exp \left( \frac{M_j(\phi)}{\|M_j(0)\|} \right) - I \right) v_j(\phi, \Pi_h).$$

By Lemma 5.13 the matrix  $\exp(M_j(\phi)/\|M_j(0)\|) - I$  is non-degenerate for  $\phi = 0$ . Then it is non-degenerate in a neighborhood of  $\phi = 0$  and

$$\|v_j(\phi, \Pi_h)\| \leq \left\| \left( \exp(M_j(\phi)/\|M_j(0)\|) - I \right)^{-1} \right\| \|E_j(\phi, |\phi|^{-j} \|M_j(0)\|^{-1}, \Pi_h)\| \leq C'' |\phi|^{\min\{P, Q+1-j\}}.$$

Now (7.9) is by the definition of  $p_j$ .  $\square$

**Lemma 7.6.** *If  $p_0 = 0$  or  $p_1 = 0$ , then the scheme does not possess any order of accuracy. Otherwise the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  where*

$$P = \min_{j \in \mathbb{N}} (\max\{p_j + j - 1, p_j\}) = \min\{p_0, \min_{j \geq 1} (p_j + j - 1)\}, \quad (7.11)$$

$$Q = \min_{j \in \mathbb{N}: p_j < P} (p_j + j - 1). \quad (7.12)$$

*Minimum over the empty set is assumed to be  $+\infty$ . The values  $P$  and  $Q$  given by (7.11) and (7.12) are optimal in the sense of Definition 5.*

*Proof.* First consider the case when  $p_0 = 0$  or  $p_1 = 0$ . Assume that the scheme possess the formal order of accuracy  $P$  and the long-time simulation order  $Q \geq P$ . Then by Lemma 7.5 there holds  $p_0 \geq \min\{P, Q + 1\}$ ,  $p_1 \geq \min\{P, Q\}$ . Since  $Q \geq P$ , each of the assumptions  $p_0 = 0$  or  $p_1 = 0$  immediately gives  $P \leq 0$ .

Now assume  $p_0 \neq 0$ ,  $p_1 \neq 0$ . This yields  $P \neq 0$ . Let us show that  $Q \geq P$ . By (7.11) there holds  $p_0 \geq P$ . Thus

$$Q = \min_{j: p_j < P} (p_j + j - 1) = \min_{j \geq 1: p_j < P} (p_j + j - 1) \geq \min_{j \geq 1} (p_j + j - 1) \geq P.$$

Now we prove that the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$ . Indeed, by (7.8) in a neighborhood of  $\phi = 0$  there holds

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq \hat{C} \sum_{j \in \mathbb{N}} \|E_j(\phi, \nu, \Pi_h)\| \leq C \sum_{j \in \mathbb{N}} \min\{|\phi|^{p_j}, \nu |\phi|^{p_j+j}\}. \quad (7.13)$$

Therefore,

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq C \sum_{j: p_j < P} \nu |\phi|^{p_j+j} + C \sum_{j: p_j \geq P} |\phi|^{p_j} \leq C' \nu |\phi|^{Q+1} + C' |\phi|^P. \quad (7.14)$$

By Theorem 6.17 we get the desired estimate (2.17).

Now we show that  $P$  and  $Q$  are optimal. Let estimate (2.17) be valid under the substitutions  $Q'$  for  $Q$  and  $P'$  for  $P$  where  $Q' \geq P' > 0$ . Then for each error component given by (7.7) there holds (7.10) with the same substitution. By Lemma 7.5 for each  $j$  there holds one of the inequalities

$$\begin{cases} p_j \geq P', \\ p_j + j - 1 \geq Q'. \end{cases} \quad (7.15)$$

Since  $Q' \geq P'$ , (7.15) for each  $j$  yields  $P' \leq \max\{p_j, p_j + j - 1\}$ . Taking minimum over  $j$ , using (7.11) we get  $P' \leq P$ .

Now assume that  $P' = P$ . Then (7.15) yields

$$Q' \leq \min_{p_j < P'} (p_j + j - 1) = \min_{p_j < P} (p_j + j - 1) = Q.$$

Thus estimate (2.17) under the substitutions  $Q'$  for  $Q$  and  $P'$  for  $P$  yields the alternative

$$\begin{cases} P' < P, \\ P' = P, & Q' \leq Q. \end{cases}$$

This means that  $P$  and  $Q$  are optimal in the sense of Definition 5. □

#### 7.4. The main result

**Theorem 7.7.** *Let  $\Pi_h$  and  $\mathcal{P}_h$  be local mappings with the kernels  $\mu_\xi, \hat{\mu}_\xi$ , correspondingly, where  $\langle \hat{\mu}_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ . Suppose the scheme possesses the formal order of accuracy  $P \in \mathbb{N}$  and long-time simulation order  $Q \in \mathbb{N}$  in the sense of  $\Pi_h$ . Then there exist diagonal matrices  $\mathfrak{C}^{(m)} \in \mathbb{R}^{M^0}$ ,  $m = P, \dots, Q$ , such that the scheme possesses the truncation error of order  $Q$  in the sense of  $\Pi_h^{(P,Q)}$  given by (2.12).*



*Proof.* By Lemma 6.11 there exist  $C_1, C_2 \geq 0$  such that in a neighborhood of  $\phi = 0$  there holds

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| = \left\| \left( e^{\nu A(\phi)} - I \right) v(\phi, \Pi_h) \right\| \leq C_1 |\phi|^P + C_2 \nu |\phi|^{Q+1}.$$

Using representation (7.1), (7.2) provided by Theorem 7.1, in a neighborhood of  $\phi = 0$  define the holomorphic function

$$w(\phi) = S(\phi) \begin{pmatrix} \delta_0 I & 0 & \dots & 0 & 0 \\ 0 & \delta_1 I & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \delta_m I & 0 \\ 0 & 0 & \dots & 0 & \delta_\infty I \end{pmatrix} S^{-1}(\phi) v(\phi, \Pi_h),$$

where the blocks correspond to the blocks of  $M(\phi)$ ,  $\delta_j = 0$  if  $p_j \geq P$  and  $\delta_j = 1$  if  $p_j < P$ . Since  $S(\phi)$ ,  $S^{-1}(\phi)$ , and  $v(\phi, \Pi_h)$  are real-valued for  $i\phi \in \mathbb{R}$ , so is  $w(\phi)$ . Then  $\overline{w(-\phi)} = w(\bar{\phi})$ . By construction  $w(\phi)$  satisfies

$$\|w(\phi) - v(\phi, \Pi_h)\| = O(|\phi|^P).$$

Since

$$A(\phi)w(\phi) = S(\phi) \begin{pmatrix} \delta_0 M_0(\phi) & 0 & \dots & 0 & 0 \\ 0 & \delta_1 \phi M_1(\phi) & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \delta_m \phi^m M_m(\phi) & 0 \\ 0 & 0 & \dots & 0 & 0 \end{pmatrix} S^{-1}(\phi) v(\phi, \Pi_h),$$

by Lemma 7.5 in a neighborhood of  $\phi = 0$  there holds

$$\|A(\phi)w(\phi)\| \leq \tilde{C} \max_{j: p_j < P} |\phi|^j \|M_j(\phi)\| \|v_j(\phi, \Pi_h)\| = O(|\phi|^{Q+1}). \quad (7.16)$$

By Lemma 6.7 there exist real-valued diagonal matrices  $\mathfrak{C}^{(m)}$ ,  $m = P, \dots, Q$ , such that  $v(\phi, \Pi_h^{(P,Q)}) = w(\phi) + O(|\phi|^{Q+1})$ . Now (7.16) yields

$$\hat{\varepsilon}(\phi, \Pi_h^{(P,Q)}) = A(\phi)v(\phi, \Pi_h^{(P,Q)}) = O(|\phi|^{Q+1}).$$

By Theorem 6.15 the scheme possesses the truncation error of order  $Q$  in the sense of  $\Pi_h^{(P,Q)}$ . □

### 7.5. Construction of a scheme with specified properties

**Proposition 7.8.** Consider the transport equation  $\partial v / \partial t + \omega \partial v / \partial x = 0$ ,  $\omega \neq 0$ . Let  $\bar{\mathbb{N}}$  be a finite subset of  $\mathbb{N} \cup \{0\}$ . Then for any set of  $p_j \in \mathbb{N} \cup \{0, \infty\}$ ,  $j \in \bar{\mathbb{N}}$ , containing at least one zero there exists a scheme of the form (2.8) and a local mapping such that the values  $p_j$  coincide with those defined by (7.3).

510

*Proof.* Let  $c_{m,\eta}$ ,  $m \in \mathbb{N}$ , be coefficients of a finite-difference approximation of the first derivative on the uniform mesh with unit step possessing exactly the order  $m - 1$ , i. e.

$$i\phi - \sum_{\eta \in \mathcal{S}_m} c_{m,\eta} e^{i\phi\eta} = i\phi - \sum_{r=0}^{\infty} \frac{(i\phi)^r}{r!} \sum_{\eta \in \mathcal{S}_m} c_{m,\eta} \eta^r = \gamma_m \phi^m + O(\phi^{m+1}) \quad (7.17)$$

where  $\gamma_m \neq 0$ , and  $\mathcal{S}_m \subset \mathbb{Z}$  is a finite set. For  $m = 0$  put  $\mathcal{S}_0 = \{0\}$ ,  $c_{0,0} = \text{sign } \omega$ , so (7.17) also holds. For  $\eta \notin \mathcal{S}_m$  put  $c_{m,\eta} = 0$ .

Now we define a scheme of the form (2.8) and a local mapping. Put  $M^0 = \bar{\mathbb{N}}$ ,

$$\mathcal{S} = \bigcup_{m \in \bar{\mathbb{N}}} \mathcal{S}_m; \quad Z_0 = I; \quad Z_\eta = 0, \eta \neq 0; \quad L_\eta = \text{diag} \{ \omega c_{\xi,\eta}, \xi \in M^0 \};$$

$$(\Pi_h f)_{\eta, \xi} = h^{p_\xi} \frac{d^{p_\xi} f}{dx^{p_\xi}}(\eta h), \quad p_\xi < \infty; \quad (\Pi_h f)_{\eta, \xi} = 0, \quad p_\xi = \infty.$$

Then  $(v(\phi, \Pi_h))_\xi = (i\phi)^{p_\xi}$  and

$$A(\phi) = i\omega\phi - \sum_{\eta \in \mathcal{S}} L_\eta e^{i\phi\eta} = \text{diag} \left\{ i\omega\phi - \omega \sum_{\eta \in \mathcal{S}} c_{\xi, \eta} e^{i\phi\eta}, \xi \in M^0 \right\} = \text{diag} \left\{ \omega\gamma_\xi \phi^\xi + O(\phi^{\xi+1}), \xi \in M^0 \right\}.$$

Thus we have a representation of the form (7.1)–(7.2), with the blocks  $j \in \bar{\aleph}$  of the size 1 and no other blocks, and the values  $p_j$  coincide with those defined by (7.3). If the finite-difference schemes with the coefficients  $c_{m, \eta}$  are stable then the scheme we constructed is also stable.  $\square$

### 7.6. One-dimensional case: summary

Consider a stable scheme of the form (2.8) and a local mapping  $\Pi_h$ . In the sense of  $\Pi_h$ , the truncation error and the solution error have spectral representations (7.6) and (7.13), correspondingly. Then by Lemma 6.10 for  $v_0 = \exp(i\alpha x)$ ,  $\alpha/\pi \in \mathbb{Q}$ , there holds

$$\begin{aligned} \|\epsilon_h(v_0, \Pi_h)\| &\leq C \sum_{j \in \aleph} \|\nabla^{p_j+j} v_0\| h^{p_j+j-1} \leq C' \|\nabla^{P_A+1} v_0\| h^{P_A}, \\ \|\epsilon_h(t, v_0, \Pi_h)\| &\leq C \sum_{j \in \aleph} \min \{ \|\nabla^{p_j} v_0\| h^{p_j}, \|\nabla^{p_j+j} v_0\| h^{p_j+j-1} t \} \leq C' \|\nabla^P v_0\| h^P + C' \|\nabla^{Q+1} v_0\| h^Q t, \end{aligned} \quad (7.18)$$

where  $\aleph$  and  $p_j$  are defined in Section 7.1 and  $P_A$ ,  $P$ , and  $Q$  are given by (7.5), (7.11), and (7.12):

$$P_A = \min_{j \in \aleph} \{p_j + j - 1\}, \quad P = \min \{p_0, \min_{j \geq 1} (p_j + j - 1)\}, \quad Q = \min_{j \in \aleph: p_j < P} (p_j + j - 1).$$

Theorem 7.7 states that there exists  $\Pi_h^{(P, Q)}$  in the sense of which the truncation error satisfies

$$\|\epsilon_h(v_0, \Pi_h^{(P, Q)})\| \leq \tilde{C} \|\nabla^{Q+1} v_0\| h^Q.$$

By the Lax – Ryabenkii theorem (Theorem 4.1) the solution error satisfies

$$\|\epsilon_h(t, v_0, \Pi_h^{(P, Q)})\| \leq \tilde{C} K \|Z^{-1}\| \|\nabla^{Q+1} v_0\| h^Q t.$$

Theorems 6.15 and 6.17 allow to extend these estimates to  $v_0$  smooth enough.

*Example.* Put  $\bar{\aleph} = \{3, 5\}$ ,  $p_3 = 1$ , and  $p_5 = 0$ . A scheme with these parameters exists by Proposition 7.8. By Lemma 7.4 this scheme possesses the truncation error of order  $P_A = 3$ . For  $v_0 = \exp(i\alpha x)$ ,  $\alpha/\pi \in \mathbb{Q}$ , the error estimate (7.18) reads as

$$\|\epsilon_h(t, v_0, \Pi_h)\| \leq C_1 \min \{ h \|\nabla v_0\|, h^3 t \|\nabla^4 v_0\| \} + C_2 \min \{ \|v_0\|, h^4 t \|\nabla^5 v_0\| \}.$$

By Lemma 7.6 the scheme possesses the formal order of accuracy  $P = 3$  and the long-time simulation order  $Q = 3$ , these values being optimal in the sense of Definition 5. However we can write the estimate

$$\|\epsilon_h(t, v_0, \Pi_h)\| \leq C_1 h \|\nabla v_0\| + C_2 h^4 t \|\nabla^5 v_0\|,$$

so the scheme possesses the 1st formal order of accuracy and the 4th order in the long-time simulation. The values  $P = 1$  and  $Q = 4$  are not optimal in the sense of our definition. Note that due to Theorem 7.7 there exists a local mapping that differs from  $\Pi_h$  by  $O(h)$  and gives the 4th order of the truncation error, the 4th formal order of accuracy and the 4th long-time simulation order.

## 8. The general case

In this section we prove Theorem 2 and state some of its corollaries.

525 8.1. The method of auxilliary mapping

The following lemma describes the method of auxiliary mapping (see Proposition 4.2) using the spectral representation of a scheme.

**Lemma 8.1.** Consider a scheme of the form (2.8) stable with a constant  $K$  and a bounded homogeneous mapping  $\Pi_h$  of  $H_{per}^q(\mathbb{R}^d)$  to  $V_{per}$  or a local mapping  $\Pi_h$  with  $\mu_\xi \in (C^{[q]}(G))^*$ . Suppose there exists a function  $w(\phi)$  and  $c_1, c_2 > 0$  such that in a neighborhood of  $\phi = 0$  there holds

$$\|v(\phi, \Pi_h) - w(\phi)\| \leq c_1 |\phi|^P, \quad \|A(\phi)w(\phi)\| \leq c_2 |\phi|^{Q+1}, \quad (8.1)$$

where  $0 < P \leq Q \leq \infty$  (we assume  $|\phi|^\infty \equiv 0$ ). Let  $r \geq \max\{Q+1, q\}$  (if  $P \leq Q < \infty$ ) or  $r \geq \max\{P, q\}$ , (if  $P < Q = \infty$ ) or  $r \geq q$  (if  $P = Q = \infty$ ). Then for  $\hat{\varepsilon}(\phi, \nu, \Pi_h)$  defined by (6.16) there holds

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq (K+1)c_1 |\phi|^P + \nu K c_2 |\phi|^{Q+1} \quad (8.2)$$

and the scheme possesses the error estimate (2.17) on  $H_{per}^r(\mathbb{R}^d)$  (or  $C_{per}^{[r]}(\mathbb{R}^d)$ ) with the same  $P$  and  $Q$ ,  $C_1 = (K+1)c_1$ ,  $C_2 = Kc_2$ , and  $C_3 > 0$  depending only on  $\Pi_h, r, T, K$ , the neighborhood of  $\phi = 0$  in the conditions of the lemma, and the norm on  $\mathbb{C}^{M^0}$ .

*Proof.* We have

$$\begin{aligned} \|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| &= \left\| \left( e^{\nu A(\phi)} - I \right) v(\phi, \Pi_h) \right\| \leq \left\| \left( e^{\nu A(\phi)} - I \right) (w(\phi) - v(\phi, \Pi_h)) \right\| + \left\| \left( e^{\nu A(\phi)} - I \right) w(\phi) \right\| \leq \\ &\leq \|\exp(\nu A(\phi)) - I\| \|w(\phi) - v(\phi, \Pi_h)\| + \left\| \int_0^\nu e^{\tau A(\phi)} d\tau A(\phi) w(\phi) \right\| \leq \\ &\leq (\|\exp(\nu A(\phi))\| + 1) \|w(\phi) - v(\phi, \Pi_h)\| + \int_0^\nu \|e^{\tau A(\phi)}\| d\tau \|A(\phi)w(\phi)\| \leq (K+1)c_1 |\phi|^P + \nu K c_2 |\phi|^{Q+1}. \end{aligned}$$

Thus we get (8.2). It remains to use Theorem 6.17.  $\square$

8.2. The existence of an auxiliary mapping

Let  $\mathbb{C}^{n \times n}$  be the space of complex matrices of size  $n$ . We need the following result.

**Theorem 8.2** ([28]). Let  $A \in \mathbb{C}^{n \times n}$  and  $v \in \mathbb{C}^n$ . Suppose for each  $\nu \geq 0$  there hold  $\|e^{\nu A}\| \leq K$  and

$$\|(e^{\nu A} - I)v\| \leq (\tilde{C}_1 + \tilde{C}_2 \nu) \|v\|. \quad (8.3)$$

Then for  $w = (A^* A + \varepsilon^2)^{-1} \varepsilon^2 v$  where  $\varepsilon = \tilde{C}_2 / \tilde{C}_1$  (if  $\tilde{C}_1 = 0$  or  $\tilde{C}_2 = 0$  then the vector  $w$  is the corresponding limit) there hold

$$\|v - w\| \leq \delta \tilde{C}_1 \|v\|, \quad \|Aw\| \leq \delta \tilde{C}_2 \|v\|, \quad (8.4)$$

where  $\delta$  depends on  $n$  and  $K$  only.

Note that  $w = \operatorname{argmin}(\tilde{C}_2^2 \|v - w\|^2 + \tilde{C}_1^2 \|Aw\|^2)$ .

Denote

$$\mathcal{W}^{P,Q}(\phi) = \left( A^*(\phi)A(\phi) + |\phi|^{2(Q+1-P)} \right)^{-1} |\phi|^{2(Q+1-P)} v(\phi, \Pi_h), \quad (8.5)$$

where  $v(\phi, \Pi_h)$  and  $A(\phi)$  are defined by (6.7) and (6.14).

**Corollary 8.3.** Let the scheme (2.8) be stable,  $\Pi_h$  be a local mapping, and (2.17) hold. In a neighborhood of  $\phi = 0$  there holds

$$\|v(\phi, \Pi_h) - \mathcal{W}^{P,Q}(\phi)\| \leq \delta C_1 |\phi|^P, \quad \|A(\phi)\mathcal{W}^{P,Q}(\phi)\| \leq \delta C_2 |\phi|^{Q+1}, \quad (8.6)$$

where  $\delta$  depends only on  $M^0$ ,  $\|\Pi_h 1\|$ , and the stability constant.

*Proof.* We have

$$\|v(\phi, \Pi_h)\| = \|(\Pi_1 \exp(i\phi \cdot \mathbf{r}))_0\| \rightarrow \|\Pi_h 1\| \quad \text{as } \phi \rightarrow 0.$$

From (2.17) by Lemma 6.11 in a neighborhood of  $\phi = 0$  there holds

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq \frac{3}{2} C_1 |\phi|^P + \nu C_2 |\phi|^{Q+1}.$$

By (6.16) in a neighborhood of  $\phi = 0$  there holds

$$\|(e^{\nu A(\phi)} - I)v(\phi, \Pi_h)\| \leq (2C_1 |\phi|^P \|\Pi_h 1\|^{-1} + 2C_2 |\phi|^{Q+1} \|\Pi_h 1\|^{-1} \nu) \|v(\phi, \Pi_h)\|.$$

It remains to use Theorem 8.2. □

Now we are ready to prove Theorem 2, which we repeat here.

**Theorem 2.** *Let  $\Pi_h$  be a local mapping with  $\mu_\xi \in (W_2^q(G))^*$  or  $(C^q(G))^*$ . Let  $P, Q > 0$ ,  $r \geq \max\{P, q\}$ , and  $R \geq \max\{Q + 1, r\}$ . Let the scheme (2.8) be stable and possess the error estimate (2.17) on  $H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^{[R]}(\mathbb{R}^d)$ ). Then there exists a homogeneous mapping  $\tilde{\Pi}_h : L_{2,per}(\mathbb{R}^d) \rightarrow V_{per}$  such that*

$$\|\tilde{\Pi}_h f - \Pi_h f\| \leq C(h^P \|\nabla^P f\| + h^r \|\nabla^r f\|_*), \quad \|\epsilon_h(f, \tilde{\Pi}_h)\| \leq Ch^Q \|\nabla^{Q+1} f\| \quad (8.7)$$

540 *for each  $h$  and  $f \in H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^{[R]}(\mathbb{R}^d)$ ), where  $C$  does not depend on  $h$  and  $f$ , and  $\|\cdot\|_*$  means either  $\|\cdot\|$  or  $\|\cdot\|_\infty$  depending on the case.*

*Proof.* First we consider the case  $\mu_\xi \in (W_2^q(G))^*$ . Recall the notation  $\mathcal{A}_\beta = \{\alpha \in \mathcal{A} : T^* \alpha \in [-\beta, \beta]^d\}$  for  $\beta > 0$ . By Corollary 8.3 there exists  $\beta > 0$  such that (8.6) holds for  $\phi \in \mathcal{A}_\beta$ . Since  $v(-\phi, \Pi_h) = v(\phi, \tilde{\Pi}_h)$  and  $A(-\phi) = \overline{A(\phi)}$ , by construction we have  $\mathcal{W}^{P,Q}(-\phi) = \overline{\mathcal{W}^{P,Q}(\phi)}$ . Using the continuity of  $v(\phi, \Pi_h)$  and the first inequality in (8.6) we get

$$\sup_{\phi \in \mathcal{A}_\beta} \|\mathcal{W}^{P,Q}(\phi)\| \leq \sup_{\phi \in \mathcal{A}_\beta} \|v(\phi, \Pi_h)\|$$

By Lemma 6.6 with  $G = \mathcal{A}_\beta$  there exists a bounded homogeneous mapping  $\tilde{\Pi}_h$  such that  $\mathcal{W}^{P,Q}(\phi) = v(\phi, \tilde{\Pi}_h)$ . Let  $F_h = \tilde{\Pi}_h - \Pi_h$ . For  $\alpha \in \mathcal{A}_{\beta/h}$  using the first inequality in (8.6) we obtain

$$\|F_h e^{i\alpha \cdot \mathbf{r}}\| = \|(F_h e^{i\alpha \cdot \mathbf{r}})_0\| = \|(\tilde{\Pi}_h e^{i\alpha \cdot \mathbf{r}})_0 - (\Pi_h e^{i\alpha \cdot \mathbf{r}})_0\| = \|\mathcal{W}^{P,Q}(\alpha h) - v(\alpha h, \Pi_h)\| \leq \delta C_1 h^P |\alpha|^P. \quad (8.8)$$

Applying Lemma 6.14 we obtain the first inequality in (8.7). The second inequality in (8.6) yields the second inequality in (8.7) by Theorem 6.15.

Now consider the case  $\mu_\xi \in (C^q(G))^*$ . By Lemma 5.9 there exists a local mapping  $\hat{\Pi}_h$  with the kernel  $\hat{\mu}_\xi \in L_2(G + B_1(0))$  such that

$$\|\Pi_h f - \hat{\Pi}_h f\| \leq Ch^r \|\nabla^r f\|_\infty.$$

By the above proof we have (8.7) with the substitution of  $\hat{\Pi}_h$  for  $\Pi_h$ , and it remains to use the triangle inequality. □

### 545 8.3. An analytical criterion

For a given scheme (2.8), a mapping  $\Pi_h$ , and two numbers  $0 < P \leq Q < \infty$  denote

$$\mathcal{F}^{P,Q}(\phi) = |\phi|^{-2P} v^*(\phi, \Pi_h) \left( A^*(\phi) A(\phi) + |\phi|^{2(Q+1-P)} \right)^{-1} A^*(\phi) A(\phi) v(\phi, \Pi_h). \quad (8.9)$$

**Lemma 8.4.** *For  $\mathcal{W}^{P,Q}$  given by (8.5) there holds*

$$\mathcal{F}^{P,Q}(\phi) = |\phi|^{-2P} \|v(\phi, \Pi_h) - \mathcal{W}^{P,Q}(\phi)\|^2 + |\phi|^{-2(Q+1)} \|A(\phi) \mathcal{W}^{P,Q}(\phi)\|^2. \quad (8.10)$$

*Proof.* We have

$$|\phi|^{-2(Q+1-P)} A^*(\phi) A(\phi) \mathcal{W}^{P,Q}(\phi) = v(\phi, \Pi_h) - \mathcal{W}^{P,Q}(\phi).$$

Thus

$$\begin{aligned} & |\phi|^{-2(Q+1-P)} \|A(\phi) \mathcal{W}^{P,Q}(\phi)\|^2 = \\ & = |\phi|^{-2(Q+1-P)} (A^*(\phi) A(\phi) \mathcal{W}^{P,Q}(\phi), \mathcal{W}^{P,Q}(\phi)) = (v(\phi, \Pi_h) - \mathcal{W}^{P,Q}(\phi), \mathcal{W}^{P,Q}(\phi)). \end{aligned}$$

On the other hand,

$$|\phi|^{2P} \mathcal{F}^{P,Q}(\phi) = (|\phi|^{-2(Q+1-P)} A^*(\phi) A(\phi) \mathcal{W}^{P,Q}(\phi), v(\phi, \Pi_h)) = (v(\phi, \Pi_h) - \mathcal{W}^{P,Q}(\phi), v(\phi, \Pi_h)).$$

Now (8.10) is obvious.  $\square$

The following theorem establishes an analytical criterion which allows us (at least in theory) to find the optimal values of the formal order of accuracy and the long-time simulation order.

**Theorem 8.5.** *Let the scheme (2.8) be stable,  $\Pi_h$  be a local mapping with kernel  $\mu_\xi \in (W_2^q(G))^*$  (or  $(C^q(G))^*$ ),  $0 < P \leq Q < \infty$ ,  $R \geq \max\{Q, q+1\}$ . The scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^{[R]}(\mathbb{R}^d)$ ) if and only if  $\mathcal{F}^{P,Q}(\phi)$  is bounded in a neighborhood of  $\phi = 0$ .*

*Proof.* Assume the error estimate (2.17). Then in an neighborhood of  $\phi = 0$  we have (8.6) and thus by (8.10) there holds

$$\mathcal{F}^{P,Q}(\phi) \leq |\phi|^{-2P} (\delta^2 C_1^2 |\phi|^{2P}) + |\phi|^{-2(Q+1)} (\delta^2 C_2^2 |\phi|^{2(Q+1)}) = \delta^2 (C_1^2 + C_2^2).$$

On the other hand, assuming that  $\mathcal{F}(\phi) \leq C^2$  from (8.10) we get

$$\|v(\phi, \Pi_h) - \mathcal{W}^{P,Q}(\phi)\| \leq C|\phi|^P, \quad \|A(\phi) \mathcal{W}^{P,Q}(\phi)\| \leq C|\phi|^{Q+1}.$$

By Lemma 8.1 the scheme possesses error estimate (2.17) on  $H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^{[R]}(\mathbb{R}^d)$ ).  $\square$

**Lemma 8.6.** *Let the scheme (2.8) be stable,  $\Pi_h$  be a local mapping with kernel  $\mu_\xi \in (W_2^q(G))^*$  (or  $(C^q(G))^*$ ),  $0 < P \leq Q < \infty$ ,  $R \geq \max\{Q, q+1\}$ . Let  $\tilde{e} \in \mathring{\Omega}$ . The scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per,\tilde{e}}^R(\mathbb{R}^d)$  (or  $C_{per,\tilde{e}}^R(\mathbb{R}^d)$ ) if and only if  $\mathcal{F}^{P,Q}(\tilde{e}\psi)$  is bounded as  $\psi$  tends to zero.*

The proof is similar to Theorem 8.5.

**Proposition 8.7.** *Let  $\Pi_h$  and  $\mathcal{P}_h$  be local mappings with kernels  $\mu_\xi, \hat{\mu}_\xi \in (W_2^q(G))^*$  (or  $(C^q(G))^*$ ), and  $\langle \hat{\mu}_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ . Let  $e \in \mathring{\Omega}$ . Suppose a stable scheme of the form (2.8) possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q < \infty$  in the sense of  $\Pi_h$  on  $H_{per,e}^R(\mathbb{R}^d)$  (or  $C_{per,e}^{[R]}(\mathbb{R}^d)$ ), where  $R = \max\{Q+1, q\}$ . Then there exist real-valued diagonal matrices  $\mathfrak{C}_e^{(m)}$ ,  $m = \lceil P \rceil, \dots, \lceil Q \rceil$ , such that the scheme possesses the truncation error of order  $\lceil Q \rceil$  in the sense of  $\Pi_{h,e}^{(\lceil P \rceil, \lceil Q \rceil)}$  given by (2.12) on the direction  $e$ . As a corollary, the optimal values of the order of accuracy and of the long-time simulation order on  $H_{per,e}^R(\mathbb{R}^d)$  (or  $C_{per,e}^{[R]}(\mathbb{R}^d)$ ) are integers.*

*Proof.* Let  $\mathcal{W}^{P,Q}(\phi)$  be given by (8.5) and  $W(\psi) = \mathcal{W}^{P,Q}(e\psi)$ . By Lemma 8.6 the function  $\mathcal{F}^{P,Q}(e\psi)$ , where  $\mathcal{F}^{P,Q}$  is defined by (8.9), is bounded as  $\psi \rightarrow 0$ . Using (8.10) for  $\phi = e\psi$  we have

$$\|v(\psi e, \Pi_h) - W(\psi)\| = O(|\psi|^P), \quad \|A(\psi e)W(\psi)\| = O(|\psi|^{Q+1}). \quad (8.11)$$

By construction each component of  $W(\psi)$  is a ratio of two analytical functions of the real argument  $\psi$ ; by (8.11) it is bounded as  $\psi \rightarrow 0$ , thus  $W(\psi)$  is an analytical function. Thus both expressions under norm signs in (8.11) are analytical at  $\psi = 0$  and so  $P$  and  $Q$  can be replaced by  $\lceil P \rceil$  and  $\lceil Q \rceil$ , correspondingly. By Lemma 6.8 there exist

real-valued diagonal matrices  $\mathfrak{C}_e^{(m)}$ ,  $m = \lceil P \rceil, \dots, \lceil Q \rceil$ , such that for the mapping  $\Pi_{h,e}^{(\lceil P \rceil, \lceil Q \rceil)}$  given by (2.13) there holds

$$v(\psi e, \Pi_{h,e}^{(\lceil P \rceil, \lceil Q \rceil)}) = W(\psi) + O(|\psi|^{\lceil Q \rceil+1}),$$

thus by the triangle inequality we have

$$\|\hat{e}(\psi e, \Pi_{h,e}^{(\lceil P \rceil, \lceil Q \rceil)})\| = \|A(\psi e)v(\psi e, \Pi_{h,e}^{(\lceil P \rceil, \lceil Q \rceil)})\| = O(|\psi|^{\lceil Q \rceil+1}).$$

It remains to use Theorem 6.15. □

Given a scheme with enhanced accuracy in the long-time simulation, Theorem 2 states that there exists an auxiliary mapping  $\tilde{\Pi}_h$  explaining this fact. This mapping is generally not of the form (2.11). For the 1D case, Proposition 8.7 states the existence of an auxiliary mapping of the form (2.11). However, in the multidimensional case, an auxiliary mapping of the form (2.11) generally does not exist, see a counter-example in Section 12.3.

## 9. The good

Throughout this section we consider a stable scheme of the form (2.8) with stability constant  $K$  and a local mapping  $\Pi_h$ . In this section we will prove Theorem 1 and Theorem 3, which is splitted into Theorems 3A and 3B.

### 9.1. The order of accuracy

In this subsection we establish the optimal value of the formal order of accuracy.

**Lemma 9.1.** *Let the matrix  $A(\phi)$  be given by (6.14). Then the eigenvalue  $\lambda = 0$  of  $A(0)$  is semisimple, i. e. there are no Jordan cells of size greater than one corresponding to  $\lambda = 0$ .*

*Proof.* Assume the converse. Represent the matrix  $A(0)$  in the form  $A(0) = SJS^{-1}$  where  $J$  is its Jordan normal form. Then  $e^{\nu A(0)} = Se^{\nu J}S^{-1}$ . The explicit expression for  $e^{\nu J}$  shows that  $\|e^{\nu J}\|$  grows unlimitedly as  $\nu \rightarrow \infty$  and so does  $\|e^{\nu A(0)}\|$ . This contradicts the stability condition (6.21). □

Consider  $\hat{e}(\phi, \nu, \Pi_h)$  given by (6.16). By Theorem 5.19 in a neighborhood of  $\phi = 0$  there holds  $A(\phi) = S(\phi)M(\phi)S^{-1}(\phi)$  where  $M(\phi)$  has the form (5.10). Let  $v_j(\phi, \Pi_h)$  be the components of  $S^{-1}(\phi)v(\phi, \Pi_h)$  corresponding to the blocks  $M^{(j)}(\phi)$ . Denote

$$E_j(\phi, \nu, \Pi_h) = (e^{\nu M^{(j)}(\phi)} - I)v_j(\phi, \Pi_h), \quad (9.1)$$

then

$$\hat{e}(\phi, \nu, \Pi_h) = S(\phi)(e^{M(\phi)\nu} - I)S^{-1}(\phi)v(\phi, \Pi_h) = S(\phi) \begin{pmatrix} \vdots \\ E_j(\phi, \nu, \Pi_h) \\ \vdots \end{pmatrix}. \quad (9.2)$$

Note that the notation is similar to the 1D case but the decomposition is different. We will also use the notation  $M^{(*)}(\phi)$  for the submatrix of  $M(\phi)$  containing the blocks  $M^{(j)}(\phi)$ ,  $j \neq 0$ , and the notation  $v_*(\phi, \Pi_h)$ ,  $E_*(\phi, \nu, \Pi_h)$  for the unions of all blocks excluding  $j = 0$  of the corresponding vectors. Below we assume that the norms of subvectors are inherited from  $\mathbb{C}^{M^0}$ , i. e., for example,  $\|v_j\| = \|(0, \dots, 0, v_j, 0, \dots, 0)^T\|$ .

Since  $M^{(0)}(0) = 0$ , there holds  $M^{(0)}(\phi) = \sum_{k=1}^d M_k^{(0)}(\phi)\phi_k$  where  $M_k^{(0)}(\phi)$  are holomorphic. Thus

$$M^{(0)}(\alpha h) \frac{t}{h} = \sum_{k=1}^d M_k^{(0)}(\alpha h) \alpha_k t$$

is a holomorphic function of  $t$ ,  $h$ , and  $\alpha$ . Particularly, if  $A(0) = 0$  then  $\hat{e}(\alpha h, t/h, \Pi_h)$  is a holomorphic function of  $t$ ,  $h$ , and  $\alpha$  for each  $t$  in a neighborhood of  $h = 0$  and  $\alpha = 0$ .

**Lemma 9.2.** Suppose the scheme (2.8) possesses the estimate (2.17) in the sense of a local mapping  $\Pi_h$ . Then in a neighborhood of  $\phi = 0$  for each  $\nu \geq 0$  there holds

$$\|v_*(\phi, \Pi_h)\| \leq \tilde{\delta} C_1 |\phi|^{\lceil P \rceil}, \quad \|E_*(\phi, \nu, \Pi_h)\| \leq (K+1) \tilde{\delta} C_1 |\phi|^{\lceil P \rceil}. \quad (9.3)$$

Moreover, if  $P$  is integer, then  $\tilde{\delta}$  depends only on  $|M^0|$ ,  $K$ , and the norm on  $\mathbb{C}^{M^0}$ .

*Proof.* By Lemma 6.11 in a neighborhood of  $\phi = 0$  there holds

$$\|\hat{e}(\phi, \nu, \Pi_h)\| \leq 2C_1 |\phi|^P + C_2 \nu |\phi|^{Q+1}.$$

Thus for each  $j$  we have

$$\|E_j(\phi, \nu, \Pi_h)\| \leq c_n \|S^{-1}(\phi)\| (2C_1 |\phi|^P + C_2 \nu |\phi|^{Q+1}),$$

where  $c_n$  depends on the norms in use. Directly from (9.1) for each  $j \neq 0$  and each  $\nu \geq 0$  we get

$$\|v_j(\phi, \Pi_h)\| \leq \left\| \left[ \exp(M^{(j)}(\phi)\nu) - I \right]^{-1} \right\| \|E_j(\phi, \nu, \Pi_h)\|.$$

Put  $\nu = 1/(2\|M^{(j)}(0)\|)$ . Then by Lemma 5.13 in a neighborhood of  $\phi = 0$  we have

$$\begin{aligned} \|v_j(\phi, \Pi_h)\| &\leq 4 \left\| \left( \frac{M^{(j)}(\phi)}{2\|M^{(j)}(0)\|} \right)^{-1} \right\| \|E_j(\phi, \nu, \Pi_h)\| \leq 10\kappa(M^{(j)}(0)) \left\| E_j \left( \phi, \frac{1}{2\|M^{(j)}(0)\|}, \Pi_h \right) \right\| \leq \\ &\leq 20c_n \delta^2 \left( 2C_1 |\phi|^P + C_2 \frac{1}{2\|M^{(j)}(0)\|} |\phi|^{Q+1} \right) \leq 100\delta^2 c_n C_1 |\phi|^P. \end{aligned}$$

585 Since  $v_j(\phi, \Pi_h)$  is holomorphic then the degree  $P$  in the right-hand side can be improved to  $\lceil P \rceil$ . This leads to the first inequality in (9.3). The second inequality follows from the first by stability.  $\square$

**Lemma 9.3.** Let the scheme (2.8) possess the estimate (2.17) in the sense of a local mapping  $\Pi_h$ . Then in a neighborhood of  $\phi = 0$  there holds

$$\|M^{(0)}(\phi)v_0(\phi, \Pi_h)\| \leq C|\phi|^{\lceil P \rceil+1}. \quad (9.4)$$

*Proof.* Similar to the previous lemma, we have

$$\|E_0(\phi, \nu, \Pi_h)\| \leq c_n \|S^{-1}(\phi)\| (2C_1 |\phi|^P + C_2 \nu |\phi|^{Q+1}). \quad (9.5)$$

Using the function  $f$  given by (5.7) ( $f(x) = (e^x - 1)/x$ ), formula (9.1) is equivalent to

$$E_0(\phi, \nu, \Pi_h) = f(M^{(0)}(\phi)\nu) M^{(0)}(\phi) \nu v_0(\phi, \Pi_h).$$

Assuming  $M^{(0)}(\phi) \neq 0$  put  $\nu = 1/\|M^{(0)}(\phi)\|$ . By Lemma 5.12 we get  $\|(f(\nu M^{(0)}(\phi)))^{-1}\| \leq 4$ , thus

$$\left\| M^{(0)}(\phi)v_0(\phi, \Pi_h) \right\| \leq \frac{4}{\nu} \|E_0(\phi, \nu, \Pi_h)\| \leq C'_1 |\phi|^P \|M^{(0)}(\phi)\| + C'_2 |\phi|^{Q+1}.$$

If  $M^{(0)}(\phi) = 0$  then  $M^{(0)}(\phi)v_0(\phi, \Pi_h) = 0$ . Since  $M^{(0)}(\phi) = O(|\phi|)$  as  $\phi \rightarrow 0$  and  $Q \geq P$  we get  $\|M^{(0)}(\phi)v_0(\phi, \Pi_h)\| = O(|\phi|^{P+1})$  as  $\phi \rightarrow 0$ . Taking into account that  $M^{(0)}(\phi)$  and  $v_0(\phi, \Pi_h)$  are holomorphic, we arrive at (9.4).  $\square$

590 Now we are ready to prove Theorem 1, which we repeat here. Recall that for  $n \in \mathbb{N}$  by  $\hat{\Omega}_n$  we denote a set of vectors  $\{e_k \in \hat{\Omega}, k = 1, \dots, C_{n+d-1}^{d-1}\}$  such that  $\{(e_k \cdot r)^n\}$  form a basis in the set of homogeneous polynomials of order  $n$  (see Lemma A.6).

**Theorem 1.** Consider a scheme of the form (2.8), stable with a constant  $K$ , and local mappings  $\Pi_h, \mathcal{P}_h$  with kernels  $\mu_\xi, \hat{\mu}_\xi \in (W_2^q(G))^*$  (or  $(C^q(G))^*$ ). Let  $P_A$  and  $P$  be the optimal orders of the truncation error and accuracy, correspondingly, in the sense of  $\Pi_h$ . Let  $\hat{\kappa} = \min_{\xi \in M^0} |\langle \hat{\mu}_\xi, 1 \rangle| > 0$  and  $R = \max\{q, P\} + 1$ . Then the following holds.

1.  $P_A$  and  $P$  are integers.
2. Either  $P = P_A$  or  $P = P_A + 1$ .
3. If  $P = P_A + 1$  then there exist real-valued diagonal matrices  $\mathfrak{C}^{(\mathbf{m})}$ ,  $|\mathbf{m}| = P$ , such that the scheme possesses the truncation error of order  $P$  in the sense of  $\Pi_h^{(P,P)}$  given by (2.12). Moreover,  $\|\mathfrak{C}^{(\mathbf{m})}\| \leq \tilde{\delta} C_1$ , where  $C_1$  is the constant in estimate (2.17) and  $\tilde{\delta}$  depends only on  $K$ ,  $\hat{\kappa}$ ,  $P$ ,  $|M^0|$ , and the norm on  $\mathbb{C}^{M^0}$ .
4. If  $P = P_A$  then there exists no set of matrices  $\{\mathfrak{C}^{(\mathbf{m})}, |\mathbf{m}| = P_A + 1\}$ , such that the scheme possesses the truncation error of order  $P_A + 1$  in the sense of  $\Pi_h^{(P_A+1, P_A+1)}$  given by (2.12).
5. If  $L(0) = 0$  then  $P = P_A$ .
6. If the scheme possesses the formal order of accuracy  $P_A + 1$  on  $H_{per, \mathbf{e}}^R(\mathbb{R}^d)$  (or  $C_{per, \mathbf{e}}^R(\mathbb{R}^d)$ ) for each  $\mathbf{e} \in \mathring{\Omega}_{P_A+1}$ , then it possesses the formal order of accuracy  $P = P_A + 1$  on  $H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^R(\mathbb{R}^d)$ ).
7.  $P$  coincides with the optimal order of accuracy in the weak sense.

*Proof.* 1. Let the scheme possess the order of the truncation error  $P_A$ . Then by Corollary 6.21 there holds (6.31). Taking  $v_0 = \exp(i\alpha \cdot \mathbf{r})$ , by Lemma 6.10 in a neighborhood of  $\phi = 0$  we get  $\|\hat{\epsilon}(\phi, \Pi_h)\| \leq c|\phi|^{P_A+1}$ . Since  $\hat{\epsilon}(\phi, \Pi_h)$  is a holomorphic function, then we have  $\|\hat{\epsilon}(\phi, \Pi_h)\| \leq \tilde{c}|\phi|^{\lceil P_A \rceil + 1}$ . By Theorem 6.15 the scheme possesses the truncation error of order  $\lceil P_A \rceil$ . Thus the optimal value of the order of the truncation error is an integer.

Let the scheme possess the order of accuracy  $p$  (possibly not optimal). Put by definition

$$w(\phi) = S(\phi) \begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix} S^{-1}(\phi) v(\phi, \Pi_h), \quad (9.6)$$

where the identity matrix corresponds to the block  $M^{(0)}(\phi)$ . By Lemma 9.2 we have

$$\|w(\phi) - v(\phi, \Pi_h)\| = \left\| S(\phi) \begin{pmatrix} v_*(\phi, \Pi_h) \\ 0 \end{pmatrix} \right\| \leq c_n \|S(\phi)\| \|v_*(\phi, \Pi_h)\| \leq \delta C_1 |\phi|^{\lceil p \rceil}. \quad (9.7)$$

If  $p$  is integer,  $\delta$  in (9.7) depends only on  $K$ ,  $|M^0|$  and the choice of norms. By Lemma 9.3 there holds

$$A(\phi)w(\phi) = S(\phi) \begin{pmatrix} 0 \\ M^{(0)}(\phi)v_0(\phi, \Pi_h) \end{pmatrix} = O(|\phi|^{\lceil p \rceil + 1}). \quad (9.8)$$

By Lemma 8.1 the scheme possesses the formal order of accuracy  $\lceil p \rceil$ . This proves that the optimal order of accuracy is an integer.

2. By Lemma 9.2 and Lemma 9.3 there holds

$$\|\hat{\epsilon}(\phi, \Pi_h)\| = \left\| S(\phi) \begin{pmatrix} M^{(*)}(\phi)v_*(\phi, \Pi_h) \\ M^{(0)}(\phi)v_0(\phi, \Pi_h) \end{pmatrix} \right\| = O(|\phi|^P).$$

By Theorem 6.15 the scheme possesses the order of the truncation error  $P - 1$ . Thus  $P \leq P_A + 1$ . On the other hand, by Theorem 4.1 we have  $P \geq P_A$ . Since both  $P$  and  $P_A$  are integers there holds either  $P = P_A$  or  $P = P_A + 1$ .

3. Now assume  $P = P_A + 1$ . Let  $w$  be given by (9.6). Since  $v(\phi, \Pi_h) = \overline{v(-\phi, \Pi_h)}$  and  $A(\phi) = \overline{A(-\phi)}$ , for  $W(\phi) = (w(\phi) + \overline{w(-\phi)})/2$  by (9.7) and (9.8) with  $p = P$  we have

$$\|W(\phi) - v(\phi, \Pi_h)\| \leq \delta c_n C_1 |\phi|^P, \quad A(\phi)W(\phi) = O(|\phi|^{P+1}). \quad (9.9)$$

By Lemma 6.7 there exist  $\mathfrak{C}^{(\mathbf{m})}$ ,  $|\mathbf{m}| = P$ , such that for  $\Pi_h^{(P,P)}$  there holds

$$v(\phi, \Pi_h^{(P,P)}) = W(\phi) + O(|\phi|^{P+1}).$$



By the triangle inequality we get

$$\hat{\epsilon}(\phi, \Pi_h^{(P,P)}) = A(\phi)v(\phi, \Pi_h^{(P,P)}) = O(|\phi|^{P+1}),$$

so the coefficients  $\mathfrak{C}^{(m)}$  satisfy the statement of the theorem. By construction (see (6.12) and (6.13)),

$$(\mathfrak{C}^{(m)})_{\xi, \xi} = \frac{1}{i^{|m|} m!} (v_\xi(0, \mathcal{P}_h))^{-1} \left( D^m (W(\phi) - v(\phi, \Pi_h))|_{\phi=0} \right)_\xi.$$

Using (9.9) and Lemma 5.2 we get  $|(\mathfrak{C}^{(m)})_{\xi, \xi}| \leq \hat{\kappa}^{-1} \delta c_n c_P C_1$ , where  $c_P$  depends on  $P$  only. From here, the estimate for  $\|\mathfrak{C}^{(m)}\|$  is obvious.

4. Suppose there exists a set of matrices  $\mathfrak{C}^{(m)}$  such that the scheme possesses the truncation error of order  $P_A + 1$  in the sense of  $\Pi_h^{(P_A+1, P_A+1)}$ . Then for  $w(\phi) = v(\phi, \Pi_h^{(P_A+1, P_A+1)})$  we have (8.1) with  $P = Q = P_A + 1$  and thus by Lemma 8.1 the scheme possesses the formal order of accuracy  $P_A + 1$  and the long-time simulation order  $P_A + 1$ .

5. If  $L(0) = 0$ , the block  $M^{(*)}(\phi)$  does not exist, and by Lemma 9.3 we have  $\|\hat{\epsilon}(\phi, \Pi_h)\| = O(|\phi|^{P+1})$ . By Theorem 6.15 the scheme possesses the order of the truncation error  $P$ , i. e. there holds  $P = P_A$ .

6. Now let the scheme possess the order of accuracy  $p = P_A + 1$  on  $H_{per, e}^R(\mathbb{R}^d)$  for each  $e \in \mathring{\Omega}_p$ . Arguing as in the proof of Lemma 6.11 for  $\phi \in \text{span}\{e\}$  in a neighborhood of  $\phi = 0$  we get

$$\|\hat{\epsilon}(\phi, \nu, \Pi_h)\| \leq c_1(e)|\phi|^p + c_2(e)\nu|\phi|^{p+1}.$$

Following proofs of Lemmas 9.2 and 9.3 for  $\phi \in \text{span}\{e\}$  in a neighborhood of  $\phi = 0$  there hold (9.3) and (9.4). By Lemma 5.3 the inequalities (9.3) and (9.4) hold for each  $\phi$  in a neighborhood of  $\phi = 0$  with another multiplicative constants. Substituting this into (9.1) by stability we get  $E_*(\phi, \nu, \Pi_h) = O(|\phi|^{[p]})$  and  $E_0(\phi, \nu, \Pi_h) = O(\nu|\phi|^{[p]+1})$ , thus  $\hat{\epsilon}(\phi, \nu, \Pi_h) = O(|\phi|^{[p]} + \nu|\phi|^{[p]+1})$ . It remains to apply Theorem 6.17.

7. If the scheme possesses the order of accuracy  $P$  in the weak sense on  $H_{per}^R(\mathbb{R}^d)$ , then by Proposition 6.22 for each  $e \in \mathring{\Omega}$  it possesses the order of accuracy  $P$  in the strong sense on  $H_{per, e}^R(\mathbb{R}^d)$ , and it remains to use the previous statement.  $\square$

An auxiliary mapping that provides (8.7) can't be generally found in the form (2.12), see Section 12.2 for the example. We just proved that this is possible if  $Q = P$ . Now we consider another two cases when this is possible, namely, the simple case and the quasi-1D case.

### 9.2. The simple case

Below we assume that the scheme is *0-exact*, i. e. the scheme preserves constant solutions, or, that is the same,  $\sum L_\eta(\Pi_1 1)_0 = 0$ . Here the subscript 1 means the substitution 1 for  $h$  and the subscript 0 means the substitution 0 for  $\eta$ . Otherwise the order of accuracy is  $P_A = -1$  and by Theorem 1 the numerical solution does not converge to the exact one.

Since the scheme is 0-exact, then the operator  $A(0)$  has nontrivial kernel:  $A(0)\vec{e} = 0$ , where  $\vec{e} = v(0, \Pi_h)$ . We shall say that the scheme is *simple* if the zero eigenvalue of the operator  $A(0)$  is simple or, this is the same,  $\dim \text{Ker} L(0) = 1$ . The scheme is simple iff the vector  $v_0(\phi, \Pi_h)$  defined in Section 9.1 is of size 1.

**Proposition 9.4.** *Suppose a scheme has no steady solutions but the constant. Then it is simple.*

*Proof.* Let  $\mathfrak{C} \in V_{per}^1$  be the constant sequence formed by the element  $\vec{e}$ . By assumption the scheme has no steady solutions but  $c\mathfrak{C}$ ,  $c \in \mathbb{R}$ . So  $\text{Ker} L \cap V_{per} = \text{span}\{\mathfrak{C}\}$ . By (6.14) we have  $A(0) = -Z^{-1}(0)L(0)$ . Consider a vector  $w_0 \in \mathbb{C}^{M^0}$  such that  $w_0 \notin \text{span}\{\vec{e}\}$ . Let  $w \in V_{per}^1$  be the constant sequence all elements of which are equal to  $w_0$ . Then  $Lw \neq 0$ . By (6.4) all the block components of  $Lw$  are equal to  $L(0)w_0$ , thus we have  $L(0)w_0 \neq 0$ . Then  $A(0)w_0 = -Z^{-1}(0)L(0)w_0 \neq 0$ . Hence,  $w_0 = \vec{e}$  is the only (up to a factor) solution of  $A(0)w_0 = 0$ . Therefore by Lemma 9.1 we obtain that  $\lambda = 0$  is a simple eigenvalue of  $A(0)$ .  $\square$

The discontinuous Galerkin method with the upwind flux in 1D is a simple scheme unless the transport velocity is zero. On a multidimensional simplicial translationally-invariant (i. e. invariant with respect to the translation by the vector of any mesh edge) mesh the DG method with the upwind flux is a simple scheme unless the transport velocity  $\omega$  is parallel to one of the mesh faces.

**Lemma 9.5.** Suppose a simple scheme (2.8) possesses the estimate (2.17) in the sense of a local mapping  $\Pi_h$ . Then in a neighborhood of  $\phi = 0$  there holds

$$\|M^{(0)}(\phi)\| \leq C|\phi|^{\lceil Q \rceil + 1}. \quad (9.10)$$

*Proof.* Denote by  $\lambda(\phi)$  the only element of the matrix  $M^{(0)}(\phi)$ . It is a holomorphic function. By Lemma 9.2 in a neighborhood of  $\phi = 0$  we have  $\|v_*(\phi, \Pi_h)\| \leq C|\phi|^P$ . Then  $v_0(0, \Pi_h) \neq 0$ ; assuming the converse, we immediately get  $(\Pi_1 1)_0 = v(0, \Pi_h) = 0$ , which contradicts the definition of the local mapping. Recall that  $v_0(\phi, \Pi_h)$  is a 1-component quantity. Hence for small  $|\phi|$  by (9.5) there holds

$$\left| e^{\nu \lambda(\phi)} - 1 \right| \leq \frac{2c_n}{|v_0(0, \Pi_h)|} (C_1 |\phi|^P + C_2 \nu |\phi|^{Q+1}).$$

Put  $\nu = |\phi|^{-\lceil Q \rceil}$ ; then the right-hand side tends to zero as  $\phi \rightarrow 0$ . Hence,  $\lambda(\phi)|\phi|^{-\lceil Q \rceil} \rightarrow 0$ . Since  $\lambda(\phi)$  is a holomorphic function of  $\phi$  at  $\phi = 0$ , then  $\lambda(\phi) = O(|\phi|^{\lceil Q \rceil + 1})$ .  $\square$

**Lemma 9.6.** Let  $\Pi_h$  be a local mapping. Let a simple scheme possess the formal order of accuracy  $P$  and the long-time simulation order  $Q$  and, at the same time, the formal order of accuracy  $P'$  and the long-time simulation order  $Q'$ . Then it possesses the formal order of accuracy  $\lceil \max\{P, P'\} \rceil$  and the long-time simulation order  $\lceil \max\{Q, Q'\} \rceil$ .

*Proof.* By Lemma 9.2 and Lemma 9.5 we have  $|\lambda(\phi)| \leq C|\phi|^{\lceil Q \rceil + 1}$  and  $\|v^*(\phi, \Pi_h)\| \leq C|\phi|^{\lceil P \rceil}$ . At the same time, we have similar inequalities with the substitution  $P'$  for  $P$  and  $Q'$  for  $Q$ . Thus we have the same inequalities with the substitution  $\max\{P, P'\}$  for  $P$  and  $\max\{Q, Q'\}$  for  $Q$ . By direct substitution into (9.2) and (9.1) we get

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq C|\phi|^{\lceil \max\{P, P'\} \rceil} + C\nu |\phi|^{\lceil \max\{Q, Q'\} \rceil + 1}.$$

It remains to apply Theorem 6.17.  $\square$

Lemma 9.6 shows that the situation mentioned in the comments to Definition 5 and showed in Section 7.6 is not possible in the simple case.

**Lemma 9.7.** Consider a simple scheme and a local mapping  $\Pi_h$ . Then the optimal value of the long-time simulation order is integer.

*Proof.* This follows from the previous lemma with  $P = P'$  and  $Q = Q'$ .  $\square$

**Lemma 9.8.** Consider a simple scheme and a local mapping  $\Pi_h$ . Let  $p \in \mathbb{N}$ . If for each  $q \in \mathbb{N}$  the scheme possesses the formal order of accuracy  $p$  and the long-time simulation order  $q$ , then it possesses the formal order of accuracy  $p$  and the long-time simulation order  $\infty$ .

*Proof.* By Lemma 9.5 for each  $s \in \mathbb{N}$  we have  $\|M^{(0)}(\phi)\| \leq C_s |\phi|^{s+1}$ ; since  $M^{(0)}(\phi)$  is holomorphic at  $\phi = 0$ , we have  $M^{(0)}(\phi) = 0$  in a neighborhood of zero. Using Lemma 9.2 and representation (9.2), (9.1), we have  $\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq c|\phi|^p$ , where  $c$  does not depend on  $\nu$ . By Theorem 6.17 the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q = \infty$ .  $\square$

**Proposition 9.9.** For a stable scheme (2.8) with  $|M^0| = 1$  and a local mapping  $\Pi_h$ , the optimal values of the order of truncation error, the order of accuracy, and the long-time simulation order coincide.

*Proof.* Let  $P$  and  $Q \geq P$  be the optimal values of the order of accuracy and of the long-time simulation order. Then by Lemma 9.5 we have  $\|M^{(0)}(\phi)\| \leq C|\phi|^{Q+1}$ . Since  $M^{(0)}$  is a  $1 \times 1$ -matrix and there is no block  $\|M^{(*)}(\phi)\|$ , we have  $\|A(\phi)\| = \|M(\phi)\| = \|M^{(0)}(\phi)\| \leq C|\phi|^{Q+1}$ . By (6.15) this yields  $\|\hat{\varepsilon}(\phi, \Pi_h)\| \leq C|\phi|^{Q+1} \|v(0, \Pi_h)\|$ . Then Theorem 6.15 states that the scheme possesses the order  $Q$  of the truncation error. By stability  $P \geq Q$ , thus  $P = Q$ . Clearly,  $Q$  is the optimal value of the order of the truncation error, otherwise by stability the scheme possesses the formal order of accuracy and the order of the long-time simulation greater than  $Q$ , which contradicts the assumption.  $\square$

**Lemma 9.10.** Let  $\Pi_h$  and  $\mathcal{P}_h$  be local mappings with kernels  $\mu_\xi, \hat{\mu}_\xi$  such that  $\langle \hat{\mu}_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ . Let a simple stable scheme of the form (2.8) possess the formal order of accuracy  $P \in \mathbb{N}$  and the long-time simulation order  $Q$  in the sense of  $\Pi_h$ . Then there exist real-valued diagonal matrices  $\mathfrak{C}^{(m)}$  such that the scheme possesses the truncation error of order  $\lceil Q \rceil$  in the sense of  $\Pi_h^{(P, \lceil Q \rceil)}$  given by (2.12).

*Proof.* We repeat a fragment of the proof of Theorem 1, using Lemma 9.5 instead of Lemma 9.3. By Lemma 9.2 and Lemma 9.5 for

$$w(\phi) = S(\phi) \begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix} S^{-1}(\phi) v(\phi, \Pi_h).$$

we have

$$\|w(\phi) - v(\phi, \Pi_h)\| = \left\| S(\phi) \begin{pmatrix} v_*(\phi, \Pi_h) \\ 0 \end{pmatrix} \right\| \leq c_n \|S(\phi)\| \|v(\phi, \Pi_h)\| = O(|\phi|^P),$$

$$A(\phi)w(\phi) = S(\phi) \begin{pmatrix} 0 \\ M^{(0)}(\phi)v_0(\phi, \Pi_h) \end{pmatrix} = O(|\phi|^{\lceil Q \rceil + 1}).$$

By Lemma 6.7 there exist  $\mathfrak{C}^{(m)}$ ,  $P \leq |m| \leq \lceil Q \rceil$ , such that for  $\Pi_h^{(P, \lceil Q \rceil)}$  there holds

$$v(\phi, \Pi_h^{(P, \lceil Q \rceil)}) = w(\phi) + O(|\phi|^{\lceil Q \rceil + 1}).$$

By the triangle inequality we get

$$\hat{\epsilon}(\phi, \Pi_h^{(P, Q)}) = A(\phi)v(\phi, \Pi_h^{(P, Q)}) = O(|\phi|^{\lceil Q \rceil + 1}).$$

By Theorem 6.15 the scheme possesses the truncation error of order  $\lceil Q \rceil$  in the sense of  $\Pi_h^{(P, Q)}$ , so the coefficients  $\mathfrak{C}^{(m)}$  satisfy the statement of the lemma.  $\square$

**Lemma 9.11.** Let the scheme (2.8) be simple,  $\Pi_h$  be a local mapping with a kernel  $\mu_\xi \in (W_2^r(G))^*$  (or  $(C^r(G))^*$ ), and let  $p, q > 0$ . Suppose for each  $e \in \dot{\Omega}_{\lceil q \rceil}$  the scheme possesses the formal order of accuracy  $p$  and the long-time simulation order  $q$  on  $H_{per, e}^R(\mathbb{R}^d)$  (or  $C_{per, e}^{\lceil R \rceil}(\mathbb{R}^d)$ ), Then the scheme possesses the formal order of accuracy  $p$  and the long-time simulation order  $q$  on  $H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^{\lceil R \rceil}(\mathbb{R}^d)$ ),  $R = \max\{r, q + 1\}$ .

*Proof.* By assumption for each  $e \in \dot{\Omega}_p$  and each  $\phi$  aligned with  $e$  we have

$$\|\hat{\epsilon}(\phi, \nu, \Pi_h)\| \leq c_1(e)|\phi|^p + c_2(e)\nu|\phi|^{q+1}.$$

Then for  $\phi$  aligned with  $e$  there hold (9.3) and (9.10). By Lemma 5.3 the inequalities (9.3) and (9.10) hold for each  $\phi$  with some other multiplicative constants. Substituting this into (9.1) by stability we get  $E_*(\phi, \nu, \Pi_h) = O(|\phi|^{\lceil p \rceil})$  and  $E_0(\phi, \nu, \Pi_h) = O(\nu|\phi|^{\lceil q \rceil + 1})$ , thus  $\hat{\epsilon}(\phi, \nu, \Pi_h) = O(|\phi|^{\lceil p \rceil} + \nu|\phi|^{\lceil q \rceil + 1})$ . It remains to apply Theorem 6.17.  $\square$

**Lemma 9.12.** Let the scheme be simple. Then the optimal values of the formal order of accuracy and of the long-time simulation order coincide with the ones in the weak sense.

*Proof.* This directly follows from Proposition 6.22 and Lemma 9.11.  $\square$

**Lemma 9.13.** Suppose the scheme is simple, and  $\Pi_h$  is a local mapping. Let  $P_A$ ,  $P > 0$ , and  $Q$  be the optimal values of the truncation error, of the formal order of accuracy and of the long-time simulation order. Then either  $P_A = P = Q$  or  $Q \geq P = P_A + 1$ .

*Proof.* If  $|M^0| = 1$ , then the statement of the lemma follows from Proposition 9.9. So without loss we assume  $|M^0| > 1$ . By Theorem 1 the values  $P$  and  $Q$  are integers and there holds either  $P = P_A$  or  $P = P_A + 1$ . So we need only to prove that  $Q \geq P_A + 1$  implies  $P = P_A + 1$ .

Since

$$\|\hat{\epsilon}(\phi, \Pi_h)\| = \left\| S(\phi) \begin{pmatrix} M^{(*)}(\phi)v_*(\phi, \Pi_h) \\ M^{(0)}(\phi)v_0(\phi, \Pi_h) \end{pmatrix} \right\| = O(|\phi|^{P_A + 1}),$$

we have  $\|M^{(*)}(\phi)v_*(\phi, \Pi_h)\| = O(|\phi|^{P_A+1})$ . Then  $\|v_*(\phi, \Pi_h)\| = O(|\phi|^{P_A+1})$  and by stability

$$\|E_*(\phi, \nu, \Pi_h)\| \leq C|\phi|^{P_A+1}.$$

Lemma 9.5 yields  $\|M^{(0)}(\phi)\| = O(|\phi|^{Q+1})$  and hence

$$\|E_0(\phi, \nu, \Pi_h)\| \leq \tilde{C}|\phi|^{Q+1}\nu.$$

Combining these estimates we get

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq C'|\phi|^{P_A+1} + C'\nu|\phi|^{Q+1}.$$

If  $Q \geq P_A + 1$  then by Theorem 6.17 the scheme possesses the formal order of accuracy  $P_A + 1$ , thus  $P = P_A + 1$ .  $\square$

**Theorem 3A.** *For a simple scheme, Theorem 3 holds.*

*Proof.* For a simple scheme, the five statements of Theorem 3 follow from Lemmas 9.7, 9.10, 9.11, 9.12, 9.13, in respective order.  $\square$

705 The following proposition may be useful in a computer-based analysis of a scheme.

**Proposition 9.14.** *Let the scheme (2.8) be simple and  $\Pi_h$  be a local mapping. Let  $Q \in \mathbb{N}$  be such that  $|\det A(\phi)| \leq c|\phi|^{q+1}$  holds for  $q = Q$  and does not hold for  $q = Q + 1$ . Let  $P_A$  be the optimal value of the order of the truncation error. Let  $P = P_A + 1$  if  $Q > P_A$  and  $P = P_A$  otherwise. Then  $P$  and  $Q$  the optimal values of the formal order of accuracy and of the long-time simulation order.*

*Proof.* Denote by  $\lambda(\phi)$  the only element of the matrix  $M^{(0)}(\phi)$ . Since the eigenvalue  $\lambda = 0$  of  $A(0)$  is simple, in a neighborhood of  $\phi = 0$  there holds  $|\lambda(\phi)| \leq 2cc_{\Pi}^{-1}|\phi|^{Q+1}$ , where  $c_{\Pi}$  is the product of all nonzero eigenvalues of  $A(0)$  taking into account their algebraic multiplicity. Substituting this estimate for  $M^{(0)}(\phi) = \lambda(\phi)$  and (9.3) into (9.2) and (9.1) we get

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq C|\phi|^p + C\nu|\phi|^{Q+1},$$

710 where  $p$  is the optimal value of the formal order of accuracy. Thus by Theorem 6.17 the scheme possesses the formal order of accuracy  $p$  and the long-time simulation order  $Q$ .

Now we claim that the values  $p$  and  $Q$  are optimal. Assume the converse. Let  $Q' > Q$  be the optimal value of the long-time simulation order. Then by Lemma 9.5 there holds  $|\lambda(\phi)| \leq C|\phi|^{Q'+1}$ , that contradicts the assumption.

The fact that  $p = P$  follows from Lemma 9.13.  $\square$

### 715 9.3. The quasi-one-dimensional case

Recall that a scheme of the form (2.8) is *quasi-1D* if the stencil  $\mathcal{S}$  of the scheme belongs to a 1D subset of  $\mathbb{Z}^d$ , i. e. there exists  $\boldsymbol{\eta} \in \mathbb{Z}^d$  such that  $\mathcal{S} \subset \{m\boldsymbol{\eta}, m = -M, \dots, M\}$ . The scheme is quasi-1D iff the matrix  $A(\phi)$  depends only on  $\phi \cdot \boldsymbol{e}$  for some  $\boldsymbol{e} \in \mathbb{R}^d$ . In this section we will use the notation  $\boldsymbol{e}$  for this vector.

720 In 1D case (when  $d = 1$ ) every scheme of the form (2.8) is quasi-1D. In the multidimensional case an example of a quasi-1D scheme is the DG method on simplicial translationally-invariant meshes when the vector of the transport velocity is collinear to one of the mesh edges.

**Theorem 9.15** ([30], §3.5, Corollary 3; [26], §2.6.2). *Let  $A(\psi)$  be a holomorphic and selfadjoint  $N \times N$ -matrix defined on the interval  $[\psi_1, \psi_2]$ . Then there exists a unitary matrix  $\mathcal{U}(\psi)$ , holomorphic on  $[\psi_1, \psi_2]$ , such that  $\mathcal{U}^{-1}(\psi)A(\psi)\mathcal{U}(\psi) = \text{diag}\{\lambda_1(\psi), \dots, \lambda_N(\psi)\}$ .*

**Lemma 9.16.** *Let  $f(\psi)$  and  $g(\phi)$  be holomorphic functions at  $\psi = 0$  and  $\phi = 0$ , correspondingly. Let  $\boldsymbol{e} \in \Omega$ . Suppose there hold  $f(\psi) \sim \psi^a$  as  $\psi \rightarrow 0$  and  $|f(\phi \cdot \boldsymbol{e})g(\phi)| \leq c|\phi|^{a+b}$ , where  $a, b \in \mathbb{N} \cup \{0\}$ . Then in a neighborhood of  $\phi = 0$  there holds*

$$|f(\phi \cdot \boldsymbol{e})g(\phi)| \leq \delta c |\phi \cdot \boldsymbol{e}|^a |\phi|^b, \quad (9.11)$$

725 where  $\delta$  depends on  $a + b$  and the space dimension  $d$  only.

*Proof.* Let  $\check{C} = \{x \in \mathbb{R}^d : e \cdot x > |x|/2\}$ . By Lemma A.6 there exist  $e_k \in \check{C} \cap \Omega$ ,  $k = 1, \dots, C_{b+d-1}^{d-1}$ , such that  $(e_k \cdot r)^b$  form a basis in the space of homogeneous polynomials of order  $b$ . For different  $e$  the set of vectors  $e_k$  can be chosen the same modulo rotation. By assumption for each  $e' \in \check{C} \cap \Omega$  in a neighborhood of  $\psi = 0$  there holds

$$|g(e'\psi)| \leq \frac{c|\psi|^{a+b}}{f((e \cdot e')\psi)} \leq 2 \frac{c|\psi|^b}{(e \cdot e')^a} \leq 2^{1+a} c |\psi|^b.$$

Thus  $|\partial^n g(0)/\partial e'^n| = 0$  for  $n < b$  and  $|\partial^b g(0)/\partial e'^b| \leq 2^{1+a} b! c$ . By Lemma 5.1 we have  $D^m g(0) = 0$  for  $|m| < b$  and  $|D^m g(0)| \leq \delta c$  for  $|m| = b$ , where  $\delta = 2^{1+a} b! \sum_k |\gamma_k^{(m)}|$  does not depend on  $e$ .  $\square$

**Lemma 9.17.** *Let the scheme be quasi-1D,  $\Pi_h$  be a local mapping,  $0 < P \leq Q < \infty$ . Then the function  $\mathcal{F}^{P,Q}(\phi)$  defined by (8.9) is bounded as  $\phi \rightarrow 0$  iff the function*

$$\tilde{\mathcal{F}}^{P,Q}(\phi) = |\phi|^{-2\lceil P \rceil} v^*(\phi, \Pi_h) \left( A^*(\phi) A(\phi) + |\phi \cdot e|^{2(\lceil Q \rceil + 1 - \lceil P \rceil)} \right)^{-1} A^*(\phi) A(\phi) v(\phi, \Pi_h).$$

is bounded as  $\phi \rightarrow 0$ .

*Proof.* Throughout this proof we will use the notation  $\psi \equiv \phi \cdot e$ . By Theorem 9.15 there exist holomorphic matrix function  $\mathcal{U}$  and holomorphic functions  $\lambda_j$  such that  $\mathcal{U}^{-1}(\psi) A^*(\phi) A(\phi) \mathcal{U}(\psi) = \text{diag}\{\lambda_j(\psi)\}$ . Therefore,

$$\begin{aligned} \mathcal{F}^{P,Q}(\phi) &= |\phi|^{-2P} \sum_{j \in M^0} |(\mathcal{U}^{-1}(\psi) v(\phi, \Pi_h))_j|^2 \frac{\lambda_j(\psi)}{\lambda_j(\psi) + |\phi|^{2(Q+1-P)}}, \\ \tilde{\mathcal{F}}^{P,Q}(\phi) &= |\phi|^{-2\lceil P \rceil} \sum_{j \in M^0} |(\mathcal{U}^{-1}(\psi) v(\phi, \Pi_h))_j|^2 \frac{\lambda_j(\psi)}{\lambda_j(\psi) + |\psi|^{2(\lceil Q \rceil + 1 - \lceil P \rceil)}}. \end{aligned}$$

First assume that  $\tilde{\mathcal{F}}^{P,Q}(\phi)$  is bounded. Then we have

$$\mathcal{F}^{\lceil P \rceil, \lceil Q \rceil}(\phi) \leq \tilde{\mathcal{F}}^{\lceil P \rceil, \lceil Q \rceil}(\phi) = \tilde{\mathcal{F}}^{P,Q}(\phi) < \infty.$$

Thus by Theorem 8.5 the scheme possesses the formal order of accuracy  $\lceil P \rceil$  and the long-time simulation order  $\lceil Q \rceil$ . By Corollary 6.18 it possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$ . Using again Theorem 8.5, we obtain that  $\mathcal{F}^{P,Q}(\phi)$  then is bounded as  $\phi$  tends to zero.

Now assume that  $\mathcal{F}^{P,Q}(\phi)$  is bounded in a neighborhood of  $\phi = 0$ . Then for each  $j \in M^0$  there holds

$$\lambda_j(\psi) |(\mathcal{U}^{-1}(\psi) v(\phi, \Pi_h))_j|^2 \leq C |\phi|^{2P} (\lambda_j(\psi) + |\phi|^{2(Q+1-P)}) \quad (9.12)$$

in a neighborhood of  $\phi = 0$ . The function  $\lambda_j(\psi)$  is holomorphic at  $\psi = 0$  and non-negative (as the eigenvalue of matrix  $A^*(\psi) A(\psi)$ ). Thus either  $\lambda_j(\psi) = c_\lambda \psi^{2q_j} (1 + O(\psi))$  as  $\psi \rightarrow 0$  with some  $c_\lambda > 0$  and  $q_j \in \mathbb{N} \cup \{0\}$  or  $\lambda_j(\psi) \equiv 0$  (in this case we put  $q_j = +\infty$ ).

The function  $(\mathcal{U}^{-1}(\phi \cdot e) v(\phi, \Pi_h))_j$  is holomorphic at  $\phi = 0$ , let  $r_j$  be the lowest order of the terms in its Taylor series. From (9.12) we have

$$2q_j + 2r_j \geq 2P + \min\{2q_j, 2(Q+1-P)\}$$

or, equivalently,

$$q_j + r_j \geq \min\{q_j + P, Q + 1\}.$$

Since  $r_j, q_j \in \mathbb{N}$ , we have

$$q_j + r_j \geq \min\{q_j + \lceil P \rceil, \lceil Q \rceil + 1\}.$$

If  $q_j + \lceil P \rceil < \lceil Q \rceil + 1$ , then  $r_j \geq \lceil P \rceil$ , so  $|(\mathcal{U}^{-1}(\psi) v(\phi, \Pi_h))_j|^2 \leq \tilde{c} |\phi|^{2\lceil P \rceil}$  and thus

$$\lambda_j(\psi) |(\mathcal{U}^{-1}(\psi) v(\phi, \Pi_h))_j|^2 \leq \tilde{c} |\phi|^{2\lceil P \rceil} \lambda_j(\psi).$$

If  $q_j + \lceil P \rceil \geq \lceil Q \rceil + 1$ , then  $r_j \geq \lceil Q \rceil + 1 - q_j$  and

$$\begin{aligned} & \lambda_j(\psi) |(\mathcal{U}^{-1}(\psi)v(\phi, \Pi_h))_j|^2 \leq \tilde{c} \psi^{2q_j} |\phi|^{2(\lceil Q \rceil + 1 - q_j)} = \\ & = \tilde{c} \psi^{2(\lceil Q \rceil + 1 - \lceil P \rceil)} \psi^{2(q_j - (\lceil Q \rceil + 1 - \lceil P \rceil))} |\phi|^{2(\lceil Q \rceil + 1 - q_j)} \leq \tilde{c} \psi^{2(\lceil Q \rceil + 1 - \lceil P \rceil)} |\phi|^{2\lceil P \rceil}. \end{aligned}$$

Combining the results for both cases, in a neighborhood of  $\phi = 0$  we get

$$\lambda_j(\psi) |(\mathcal{U}^{-1}(\psi)v(\phi, \Pi_h))_j|^2 \leq \tilde{c} |\phi|^{2\lceil P \rceil} (\lambda_j(\psi) + |\psi|^{2(\lceil Q \rceil + 1 - \lceil P \rceil)}). \quad (9.13)$$

735 From here, the statement of the lemma is obvious.  $\square$

**Lemma 9.18.** *Suppose the scheme is quasi-1D and  $\Pi_h$  is a local mapping. Then the optimal values of the formal order of accuracy and of the long-time simulation order are integers.*

*Proof.* The optimal order of accuracy  $P$  is integer by Theorem 1. Let the scheme possess the formal order of accuracy  $P$  and a long-time simulation order  $Q$ . By Theorem 8.5,  $\mathcal{F}^{P,Q}(\phi)$  is bounded as  $\phi \rightarrow 0$ . By Lemma 9.17 so is  $\tilde{\mathcal{F}}^{P,Q}(\phi)$  and, applying this Lemma again, so is  $\mathcal{F}^{P,\lceil Q \rceil}(\phi)$ . By Theorem 8.5 the scheme possesses the formal order  $P$  and the long-time simulation order  $\lceil Q \rceil$ .  $\square$

**Lemma 9.19.** *Let the scheme be quasi-1D,  $\Pi_h$  be a local mapping,  $P \in \mathbb{N}$ . Let the scheme possess the formal order of accuracy  $P$  and the long-time simulation order  $s$  for each  $s \geq P$ . Then the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $\infty$ .*

*Proof.* By Theorem 8.5 the functional  $\mathcal{F}^{P,s}(\phi)$  is bounded in a neighborhood of  $\phi = 0$ . Thus using the notation of Lemma 9.17 for each  $j \in M^0$  and  $s \geq P$  there holds

$$q_j + r_j \geq \min\{q_j + P, s + 1\}.$$

Thus for each  $j$  there holds either  $q_j = +\infty$  or  $r_j \geq P$ . Put  $w(\phi) = \mathcal{U}(\psi) \text{diag}\{\chi_j\} \mathcal{U}^{-1}(\psi) v(\phi, \Pi_h)$  where  $\chi_j = 0$  if  $r_j \geq P$  and  $\chi_j = 1$  otherwise. Then

$$\|w(\phi) - v(\phi, \Pi_h)\| \leq c |\phi|^P, \quad A(\phi)w(\phi) = 0,$$

745 and it remains to apply Lemma 8.1.  $\square$

**Lemma 9.20.** *Let  $\Pi_h$  and  $\mathcal{P}_h$  be local mappings with kernels  $\mu_\xi, \hat{\mu}_\xi \in (W_2^q(G))^*$  (or  $(C^q(G))^*$ ), such that  $\langle \hat{\mu}_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ . Suppose the scheme (2.8) is quasi-1D and possesses the formal order of accuracy  $P \in \mathbb{N}$  and the long-time simulation order  $Q \in \mathbb{N}$  in the sense of  $\Pi_h$  on  $H_{\text{per}}^R(\mathbb{R}^d)$  (or  $C_{\text{per}}^R(\mathbb{R}^d)$ ), where  $R = \max\{Q + 1, q\}$ . Then there exist real-valued diagonal matrices  $\mathfrak{C}^{(m)}$  such that the scheme possesses the truncation error of order  $Q$  in the sense of  $\Pi_h^{(P,Q)}$  given by (2.12).*

*Proof.* By Theorem 8.5,  $\mathcal{F}^{P,Q}(\phi)$  is bounded as  $\phi \rightarrow 0$ . By Lemma 9.17 so does  $\tilde{\mathcal{F}}^{P,Q}(\phi)$ . Put

$$w(\phi) = \left( A^*(\phi)A(\phi) + |\phi \cdot e|^{2(Q+1-P)} \right)^{-1} |\phi \cdot e|^{2(Q+1-P)} v(\phi, \Pi_h) \quad (9.14)$$

and note that

$$\tilde{\mathcal{F}}^{P,Q}(\phi) = |\phi|^{-2P} \|v(\phi, \Pi_h) - w(\phi)\|^2 + |\phi|^{-2P} |\phi \cdot e|^{-2(Q-P+1)} \|A(\phi)w(\phi)\|^2. \quad (9.15)$$

(the proof of this fact repeats the one of Lemma 8.4). By construction,  $w(\phi) = W(\phi)/\mathfrak{d}(\psi)$ , there  $W(\phi)$  is holomorphic at  $\phi = 0$  and  $\mathfrak{d}(\psi) = \det(A^*(\psi e)A(\psi e) + \psi^{2(Q+1-P)})$ . Let  $\mathfrak{d}(\psi) \sim c\psi^q$  as  $\psi \rightarrow 0$ , where  $c \neq 0$ . All the terms in the Taylor expansion of the numerator contain the multiplier  $\psi^q$ , otherwise  $w(\phi)$  is unbounded and so does  $\tilde{\mathcal{F}}^{P,Q}(\phi)$ , which contradicts the assumption. This proves that  $w(\phi)$  is analytical at  $\phi = 0$ .

From (9.15) and the boundedness of  $\tilde{\mathcal{F}}^{P,Q}$  it follows that  $\|v(\phi, \Pi_h) - w(\phi)\| \leq \tilde{C} |\phi|^P$  and

$$\|A(\phi)w(\phi)\| |\phi \cdot e|^P \leq \tilde{C} |\phi \cdot e|^{Q+1} |\phi|^P \leq \tilde{C} |\phi|^{Q+1+P}.$$

755 Since  $A(\phi)w(\phi)$  is holomorphic, this implies  $A(\phi)w(\phi) = O(|\phi|^{Q+1})$ .  
 By Lemma 6.7 there exist  $\mathfrak{C}^{(m)}$ ,  $P \leq |\mathbf{m}| \leq Q$ , such that for  $\Pi_h^{(P,Q)}$  there holds

$$v(\phi, \Pi_h^{(P,Q)}) = w(\phi) + O(|\phi|^{Q+1}).$$

By the triangle inequality we get

$$\hat{\epsilon}(\phi, \Pi_h^{(P,Q)}) = A(\phi)v(\phi, \Pi_h^{(P,Q)}) = O(|\phi|^{Q+1}).$$

By Theorem 6.15 the scheme possesses the truncation error of order  $Q$  in the sense of  $\Pi_h^{(P,Q)}$ , so the coefficients  $\mathfrak{C}^{(m)}$  satisfy the statement of the lemma.  $\square$

760 **Lemma 9.21.** *Let the scheme (2.8) be quasi-1D,  $\Pi_h$  be a local mapping with a kernel  $\mu_\xi \in (W_2^r(G))^*$  (or  $(C^r(G))^*$ ), and let  $0 < P \leq Q < \infty$ . Suppose for each  $\tilde{e} \in \mathring{\Omega}_{[Q]}$  the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per,\tilde{e}}^R(\mathbb{R}^d)$  (or  $C_{per,\tilde{e}}^{[R]}(\mathbb{R}^d)$ ). Then the scheme possesses the formal order of accuracy  $[P]$  and the long-time simulation order  $[Q]$  on  $H_{per}^R(\mathbb{R}^d)$  (or  $C_{per}^{[R]}(\mathbb{R}^d)$ ).*

*Proof.* We will use the notation of Lemma 9.17. By assumption, for each direction  $\tilde{e} \in \mathring{\Omega}_{[Q]}$  the scheme possesses the order of accuracy  $P$  and the long-time simulation order  $Q$  on  $H_{per,\tilde{e}}^s(\mathbb{R}^d)$  for  $s$  large enough. By Lemma 8.6, the functional  $\mathcal{F}^{P,Q}(\psi\tilde{e})$  is bounded as  $\psi$  tends to zero. Thus for each  $j \in M^0$  and each  $\tilde{e} \in \mathring{\Omega}_Q$  there holds

$$\sqrt{\lambda_j(\psi)}|(\mathcal{U}^{-1}(\psi)v(\phi, \Pi_h))_j| \leq C_{\tilde{e}}|\phi|^P(\sqrt{\lambda_j(\psi)} + |\phi|^{Q+1-P}) \leq 2C_{\tilde{e}}|\phi|^{\min\{Q+1, P+q_j\}},$$

where  $\phi$  belongs to a neighborhood of zero in  $\text{span}\{\tilde{e}\}$ . Since the function on the left-hand side is the absolute value of a function analytical at  $\phi = 0$ , by Lemma 5.3 in a neighborhood of zero we have

$$\sqrt{\lambda_j(\psi)}|(\mathcal{U}^{-1}(\psi)v(\phi, \Pi_h))_j| \leq C|\phi|^{\min\{[Q]+1, [P]+q_j\}},$$

where  $C$  does not depend on  $e$ .

If for a particular  $j$  there holds  $[P] \leq [Q] + 1 - q_j$ , then we have

$$|(\mathcal{U}^{-1}(\psi)v(\phi, \Pi_h))_j| \leq \tilde{c}|\phi|^{[P]}$$

and hence (9.13). If  $q_j > [Q] + 1 - [P]$ , then by Lemma 9.16 we have

$$\lambda_j(\psi)|(\mathcal{U}^{-1}(\psi)v(\phi, \Pi_h))_j|^2 \leq c'\psi^{2q_j}|\phi|^{2([Q]+1-q_j)} \leq c'\psi^{2([Q]+1-[P])}|\phi|^{2[P]},$$

and we again have (9.13). Thus  $\tilde{\mathcal{F}}^{[P],[Q]}(\phi)$  is bounded as  $\phi \rightarrow 0$ . By Lemma 9.17 so does  $\mathcal{F}^{[P],[Q]}(\phi)$  and by Theorem 8.5 the scheme possesses the formal order of accuracy  $[P]$  and the long-time simulation order  $[Q]$ .  $\square$

765 **Lemma 9.22.** *Suppose the scheme is quasi-1D and  $\Pi_h$  is a local mapping. Then the optimal values of the formal order of accuracy and of the long-time simulation order coincide with the ones in the weak sense.*

*Proof.* Let the scheme possess the formal order of accuracy  $P$  and the long-time simulation order  $Q$  in the weak sense. By Proposition 6.22 for each  $\tilde{e} \in \mathring{\Omega}$  the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order in the strong sense on  $H_{per,\tilde{e}}^R(\mathbb{R}^d)$ . It remains to use Lemma 9.21.  $\square$

770 **Theorem 3B.** *For a quasi-1D scheme, Theorem 3 holds.*

*Proof.* For a quasi-1D scheme, the first four statements of Theorem 3 follow from Lemmas 9.18, 9.20, 9.21, 9.22, in respective order. The fifth statement of Theorem 3, which only concerns only simple schemes, has been already proved in Theorem 3A.  $\square$

## 10. Algorithms for scheme analysis

Throughout this section we consider a scheme of the form (2.8), a local mapping  $\Pi_h$ , and assume that the scheme is stable and 0-exact. For the cases where the existence of a local mapping of the form (2.12) is proved it is possible to construct algorithms for the analysis of the scheme.

### 10.1. The order of the truncation error

**Lemma 10.1.** *Let  $P_A \in \mathbb{N} \cup \{0\}$ . Then the following statements are equivalent.*

1. *The scheme possesses the truncation error of order  $P_A$ .*
2. *There holds  $\|\hat{\epsilon}(\phi, \Pi_h)\| = O(|\phi|^{P_A+1})$  as  $\phi \rightarrow 0$ .*
3. *For each multiindex  $\mathbf{m}$  such that  $|\mathbf{m}| \leq P_A$  there holds  $(\epsilon_1(\mathbf{r}^{\mathbf{m}}/\mathbf{m}!, \Pi_1))_0 = 0$ .*

*Proof.* The equivalence of statements 1 and 2 is by Lemma 6.10 and boundedness of  $Z(\phi)$  and  $Z^{-1}(\phi)$ . Now we prove the equivalence of statements 2 and 3.

Denote by  $\tilde{\Pi}_h$  the operator taking  $f$  to  $\tilde{\Pi}_h f = h\epsilon_h(f, \Pi_h)$ . By (6.19) we have

$$Z(\phi)\hat{\epsilon}(\phi, \Pi_h) = -F[h\epsilon_h(e^{i\phi \cdot \mathbf{r}/h}, \Pi_h)](\phi) = -F[\tilde{\Pi}_h e^{i\phi \cdot \mathbf{r}/h}](\phi).$$

The mapping  $\tilde{\Pi}_h$  is of the form (2.11) and enjoys all the properties of a local mapping with the only exception that  $\tilde{\Pi}_h 1 = 0$  (recall that  $P_A \geq 0$ ). Particularly, Lemma 6.5 remains valid for  $\tilde{\Pi}_h$ . Then by Lemma 6.5 we obtain

$$\begin{aligned} Z(\phi)\hat{\epsilon}(\phi, \Pi_h) &= -v(\phi, \tilde{\Pi}_h) = -(\tilde{\Pi}_1 e^{i\phi \cdot \mathbf{r}})_0 = \\ &= - \sum_{0 \leq |\mathbf{m}| < \infty} \phi^{\mathbf{m}} \left( \tilde{\Pi}_1 \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!} \right)_0 = - \sum_{0 \leq |\mathbf{m}| < \infty} \phi^{\mathbf{m}} \left( \epsilon_1 \left( \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!}, \Pi_1 \right) \right)_0. \end{aligned}$$

Since  $Z(\phi)$  and  $Z^{-1}(\phi)$  are uniformly bounded, the equivalence of statements 2 and 3 is obvious.  $\square$

**Lemma 10.2.** *Let  $P_A \in \mathbb{N} \cup \{0\}$  and  $\mathbf{e} \in \Omega$ . Then  $\|\hat{\epsilon}(\mathbf{e}\psi, \Pi_h)\| = O(|\psi|^{P_A+1})$  as  $\psi \rightarrow 0$  iff for each  $m = 0, \dots, P_A$  there holds  $(\epsilon_1((\mathbf{r} \cdot \mathbf{e})^m/\mathbf{m}!, \Pi_1))_0 = 0$ .*

*Proof.* The proof repeats the one of the previous lemma.  $\square$

Recall the following standard algorithm.

**Algorithm 1** (to detect the optimal order of the truncation error).

1. Put  $P_A = 0$ .
2. Compute  $f^{\mathbf{m}} = -(\epsilon_1(\mathbf{r}^{\mathbf{m}}/\mathbf{m}!, \Pi_1))_0$  for each  $\mathbf{m}$  such that  $|\mathbf{m}| = P_A + 1$ .
3. If for each  $\mathbf{m}$  with  $|\mathbf{m}| = P_A + 1$  there holds  $f^{\mathbf{m}} = 0$  then put  $P_A = P_A + 1$  and return to the step 2.

**Proposition 10.3.** *If Algorithm 1 returns the value  $P_A$ , then  $P_A$  is the optimal order of the truncation error. If Algorithm 1 loops endlessly, then the scheme is exact, i. e.  $\epsilon_h(f, \Pi_h) = 0$  for each  $f$ .*

*Proof.* By Theorem 1 the optimal order of the truncation error is integer, so the first statement follows from Lemma 10.1. If the algorithm loops endlessly, then by Lemma 10.1 there holds  $\|\hat{\epsilon}(\phi, \Pi_h)\| = O(|\phi|^s)$  for each  $s \in \mathbb{N}$ . Since  $\hat{\epsilon}(\phi, \Pi_h)$  is holomorphic at  $\phi = 0$ , we have  $\hat{\epsilon}(\phi, \Pi_h) = 0$  in a neighborhood of  $\phi = 0$ , and it remains to apply Lemma 6.16.  $\square$



800 **10.2. The order of the truncation error in the sense of auxiliary mappings**

Let  $\Pi_h$  and  $\mathcal{P}_h$  be local mappings, and the kernel of  $\mathcal{P}_h$  satisfy  $\langle \hat{\mu}_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ . Let  $\mathfrak{C}^{(\mathbf{m})}$  be real-valued diagonal matrices. Now we discuss the truncation error in the sense of the mappings  $\Pi_h^{(p,q)}$  given by (2.12).

**Lemma 10.4.** *Let  $\Pi_h^{(p,q)}$  be an operator of the form (2.12) with coefficients  $\mathfrak{C}^{(\mathbf{m})}$ . The scheme possesses the truncation error of order  $P_A \in \mathbb{N}$  in the sense of  $\Pi_h^{(p,q)}$  if and only if for each  $|\mathbf{m}| \leq P_A$  there holds*

$$L(0)\mathfrak{C}^{(\mathbf{m})}(\mathcal{P}_1 1)_0 = - \left( \epsilon_1 \left( \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!}, \Pi_1^{(p, \min\{|\mathbf{m}|-1, q\})} \right) \right)_0, \quad (10.1)$$

where the mappings  $\Pi^{(p, \min\{|\mathbf{m}|-1, q\})}$  are based on the same coefficients  $\mathfrak{C}^{(\mathbf{m})}$ .

805 *Proof.* If  $q < P_A$ , put  $\mathfrak{C}^{\mathbf{m}} = 0$  for each  $q < |\mathbf{m}| \leq P_A$ , so without loss we can assume that  $q \geq P_A$ .

For each  $|\mathbf{m}| \leq P_A$  we have

$$\begin{aligned} \left( \epsilon_1 \left( \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!}, \Pi_1^{(p, P_A)} \right) \right)_0 &= \left( -Z\Pi_1^{(p, P_A)}(\omega \cdot \nabla) \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!} + L\Pi_1^{(p, P_A)} \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!} \right)_0 = \\ &= - \sum_{\eta \in \mathcal{S}} Z_\eta \left( \Pi_1^{(p, P_A)}(\omega \cdot \nabla) \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!} \right)_\eta + \sum_{\eta \in \mathcal{S}} L_\eta \left( \Pi_1^{(p, P_A)} \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!} \right)_\eta. \end{aligned}$$

The function  $(\omega \cdot \nabla)\mathbf{r}^{\mathbf{m}}$  is a polynomial of order not higher than  $|\mathbf{m}| - 1$ ; thus  $\Pi_1^{(p, P_A)}$  acts on it the same way as  $\Pi_1^{(p, |\mathbf{m}|-1)}$ . Besides,

$$\left( \Pi_1^{(p, P_A)} \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!} \right)_\eta = \left( \Pi_1^{(p, |\mathbf{m}|-1)} \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!} \right)_\eta + \mathfrak{C}^{(\mathbf{m})}(\mathcal{P}_1 1)_\eta = \left( \Pi_1^{(p, |\mathbf{m}|-1)} \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!} \right)_\eta + \mathfrak{C}^{(\mathbf{m})}(\mathcal{P}_1 1)_0.$$

Therefore, using  $L(0) = \sum_\eta L_\eta$  we have

$$\left( \epsilon_1 \left( \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!}, \Pi_1^{(p, P_A)} \right) \right)_0 = \left( \epsilon_1 \left( \frac{\mathbf{r}^{\mathbf{m}}}{\mathbf{m}!}, \Pi_1^{(p, |\mathbf{m}|-1)} \right) \right)_0 + L(0)\mathfrak{C}^{(\mathbf{m})}(\mathcal{P}_1 1)_0.$$

By Lemma 10.1 the scheme possesses the truncation error of order  $P_A$  iff the left-hand side of this identity is zero for each  $|\mathbf{m}| \leq P_A$ . Thus we get the statement of the lemma.  $\square$

In 1D case equation (10.1) takes a form

$$\begin{aligned} L(0)\mathfrak{C}^{(m)}(\mathcal{P}_1 1)_0 &= \omega \sum_{\eta \in \mathcal{S}} Z_\eta \left[ \left( \Pi_1 \frac{x^{m-1}}{(m-1)!} \right)_\eta + \sum_{n=p}^{\min\{m-1, q\}} \mathfrak{C}^{(n)} \left( \mathcal{P}_1 \frac{x^{m-n-1}}{(m-n-1)!} \right)_\eta \right] - \\ &- \sum_{\eta \in \mathcal{S}} L_\eta \left[ \left( \Pi_1 \frac{x^m}{m!} \right)_\eta + \sum_{n=p}^{\min\{m-1, q\}} \mathfrak{C}^{(n)} \left( \mathcal{P}_1 \frac{x^{m-n}}{(m-n)!} \right)_\eta \right]. \end{aligned}$$

**10.3. The order of accuracy**

**Algorithm 2** (to detect the order of accuracy).

- 810 1. Detect the order of the truncation error  $P_A$  using Algorithm 1.  
 2. Compute  $f^{\mathbf{m}} = -(\epsilon_1(\mathbf{r}^{\mathbf{m}}/\mathbf{m}!, \Pi_1))_0$  for each  $\mathbf{m}$  such that  $|\mathbf{m}| = P_A + 1$ .  
 3. If for each  $\mathbf{m}$  with  $|\mathbf{m}| = P_A + 1$  there holds  $f^{\mathbf{m}} \in \text{Im}L(0)$  then put  $P = P_A + 1$ , otherwise put  $P = P_A$ .

**Theorem 10.5.** *The value  $P$  given by Algorithm 2 coincides with the optimal value of the order of accuracy.*

815 *Proof.* Let  $\mathcal{P}_h$  be any local mapping such that  $(\mathcal{P}_1 1)_{0, \xi} \neq 0$  for each  $\xi \in M^0$ . Obviously, the algorithm returns the value  $P = P_A + 1$  iff for each multiindex  $\mathbf{m}$  such that  $|\mathbf{m}| = P_A + 1$  the system  $L(0)\mathfrak{C}^{(\mathbf{m})}(\mathcal{P}_1 1)_0 = f^{(\mathbf{m})}$  (as a system for the coefficients of  $\mathfrak{C}^{(\mathbf{m})}$ ) is consistent. By Lemma 10.4 this is equivalent to the existence of  $\mathfrak{C}^{(\mathbf{m})}$  such that the scheme possesses the order  $P_A + 1$  of the truncation error in the sense of  $\Pi_h^{(P_A+1, P_A+1)}$  given by (2.12) with the chosen mapping  $\mathcal{P}_h$ . By Theorem 1 this holds iff the scheme possesses the order of accuracy  $P_A + 1$ .  $\square$

#### 10.4. The quasi-one-dimensional case

820 **Algorithm 3** (to detect the formal and the long-time simulation orders in the quasi-1D case).

1. Detect the optimal value of the order of accuracy  $P$  using Algorithm 2.
2. Choose any local mapping  $\mathcal{P}_h$  such that  $(\mathcal{P}_1 1)_{0,\xi} \neq 0$  for each  $\xi \in M^0$ .
3. Put  $Q' = P$ .
4. While the set of equations  $(\epsilon_1(\mathbf{r}^{\mathbf{m}}/\mathbf{m}!, \Pi_1^{(P,Q'+1)}))_0 = 0$ , where  $0 \leq |\mathbf{m}| \leq Q' + 1$ , is consistent as a system for diagonal matrices  $\{\mathfrak{C}^{(\mathbf{m})}, P \leq |\mathbf{m}| \leq Q' + 1\}$ , increment  $Q'$ .

**Theorem 10.6.** *Let the scheme (2.8) be stable and quasi-1D. Let  $P$  and  $Q$  be its optimal values of the formal order of accuracy and the long-time simulation order (see Definition 5). If  $Q = \infty$ , then the algorithm loops endlessly. Otherwise the values  $P$  and  $Q'$  given by Algorithm 3 coincide with  $P$  and  $Q$ .*

*Proof.* First assume  $Q < \infty$ . By Theorem 3 there exist  $\mathfrak{C}^{(\mathbf{m})}$ ,  $P \leq |\mathbf{m}| \leq Q$ , such that the scheme possesses the truncation error of order  $Q$  in the sense of  $\Pi_h^{(P,Q)}$ . Then  $\hat{\epsilon}(\phi, \Pi_h^{(P,Q)}) = O(|\phi|^{Q+1})$  and by Lemma 10.1 the matrices  $\mathfrak{C}^{(\mathbf{m})}$  satisfy the system checked at step 4 for  $Q' = Q - 1$ . Thus the algorithm either returns  $Q' \geq Q$  or loops endlessly.

If  $Q = \infty$ , then for each  $q \in \mathbb{N}$  the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $q$ . By the argument above, the algorithm will find coefficients  $\mathfrak{C}^{(\mathbf{m})}$  for each  $P \leq |\mathbf{m}| \leq q$ . Since  $q$  may be chosen arbitrary high, the algorithm will loop endlessly.

835 Now suppose that the algorithm returns values  $P$  and  $Q' < \infty$ . Let  $\{\mathfrak{C}^{(\mathbf{m})}, P \leq |\mathbf{m}| \leq Q'\}$  be a solution of the last consistent system and let  $\Pi_h^{(P,Q')}$  be given by (2.12) with these coefficients (if no systems checked at step 4 were consistent, put  $\Pi_h^{(P,Q')} = \Pi_h$ ). By Lemma 10.1 there holds  $\hat{\epsilon}(\phi, \Pi_h^{(P,Q')}) = A(\phi)v(\phi, \Pi_h^{(P,Q')}) = O(|\phi|^{Q'+1})$ . Obviously,  $v(\phi, \Pi_h^{(P,Q')}) - v(\phi, \Pi_h) = O(|\phi|^P)$ . Thus by Lemma 8.1 the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q'$  in the sense of  $\Pi_h$ , i. e. there holds  $Q \geq Q'$ .

840 If the algorithm loops endlessly, then arguing as above we obtain that the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $s$  for each  $s \in \mathbb{N}$ ,  $s \geq P$ . By Lemma 9.19 it possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q = \infty$ .  $\square$

The following proposition simplifies the analysis in practice allowing to reduce the number of unknowns while processing Algorithm 3. Note that it assumes the additional condition  $\langle \mu_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ .

845 **Proposition 10.7.** *Let  $\Pi_h$  and  $\mathcal{P}_h$  be local mappings with kernels  $\mu_\xi, \hat{\mu}_\xi$  such that  $\langle \mu_\xi, 1 \rangle \neq 0$  and  $\langle \hat{\mu}_\xi, 1 \rangle \neq 0$  for each  $\xi \in M^0$ . Let  $p, q \in \mathbb{N}$ . Suppose there exist real-valued diagonal matrices  $\mathfrak{C}^{(\mathbf{m})}$ ,  $p \leq |\mathbf{m}| \leq q$ , such that the scheme possesses the truncation error of order  $q$  in the sense of  $\Pi_h^{(p,q)}$  given by (2.12). For each  $\mathbf{m}$ ,  $p \leq |\mathbf{m}| \leq q$ , take any  $\xi_{\mathbf{m}} \in M^0$  and  $c_{\mathbf{m}} \in \mathbb{R}$ . Then there exist real-valued diagonal matrices  $\tilde{\mathfrak{C}}^{(\mathbf{m})}$ ,  $p \leq |\mathbf{m}| \leq q$ , such that  $\tilde{\mathfrak{C}}_{\xi_{\mathbf{m}}, \xi_{\mathbf{m}}}^{(\mathbf{m})} = c_{\mathbf{m}}$  and the scheme possesses the truncation error of order  $q$  in the sense of  $\tilde{\Pi}_h^{(p,q)}$  given by (2.12) with the substitution  $\tilde{\mathfrak{C}}^{(\mathbf{m})}$  for  $\mathfrak{C}^{(\mathbf{m})}$ .*

*Proof.* Denote  $f_{p-1}(\phi) = v(\phi, \Pi_h^{(p,q)})$ . The function  $f_{p-1}(\phi)$  is holomorphic by construction; by assumption  $(f_{p-1}(0))_\xi \neq 0$  for each  $\xi \in M^0$ . Denote

$$Y_{\mathbf{m}}[f] = \frac{1}{\mathbf{m}!} \frac{1}{i^{|\mathbf{m}|}} D^{\mathbf{m}} \left( \frac{(f(\phi))_{\xi_{\mathbf{m}}}}{(v(\phi, \mathcal{P}_h))_{\xi_{\mathbf{m}}}} \right) \Big|_{\phi=0}, \quad g(\phi) = v(\phi, \Pi_h).$$

For  $k = p, \dots, q$  define inductively

$$f_k(\phi) = f_{k-1}(\phi) \left( 1 + \sum_{|\mathbf{m}|=k} a_{\mathbf{m}} (i\phi)^{\mathbf{m}} \right),$$

choosing constants  $a_{\mathbf{m}}$ ,  $|\mathbf{m}| = k$ , to satisfy

$$Y_{\mathbf{m}}[f_k - g] = c_{\mathbf{m}}. \tag{10.2}$$

Each of these equations easily reduces to the linear equation

$$\frac{(v(0, \Pi_h))_{\xi_m}}{(v(0, \mathcal{P}_h))_{\xi_m}} a_m + Y_m[f_{k-1} - g] = c_m,$$

which is uniquely solvable since by assumption  $(v(0, \Pi_h))_{\xi_m} \neq 0$ . Moreover,  $f_k$  satisfies (10.2) for each  $|\mathbf{m}| < k$ . Indeed, for  $k = p$  this is by construction. Let it be true for  $f_{k-1}$ . If  $|\mathbf{m}| \leq k - 1$  then

$$Y_m[f_k - g] = Y_m[f_{k-1} - g] + Y_m \left[ f_{k-1}(\phi) \sum_{|\mathbf{n}|=k} a_n (i\phi)^{\mathbf{n}} \right] = c_m + \sum_{|\mathbf{n}|=k} a_n Y_m[f_{k-1}(\phi)(i\phi)^{\mathbf{n}}].$$

It remains to note that  $Y_m[f(\phi)(i\phi)^{\mathbf{n}}] = 0$  if  $|\mathbf{n}| > |\mathbf{m}|$ .

Denote  $W(\phi) = f_q(\phi)$ . By construction, there holds  $W(\phi) = v(\phi, \Pi_h)(1 + O(|\phi|^p))$  and  $W(\phi) = v(\phi, \Pi_h^{(p,q)})(1 + O(|\phi|^p))$ . By Lemma 6.7 there exist real-valued diagonal matrices  $\tilde{\mathfrak{C}}^{(\mathbf{m})}$ ,  $p \leq |\mathbf{m}| \leq q$ , such that  $W(\phi) = v(\phi, \tilde{\Pi}_h^{(p,q)}) + O(|\phi|^{q+1})$ ; by construction from (10.2) (see (6.13)) we have  $\tilde{\mathfrak{C}}_{\xi_m, \xi_m}^{(\mathbf{m})} = c_m$ . Then

$$\begin{aligned} \hat{\epsilon}(\phi, \tilde{\Pi}_h^{(p,q)}) &= A(\phi)v(\phi, \tilde{\Pi}_h^{(p,q)}) = A(\phi)W(\phi) + O(|\phi|^{q+1}) = \\ &= A(\phi)v(\phi, \Pi_h^{(p,q)})(1 + O(|\phi|^p)) + O(|\phi|^{q+1}) = O(|\phi|^{q+1}). \end{aligned}$$

By Theorem 6.15 the scheme possesses the truncation error of order  $q$  in the sense of  $\tilde{\Pi}_h^{(p,q)}$ . □

### 10.5. The simple case

For a simple scheme the optimal values of the formal order of accuracy and the long-time simulation order can be found by the following algorithm.

**Algorithm 4** (to detect the formal and long-time simulation orders in the simple case).

1. Detect the order of accuracy  $P_A$  using Algorithm 1.
2. Choose any local mapping  $\mathcal{P}_h$  such that  $(\mathcal{P}_1 1)_{0, \xi} \neq 0$  for each  $\xi \in M^0$ .
3. Put  $m = P_A + 1$ .
- 860 4. Compute  $f^{\mathbf{m}} = -(\epsilon_1(\mathbf{r}^{\mathbf{m}}/\mathbf{m}!, \Pi_1^{(P_A+1, m-1)}))_0$  for each  $\mathbf{m}$  such that  $|\mathbf{m}| = m$ , substituting into  $\Pi_h^{(P_A+1, m-1)}$  the previously found coefficients  $\mathfrak{C}^{(\mathbf{n})}$ ,  $P_A + 1 \leq |\mathbf{n}| \leq m - 1$ .
5. If for each  $\mathbf{m}$  with  $|\mathbf{m}| = m$  there holds  $f^{\mathbf{m}} \in \text{Im} L(0)$ , then:
  - find any diagonal matrices  $\mathfrak{C}^{(\mathbf{m})}$  satisfying  $L(0)\mathfrak{C}^{(\mathbf{m})}(\mathcal{P}_1 1)_0 = f^{\mathbf{m}}$ ;
  - increase  $m$  by one;
  - 865 • return to the step 4.
6. Put  $Q' = m - 1$ .
7. Put  $P' = P_A$  if  $Q' = P_A$  and  $P' = P_A + 1$  otherwise.

Note that at step 5, since  $\text{Rank } L(0) = |M^0| - 1$ , the equation specifies a one-parametric family of diagonal matrices  $\mathfrak{C}^{\mathbf{m}}$ . Any matrix of this family may be chosen.

**Theorem 10.8.** Consider a simple scheme of the form (2.8) and a local mapping  $\Pi_h$ . Let  $P$  and  $Q$  be its optimal values of the formal order of accuracy and the long-time simulation order (see Definition 5). If  $Q = \infty$ , then the algorithm loops endlessly. Otherwise Algorithm 4 returns values  $P' = P$  and  $Q' = Q$ .

*Proof.* Let  $P_A$  be the optimal order of the truncation error. Let  $P'$  and  $Q'$  be the values returned by the algorithm; if it loops endlessly, formally put  $P' = P_A + 1$ ,  $Q' = \infty$ . We need to prove that  $Q' = Q$ , then the equality  $P' = P$  will follow from statement 5 of Theorem 3.

First we prove that  $Q \geq Q'$ . If  $Q' = P_A$ , this is obvious by stability. If  $P_A < Q' < \infty$ , arguing as in the proof of Theorem 10.6 we get that the scheme possesses the orders  $P'$  and  $Q'$ . By Lemma 9.6 it possesses the orders  $P$  and

880  $Q'$ , hence  $Q \geq Q'$ . If  $Q' = \infty$ , using the same argument, for each  $s \geq P'$  the scheme possesses the formal order of accuracy  $P'$  and the long-time simulation order  $s$ . By Lemma 9.8 it possesses the formal order of accuracy  $P'$  and the long-time simulation order  $Q = \infty$ .

Now we prove that  $Q' \geq Q$ . Let  $\Pi_h^{(P_A+1, Q')}$  be given by (2.12) with the coefficients  $\mathfrak{C}^{(m)}$  found by the algorithm. By construction and Lemma 10.4,

$$\left( \epsilon_1 \left( \frac{\mathbf{r}^m}{m!}, \Pi_1^{(P_A+1, Q')} \right) \right)_0 = 0$$

885 for each  $|\mathbf{m}| \leq Q'$ , i. e. the scheme possesses the truncation error of order  $Q'$  in the sense of  $\Pi_h^{(P_A+1, Q')}$ . Assume that  $Q > Q'$ . If  $Q'$  is an optimal value of the truncation error in the sense of  $\Pi_h^{(P_A+1, Q')}$ , then by statement 5 of Theorem 3 the scheme possesses the formal order of accuracy  $Q' + 1$  in the sense of  $\Pi_h^{(P_A+1, Q')}$ . By statement 3 of Theorem 1 there exist  $\mathfrak{C}^{(m)}$ ,  $|\mathbf{m}| = Q' + 1$ , such that the scheme possesses the truncation error of order  $Q' + 1$  in the sense of  $\Pi_h^{(P_A+1, Q'+1)}$ . If  $Q'$  is not an optimal value, then the scheme possesses the truncation error of order  $Q' + 1$  in the sense of  $\Pi_h^{(P_A+1, Q'+1)}$  with  $\mathfrak{C}^{(m)} = 0$ ,  $|\mathbf{m}| = Q' + 1$ . These coefficients  $\mathfrak{C}^{(m)}$ ,  $|\mathbf{m}| = Q' + 1$ , give a solution of the equation to solve at step 5, thus the algorithm makes one more step. This contradiction proves that  $Q' \geq Q$  and thus the whole Theorem.  $\square$

Note that in the quasi-1D case Algorithm 4 may give a wrong result (see Section 11.3).

## 890 11. 1D examples

Throughout this section we consider 1D transport equation (2.6) with transport velocity  $\omega = 1$ , i. e.  $\partial v / \partial t + \partial v / \partial x = 0$ , and put  $a_1 = 1$ , i. e. the parameter  $h$  coincides with the mesh step.

### 11.1. P1 discontinuous Galerkin method

Consider the discontinuous Galerkin method (its definition may be found, for example, in [31]) based on the piecewise-linear polynomials on the uniform mesh with the nodes  $x_j = jh$ . On each cell  $(x_j, x_{j+1})$  define two basis functions:

$$\phi_j^L(x) = \frac{x_{j+1} - x}{h}, \quad \phi_j^R(x) = \frac{x - x_j}{h},$$

and extend these functions by zero outside  $(x_j, x_{j+1})$ . For the numerical solution

$$u(t, x) = \sum_{j=0}^{N-1} \left( u_j^L(t) \phi_j^L(x) + u_j^R(t) \phi_j^R(x) \right),$$

the DG scheme gives

$$\left( \begin{array}{cc} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{array} \right) \frac{du_j}{dt} + \frac{1}{h} \left[ \left( \begin{array}{cc} \frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{array} \right) u_j + \left( \begin{array}{cc} 0 & -1 \\ 0 & 0 \end{array} \right) u_{j-1} \right] = 0. \quad (11.1)$$

It has the form (2.8) with  $\mathcal{S} = \{0, -1\}$ ,

$$Z_{-1} = 0, \quad Z_0 = \left( \begin{array}{cc} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{array} \right), \quad L_{-1} = \left( \begin{array}{cc} 0 & -1 \\ 0 & 0 \end{array} \right), \quad L_0 = \left( \begin{array}{cc} \frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{array} \right).$$

By definition (6.14) we have

$$A(\phi) = i\phi I - \left( \begin{array}{cc} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{array} \right)^{-1} \left( \begin{array}{cc} \frac{1}{2} & \frac{1}{2} - e^{-i\phi} \\ -\frac{1}{2} & \frac{1}{2} \end{array} \right) = \left( \begin{array}{cc} i\phi - 3 & -1 + 4e^{-i\phi} \\ 3 & i\phi - 1 - 2e^{-i\phi} \end{array} \right).$$

We see that  $A(0) \neq 0$  and

$$\det A(\phi) = (-2i\phi - 6)e^{-i\phi} + (-\phi^2 - 4i\phi + 6) = \frac{\phi^4}{12} + O(\phi^5).$$

The eigenvalues of the matrix  $A(\phi)$  are  $\lambda_{\pm}(\phi) = i\phi - 2 - e^{-i\phi} \pm \varkappa(\phi)$ , where  $\varkappa(\phi) = (e^{-2i\phi} + 10e^{-i\phi} - 2)^{1/2}$ . The branch of the root is such that  $\varkappa(0) = 3$ . The right eigenvectors matrix is

$$S(\phi) = \frac{1}{4e^{-i\phi} - 1} \begin{pmatrix} 4e^{-i\phi} - 1 & -4e^{-i\phi} + 1 \\ 1 - e^{-i\phi} + \varkappa(\phi) & -1 + e^{-i\phi} + \varkappa(\phi) \end{pmatrix} = \begin{pmatrix} \mathbf{r}_+(\phi) & \mathbf{r}_-(\phi) \end{pmatrix},$$

the left eigenvectors matrix is

$$S^{-1}(\phi) = \frac{1}{2\varkappa(\phi)} \begin{pmatrix} -1 + e^{-i\phi} + \varkappa(\phi) & 4e^{-i\phi} - 1 \\ -1 + e^{-i\phi} - \varkappa(\phi) & 4e^{-i\phi} - 1 \end{pmatrix} = \begin{pmatrix} \mathbf{l}_+(\phi) \\ \mathbf{l}_-(\phi) \end{pmatrix}.$$

The first column in  $S(\phi)$  and, correspondingly, the first row in  $S^{-1}(\phi)$  correspond to  $\lambda_+(\phi)$ , the second ones to  $\lambda_-(\phi)$ . The eigenvalues can be approximated as  $\lambda_+(\phi) \approx -\phi^4/72 + O(\phi^5)$ ,  $\lambda_-(\phi) = -6 + O(\phi)$ . Clearly,  $\lambda_+(\phi)$  corresponds to the physical mode

$$\begin{pmatrix} u_j^L(t) \\ u_j^R(t) \end{pmatrix} = \exp\left(i\frac{\phi}{h}(jh - t)\right) \exp\left(\lambda_+(\phi)\frac{t}{h}\right) \mathbf{r}_+(\phi)$$

and  $\lambda_-(\phi)$  to the spurious one (the same expression with  $\lambda_+(\phi)$ ,  $\mathbf{r}_+(\phi)$  replaced by  $\lambda_-(\phi)$ ,  $\mathbf{r}_-(\phi)$ ).

In order to speak about the accuracy we need to specify a map  $\Pi_h$ . For simplicity we use the pointwise map  $(\Pi_h f)_j^L = f(x_j)$ ,  $(\Pi_h f)_j^R = f(x_{j+1})$ . Then  $v(\phi, \Pi_h) = (1, e^{i\phi})^T$  and

$$S^{-1}(\phi)v(\phi, \Pi_h) = \frac{1}{2\varkappa(\phi)} \begin{pmatrix} 3 + e^{-i\phi} - e^{i\phi} + \varkappa(\phi) \\ 3 + e^{-i\phi} - e^{i\phi} - \varkappa(\phi) \end{pmatrix} = \begin{pmatrix} 1 + O(\phi) \\ \frac{1}{12}\phi^2 + O(\phi^3) \end{pmatrix}.$$

895 Using the notation used in Section 7 we have  $\bar{\mathbf{N}} = \{0, 4\}$ ,  $p_0 = 2$ ,  $p_4 = 0$ . Thus by Lemma 7.6 the scheme possesses the 1st order of the truncation error and the optimal error estimate has the form  $O(h^2 + h^3t)$ .

This result can be also obtained with the help of Proposition 9.14 without the calculation of eigenvalues. Since  $\dim \text{Ker } A(0) = 1$ , the scheme is simple. By construction, the scheme is exact on linear functions and thus possesses the truncation error of the first order. Using the fact that  $\det A(\phi) = c\phi^4 + O(|\phi|^5)$ ,  $c \neq 0$ , by Proposition 9.14  
900 the scheme possesses the second order of accuracy and the third order in the long-time simulation. If we additionally verify that the scheme is not exact on quadratic polynomials, then by Proposition 9.14 we get that these values are optimal.

Replacing  $\Pi_h$  by  $\tilde{\Pi}_h$  it is possible to improve  $p_0$  and hence the order of the truncation error and the formal order of accuracy. For this purpose we need to reduce the second component of  $S^{-1}(\phi)v(\phi, \tilde{\Pi}_h)$ . We can nullify it if we put  $w(\phi) = \gamma(\phi)\mathbf{r}_+(\phi)$  for some  $\gamma(\phi) \in \mathbb{C} \setminus \{0\}$  such that  $\gamma(-\phi) = \overline{\gamma(\phi)}$  (for example,  $\gamma(\phi) \equiv 1$ ) and get  $\tilde{\Pi}_h$  by Lemma 6.6 such that  $v(\phi, \tilde{\Pi}_h) = w(\phi)$  in a neighborhood of  $\phi = 0$ . This map is nonlocal; to specify a local map we need to approximate  $w(\phi)$ . For example, if we put  $\gamma(\phi) \equiv 1$ , define  $\mathcal{P}_h$  as  $(\mathcal{P}_h f)_{\eta,L} = (\mathcal{P}_h f)_{\eta,R} = f(h\eta)$ , and use Lemma 6.7 then we get the mapping  $\Pi_h^{(2,3)}$  given by

$$(\Pi_h^{(2,3)} f)_{\eta,L} = f(\eta h), \quad (\Pi_h^{(2,3)} f)_{\eta,R} = f((\eta + 1)h) + \frac{1}{6}h^2 f''(\eta h) + \frac{5}{18}h^3 f'''(\eta h). \quad (11.2)$$

In this case  $v(\phi, \Pi_h^{(2,3)}) = (1, e^{i\phi} - \phi^2/6 - 5i\phi^3/18)^T$ . The scheme (11.1) possesses the truncation error of order 3 in the sense of  $\Pi_h^{(2,3)}$ . Doing the same with  $\mathcal{P}_h = \Pi_h$  we get the mapping  $\tilde{\Pi}_h^{(2,3)}$  given by

$$(\tilde{\Pi}_h^{(2,3)} f)_{\eta,L} = f(\eta h), \quad (\tilde{\Pi}_h^{(2,3)} f)_{\eta,R} = f((\eta + 1)h) + \frac{1}{6}h^2 f''((\eta + 1)h) + \frac{1}{9}h^3 f'''((\eta + 1)h),$$

and thus  $v(\phi, \tilde{\Pi}_h^{(2,3)}) = (1, e^{i\phi}(1 - \phi^2/6 - i\phi^3/9))^T$ . Approximations of derivatives in (11.2) with enough order yield other suitable mappings.

905 Now we demonstrate the method of auxiliary mapping with no use of the eigenvectors. We have

$$L(0) = L_0 + L_{-1} = \begin{pmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{pmatrix}.$$

Since  $\dim \text{Ker } L(0) = 1$ , we will use Algorithm 4. We will use  $\mathcal{P}_h$  defined as  $(\mathcal{P}_h f)_{\eta,L} = (\mathcal{P}_h f)_{\eta,R} = f(h\eta)$ , then  $\vec{\epsilon} = v(0, \mathcal{P}_h) = (1, 1)^T$ . Substituting a linear function into the scheme we see that the scheme possesses the first order of the truncation error. Doing the same with a quadratic polynomial we see the optimal value of the order of the truncation error is equal to 1.

We need to compute  $f^2 = -(\epsilon_1(x^2/2, \Pi_1))_0$ . We have

$$f^2 = (Z\Pi_1 x - L\Pi_1(x^2/2))_0 = Z_0(\Pi_1 x)_0 - L_0(\Pi_1(x^2/2))_0 - L_{-1}(\Pi_1(x^2/2))_{-1}$$

and

$$f^2 = \begin{pmatrix} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{pmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} - \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{pmatrix} \begin{bmatrix} 0 \\ 1/2 \end{bmatrix} - \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix} \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} = \begin{bmatrix} -1/12 \\ 1/12 \end{bmatrix}.$$

The system

$$\begin{pmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{pmatrix} \mathfrak{C}^{(2)} \vec{\epsilon} = \begin{pmatrix} -1/12 \\ 1/12 \end{pmatrix}$$

has the solution

$$\mathfrak{C}^{(2)} = \begin{pmatrix} 0 & 0 \\ 0 & 1/6 \end{pmatrix} + c_2 I.$$

910 According to the algorithm the value of  $c_2$  can be set arbitrarily; put  $c_2 = 0$ . Now we compute  $f^3 = -(\epsilon_1(x^3/6, \Pi_1^{(2,2)}))_0$  where  $\Pi_h^{(2,2)}$  is prescribed by (2.12) with the coefficients  $\mathfrak{C}^{(2)}$  we just found. We have

$$f^3 = \begin{pmatrix} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{pmatrix} \left[ \begin{pmatrix} 0 \\ 1/2 \end{pmatrix} + \mathfrak{C}^{(2)} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right] - \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{pmatrix} \left[ \begin{pmatrix} 0 \\ 1/6 \end{pmatrix} + \mathfrak{C}^{(2)} \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right] - \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix} \left[ \begin{pmatrix} -1/6 \\ 0 \end{pmatrix} + \mathfrak{C}^{(2)} \begin{pmatrix} -1 \\ -1 \end{pmatrix} \right] = \begin{pmatrix} -5/36 \\ 5/36 \end{pmatrix}.$$

The system

$$\begin{pmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{pmatrix} \mathfrak{C}^{(3)} \vec{\epsilon} = \begin{pmatrix} -5/36 \\ 5/36 \end{pmatrix}$$

has the solution

$$\mathfrak{C}^{(3)} = \begin{pmatrix} 0 & 0 \\ 0 & 5/18 \end{pmatrix} + c_3 I.$$

The map  $\Pi_h^{(2,3)}$  generated by chosen  $\Pi_h$ ,  $\mathcal{P}_h$  and matrices  $\mathfrak{C}^{(2)}$  and  $\mathfrak{C}^{(3)}$  with  $c_2 = c_3 = 0$  coincides with (11.2).

Repeat the operation for  $f^4 = -(\epsilon_1(x^4/24, \Pi_1^{(2,3)}))_0$ .

$$f^4 = \begin{pmatrix} \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{3} \end{pmatrix} \left[ \begin{pmatrix} 0 \\ 1/6 \end{pmatrix} + \mathfrak{C}^{(2)} \begin{pmatrix} 1 \\ 1/2 \end{pmatrix} + \mathfrak{C}^{(3)} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right] - \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{pmatrix} \left[ \begin{pmatrix} 0 \\ 1/24 \end{pmatrix} + \mathfrak{C}^{(2)} \begin{pmatrix} 0 \\ 0 \end{pmatrix} + \mathfrak{C}^{(3)} \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right] - \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix} \left[ \begin{pmatrix} -1/24 \\ 0 \end{pmatrix} + \mathfrak{C}^{(2)} \begin{pmatrix} -1/2 \\ -1/2 \end{pmatrix} + \mathfrak{C}^{(3)} \begin{pmatrix} -1 \\ -1 \end{pmatrix} \right] = \begin{pmatrix} -127/432 \\ 67/432 \end{pmatrix}.$$

The system  $L(0)\mathfrak{C}^{(4)}\vec{\epsilon} = f^4$  is inconsistent. Thus  $P = 2$  and  $Q = 3$  are the optimal values of the formal order of accuracy and of the long-time simulation order.

## 915 11.2. Arbitrary order DG method

In this subsection we demonstrate Algorithm 4 on the discontinuous Galerkin method based on the  $p$ -th order polynomials, where  $p \in \mathbb{N}$ . The discontinuous Galerkin method gives a scheme of the form

$$\sum_{\xi=0}^p \frac{du_{\eta,\xi}}{dt} \int_0^1 \phi_\xi(x) \psi_j(x) dx - \sum_{\xi=0}^p u_{\eta,\xi} \int_0^1 \phi_\xi(x) \psi_j'(x) dx + \sum_{\xi=0}^p u_{\eta,\xi} \phi_\xi(1) \psi_j(1) - \sum_{\xi=0}^p u_{\eta-1,\xi} \phi_\xi(1) \psi_j(0) = 0,$$

where  $\phi_\xi, \xi = 0, \dots, p$  and  $\psi_j, j = 0, \dots, p$  are two bases in the space of polynomials of order  $p$ .

In this scheme we can choose the set of basis functions in the polynomial space and the initial local mapping  $\Pi_h$  at our convenience. Let  $x_k$  and  $w_k$  be the nodes and weights of the Gauss – Jacobi quadrature rule with  $p + 1$  nodes, i. e. the quadrature rule on  $[0, 1]$  exact on the polynomials of order  $2p$  such that  $x_p = 1$ . Let  $\{\phi_\xi(x), x \in [0, 1], \xi = 0, \dots, p\}$  be the collocation basis with the nodes  $x_\xi$ , i. e.  $\phi_\xi(x_\zeta) = 1$  for  $\xi = \zeta$  and 0 otherwise. Let the test functions be  $\psi_j(x) = x^j/j!, j = 0, \dots, p$ . Then we have

$$\int_0^1 \phi_\xi(x) \psi_j(x) dx = \sum_{k=0}^p w_k \phi_\xi(x_k) \psi_j(x_k) = w_\xi \frac{x_\xi^j}{j!}, \quad \int_0^1 \phi_\xi(x) \psi'_j(x) dx = \sum_{k=0}^p w_k \phi_\xi(x_k) \psi'_j(x_k) = w_\xi \frac{x_\xi^{j-1}}{(j-1)!},$$

and  $\phi_\xi(1) = 1$  for  $\xi = p$  and 0 otherwise. So the scheme is of the form (2.8), namely,

$$\begin{aligned} Z_0 \frac{du_\eta}{dt} + L_0 u_\eta + L_{-1} u_{\eta-1} &= 0, \\ Z_0 &= \begin{pmatrix} w_0 & \dots & w_{p-1} & w_p \\ w_0 x_0 & \dots & w_{p-1} x_{p-1} & w_p x_p \\ w_0 x_0^2/2 & \dots & w_{p-1} x_{p-1}^2/2 & w_p x_p^2/2 \\ \vdots & \ddots & \vdots & \vdots \\ w_0 x_0^p/p! & \dots & w_{p-1} x_{p-1}^p/p! & w_p x_p^p/p! \end{pmatrix}, \quad L_{-1} = \begin{pmatrix} 0 & \dots & 0 & -1 \\ 0 & \dots & 0 & 0 \\ 0 & \dots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & 0 \end{pmatrix}, \\ L_0 &= \begin{pmatrix} 0 & \dots & 0 & 1 \\ -w_0 & \dots & -w_{p-1} & 1 - w_p \\ -w_0 x_0 & \dots & -w_{p-1} x_{p-1} & 1/2 - w_p x_p \\ \vdots & \ddots & \vdots & \vdots \\ -w_0 x_0^{p-1}/(p-1)! & \dots & -w_{p-1} x_{p-1}^{p-1}/(p-1)! & 1/p! - w_p x_p^{p-1}/(p-1)! \end{pmatrix}. \end{aligned} \quad (11.3)$$

We equip the scheme with the mapping  $\Pi_h$  defined by  $(\Pi_h f)_{\eta, \xi} = f(h(\eta + x_\xi))$ .

The bottom-left  $(p \times p)$ -submatrix of  $L(0) = L_{-1} + L_0$  is the Vandermonde matrix with scaled columns, so it is nonsingular. Thus  $\dim \text{Ker } L(0) = 1$  and  $\text{Im } L(0)$  consists of the vectors with the first component equal to zero. Denote  $\vec{\epsilon} = (1, \dots, 1)^T$ . We claim that  $L(0)\vec{\epsilon} = 0$  and thus  $\text{Ker } L(0) = \text{span}\{\vec{\epsilon}\}$ . Indeed, for  $j = 0$  the equality  $(L(0)\vec{\epsilon})_j = 0$  is obvious; for  $j > 0$  we have

$$(L(0)\vec{\epsilon})_j = \frac{1}{j!} - \sum_{\xi=0}^p w_\xi \frac{x_\xi^{j-1}}{(j-1)!} = 0.$$

920 The scheme is simple ( $\dim \text{Ker } L(0) = 1$ ) and we can use Algorithm 4. We choose the mapping  $\mathcal{P}_h$  defined by  $(\mathcal{P}_h f)_{\eta, \xi} = f(h\eta)$ ,  $\xi = 0, \dots, p$ . Step 5 of the algorithm defines the diagonal matrix  $\mathfrak{C}^{(m)}$  as a solution of  $L(0)\mathfrak{C}^{(m)}\vec{\epsilon} = f^m$ . Since  $L(0)\vec{\epsilon} = 0$ , the general solution of this equation has the form  $\mathfrak{C}^{(m)} + \alpha I$ ,  $\alpha \in \mathbb{R}$ , so we will specify  $\mathfrak{C}^{(m)}$  by  $\mathfrak{C}_{p,p}^{(n)} = 0$ .

Step 1. We need to get the order of the truncation error. We have

$$-\left(\epsilon_1 \left(\frac{x^m}{m!}, \Pi_1\right)\right)_0 = Z_0 \left(\Pi_1 \frac{x^{m-1}}{(m-1)!}\right)_0 - L_0 \left(\Pi_1 \frac{x^m}{m!}\right)_0 - L_{-1} \left(\Pi_1 \frac{x^m}{m!}\right)_{-1}.$$

Componentwise,

$$-\left(\epsilon_1 \left(\frac{x^m}{m!}, \Pi_1\right)\right)_{0,j} = \sum_{\xi} (Z_0)_{j,\xi} \left(\Pi_1 \frac{x^{m-1}}{(m-1)!}\right)_{0,\xi} - \sum_{\xi} (L_0)_{j,\xi} \left(\Pi_1 \frac{x^m}{m!}\right)_{0,\xi} - \sum_{\xi} (L_{-1})_{j,\xi} \left(\Pi_1 \frac{x^m}{m!}\right)_{-1,\xi}.$$

For  $m = 0$  this is zero, so assume  $m > 0$ . For  $j = 0$  we have

$$-\left(\epsilon_1 \left(\frac{x^m}{m!}, \Pi_1\right)\right)_{0,0} = \sum_{\xi} w_\xi \frac{x_\xi^{m-1}}{(m-1)!} - \frac{1}{m!}$$

and for  $j = 1, \dots, p$

$$-\left(\epsilon_1\left(\frac{x^m}{m!}, \Pi_1\right)\right)_{0,j} = \sum_{\xi} w_{\xi} \frac{x_{\xi}^j}{j!} \frac{x_{\xi}^{m-1}}{(m-1)!} + \sum_{\xi} w_{\xi} \frac{x_{\xi}^{j-1}}{(j-1)!} \frac{x_{\xi}^m}{m!} - \frac{1}{j!m!} = \frac{m+j}{j!m!} \left( \sum_{\xi} w_{\xi} x_{\xi}^{j+m-1} - \frac{1}{m+j} \right).$$

Since  $x_{\xi}$  and  $w_{\xi}$  are nodes and weights of the quadrature rule of order  $2p$ , then

$$\left(\epsilon_1\left(\frac{x^m}{m!}, \Pi_1\right)\right)_{0,j} = 0, \quad j+m-1 \leq 2p. \quad (11.4)$$

Since the maximal value for  $j$  is  $p$ , then the truncation error is zero for  $m \leq p+1$ . Thus the scheme possesses the truncation error of order  $P_A = p+1$ . The value  $P_A = p+1$  is optimal, otherwise the quadrature formula would be exact on all polynomials of order  $2p+1$ .

Step 2. Put  $m = P_A + 1 = p+2$ .

Steps 3, 4. Let  $\mathfrak{C}^{(n)}$ ,  $n = p+2, \dots, m-1$  be the coefficients already found by the algorithm by solving

$$L(0)\mathfrak{C}^{(n)}\vec{\mathfrak{e}} = f^n, \quad f^n = -\left(\epsilon_1\left(\frac{x^n}{n!}, \Pi_1^{(p+2, n-1)}\right)\right)_0 \quad (11.5)$$

for  $\mathfrak{C}^{(n)}$ . We claim that the first  $2p+2-n$  components of the right-hand side of (11.5) are equal to zero. Indeed, for  $n = p+2$  this follows from (11.4). Assume this holds for all  $n' = p+2, \dots, n-1$ . Then

$$\begin{aligned} f^n = & -\left(\epsilon_1\left(\frac{x^n}{n!}, \Pi_1^{(p+2, n-1)}\right)\right)_0 = -\left(\epsilon_1\left(\frac{x^n}{n!}, \Pi_1\right)\right)_0 + Z_0 \sum_{n'=p+2}^{n-1} \mathfrak{C}^{(n')} \left(\mathcal{P}_1 \frac{x^{n-1-n'}}{(n-1-n')!}\right)_0 - \\ & -L_0 \sum_{n'=p+2}^{n-1} \mathfrak{C}^{(n')} \left(\mathcal{P}_1 \frac{x^{n-n'}}{(n-n')!}\right)_0 - L_{-1} \sum_{n'=p+2}^{n-1} \mathfrak{C}^{(n')} \left(\mathcal{P}_1 \frac{x^{n-n'}}{(n-n')!}\right)_{-1}. \end{aligned}$$

Recall that  $\mathcal{P}_1$  takes the point value at  $x = 0$  and  $\mathfrak{C}_{p,p}^{(n)} = 0$  for each  $n$ . Thus

$$f^n = -\left(\epsilon_1\left(\frac{x^n}{n!}, \Pi_1\right)\right)_0 + Z_0 \mathfrak{C}^{(n-1)}\vec{\mathfrak{e}}.$$

For the first term on the right-hand side the first  $2p+2-n$  components are zero (see (11.4)). For the second term on the right-hand side by (11.3) and  $\mathfrak{C}_{p,p}^{(n')} = 0$  for each  $n'$  we have

$$(Z_0 \mathfrak{C}^{(n-1)}\vec{\mathfrak{e}})_j = (-L_0 \mathfrak{C}^{(n-1)}\vec{\mathfrak{e}})_{j+1}, \quad j = 0, \dots, p-1.$$

So by induction assumption the first  $2p+2-n$  components of  $Z_0 \mathfrak{C}^{(n-1)}\vec{\mathfrak{e}}$  are also zero. Thus the first  $2p+2-n$  components of  $f^n$  are equal to zero. In particular, for each  $m \leq 2p+1$  there holds  $f_0^m = 0$ , i. e.  $f^m \in \text{Im}L(0)$ . For  $m = 2p+2$  this possibly will not hold (in fact – not possibly but definitely, but we have no proof for this).

Step 5. Put  $Q' = m-1$ . We know that at the last stage  $m$  was greater than or equal to  $2p+2$ , thus  $Q' \geq 2p+1$ .

Step 6. Since  $Q' \geq 2p+1 > P_A = p+1$ , we have  $P' = P_A + 1 = p+2$ .

Since the algorithm returns  $Q' \geq 2p+1$ , by Theorem 10.8 the DG scheme based on  $p$ -th order polynomials possesses the long-time simulation order  $2p+1$ . This proof can be considered as a variation of the proof in [15].

### 11.3. Alternating central difference scheme

In this section we demonstrate Algorithm 3.

Let the set of the DOFs be  $M^0 = \{\text{“L”}, \text{“R”}\}$ ,  $|M^0| = 2$ . We will treat the first DOF as a value in the node and the second one as a value in the cell center. Correspondingly we define the pointwise map  $\bar{\Pi}_h$  as  $(\bar{\Pi}_h f)_{j,L} = f(jh)$ ,  $(\bar{\Pi}_h f)_{j,R} = f(jh + h/2)$ ,  $j \in \mathbb{Z}$ .



Derivatives at nodes are approximated by the 2-nd order central difference and derivatives at cell centers are approximated by the 4-th order central difference. For brevity denote  $(u_j)_L = v_j$ ,  $(u_j)_R = v_{j+1/2}$ , then

$$\begin{aligned} \frac{dv_j}{dt} + \frac{v_{j+1/2} - v_{j-1/2}}{h} &= 0, \quad j \in \mathbb{Z}; \\ \frac{dv_{j+1/2}}{dt} + \frac{4}{3} \frac{v_{j+1} - v_j}{h} - \frac{1}{3} \frac{v_{j+3/2} - v_{j-1/2}}{2h} &= 0, \quad j \in \mathbb{Z}. \end{aligned} \quad (11.6)$$

940 In the sense of  $\bar{\Pi}_h$  the scheme (11.6) possesses the second order of truncation error.

The scheme can be rewritten in the block form

$$\begin{aligned} \frac{du_j}{dt} + L_{-1}u_{j-1} + L_0u_j + L_1u_{j+1} &= 0, \\ L_{-1} &= \begin{pmatrix} 0 & -1 \\ 0 & 1/6 \end{pmatrix}, \quad L_0 = \begin{pmatrix} 0 & 1 \\ -4/3 & 0 \end{pmatrix}, \quad L_1 = \begin{pmatrix} 0 & 0 \\ 4/3 & -1/6 \end{pmatrix}. \end{aligned}$$

The matrix  $L(\phi) = e^{-i\phi}L_{-1} + L_0 + e^{i\phi}L_1$  is

$$L(\phi) = \begin{pmatrix} 0 & 1 - e^{-i\phi} \\ \frac{4}{3}(-1 + e^{i\phi}) & -\frac{1}{6}(e^{i\phi} - e^{-i\phi}) \end{pmatrix} = \begin{pmatrix} 0 & 2ie^{-i\phi/2} \sin(\phi/2) \\ \frac{8}{3}ie^{i\phi/2} \sin(\phi/2) & -\frac{1}{3}i \sin(\phi) \end{pmatrix}.$$

The matrix  $A(\phi)$  defined by (6.14) is

$$A(\phi) = \begin{pmatrix} i\phi & -2ie^{-i\phi/2} \sin(\phi/2) \\ -\frac{8}{3}ie^{i\phi/2} \sin(\phi/2) & i\phi + \frac{1}{3}i \sin(\phi) \end{pmatrix}.$$

The eigenvalues of  $A(\phi)$  are equal to

$$\lambda_{\pm}(\phi) = i\phi + \frac{1}{6}i \sin(\phi) \pm \frac{1}{6}i \sin(\phi/2) \varkappa(\phi), \quad \varkappa(\phi) = \sqrt{194 + 2 \cos \phi} > 0.$$

The eigenvalues  $\lambda_{\pm}(\phi)$  are pure imaginary. For  $\phi = 0$  there holds  $A(\phi) = 0$ . For  $\phi \neq 0$  there holds  $\lambda_+(\phi) \neq \lambda_-(\phi)$ . Thus for each  $\phi$  the matrix  $A(\phi)$  has two linearly independent eigenvectors. The matrices of the right and left eigenvectors are

$$S(\phi) = \frac{1}{4\varkappa(\phi)} \begin{pmatrix} 2 \exp(-i\phi/2) & 2 \exp(-i\phi/2) \\ i\lambda_+/\sin(\phi/2) & i\lambda_-/\sin(\phi/2) \end{pmatrix}, \quad S^{-1}(\phi) = \begin{pmatrix} \exp(i\phi/2)(\varkappa(\phi) - 2 \cos(\phi/2)) & -12 \\ \exp(i\phi/2)(\varkappa(\phi) + 2 \cos(\phi/2)) & 12 \end{pmatrix}.$$

For  $\phi = 0$  we take the corresponding limits. The first column of  $S(\phi)$  and the first row of  $S^{-1}(\phi)$  correspond to  $\lambda_+(\phi)$ , the second ones to  $\lambda_-(\phi)$ .

Since matrices  $S(\phi)$  and  $S^{-1}(\phi)$  are bounded uniformly in  $\phi \in \mathbb{R}$ , this yields  $\|\exp(A(\phi)\nu)\| \leq K$  for  $\nu \geq 0$ . Thus the scheme (11.6) is stable.

The Taylor expansions of the eigenvalues are

$$\lambda_+(\phi) = \frac{7}{3}i\phi + O(\phi^3), \quad \lambda_-(\phi) = \frac{1}{42}i\phi^3 + O(\phi^5).$$

The Taylor expansions of the matrix  $S^{-1}$  is

$$S^{-1}(\phi) = \begin{pmatrix} 12 + 6i\phi - \frac{9}{7}\phi^2 + \dots & -12 \\ 16 + 8i\phi - \frac{16}{7}\phi^2 + \dots & 12 \end{pmatrix}.$$

Now consider the “block” map  $\Pi_h$  defined as

$$(\Pi_h f)_{j,L} = (\Pi_h f)_{j,R} = f(jh).$$

Then we have  $v(\phi, \Pi_h) = (1, 1)^T$  and

$$S^{-1}(\phi)v(\phi, \Pi_h) = \begin{pmatrix} 6i\phi + O(\phi^2) \\ 28 + O(\phi) \end{pmatrix}.$$

945 Thus we get  $\bar{\mathbf{N}} = \{1, 3\}$ ,  $p_1 = 1, p_3 = 0$ . By Lemma 7.6 this yields  $\|\varepsilon_h(t, v_0, \Pi_h)\| \leq C_1 h \|\nabla v_0\|_\infty + C_2 t h^2 \|\nabla^3 v_0\|_\infty$  for each  $v_0 \in C_{per}^3(\mathbb{R})$ .

Now we demonstrate Algorithm 3. Using Algorithm 1 it is easy to check in the sense of  $\Pi_h$  the scheme possesses the truncation error of order  $P_A = 1$  and not of order 2. Since  $A(0) = 0$ , by statement 5 of Theorem 1 there holds  $P = P_A = 1$ , and we can skip step 1 of the algorithm. Now put  $\mathcal{P}_h = \Pi_h$  and denote  $\vec{\epsilon} = (1, 1)^T$ . Consider the system

$$\begin{cases} (\epsilon_1(x, \Pi_1^{(1,2)}))_0 = 0; \\ (\epsilon_1(x^2/2, \Pi_1^{(1,2)}))_0 = 0. \end{cases}$$

Since  $(\epsilon_1(x, \Pi_1))_0 = 0$  and  $L(0) = 0$ , we have  $(\epsilon_1(x, \Pi_1^{(1,2)}))_0 = (\epsilon_1(x, \Pi_1))_0 + L(0)\mathfrak{C}^{(1)}\vec{\epsilon} = 0$ , so the first equation holds for each  $\mathfrak{C}^{(1)}$  and  $\mathfrak{C}^{(2)}$ . The second equation has the form

$$- \left( \Pi_1^{(1,2)}(x) \right)_0 + L_{-1} \left( \Pi_1^{(1,2)}(x) \right)_{-1} + L_0 \left( \Pi_1^{(1,2)}(x) \right)_0 + L_1 \left( \Pi_1^{(1,2)}(x) \right)_1 = 0,$$

which expands to

$$\begin{aligned} & -\mathfrak{C}^{(1)} \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \begin{pmatrix} 0 & -1 \\ 0 & 1/6 \end{pmatrix} \left[ \begin{pmatrix} 1/2 \\ 1/2 \end{pmatrix} + \mathfrak{C}^{(1)} \begin{pmatrix} -1 \\ -1 \end{pmatrix} + \mathfrak{C}^{(2)}\vec{\epsilon} \right] + \\ & + \begin{pmatrix} 0 & 1 \\ -4/3 & 0 \end{pmatrix} \mathfrak{C}^{(2)}\vec{\epsilon} + \begin{pmatrix} 0 & 0 \\ 4/3 & -1/6 \end{pmatrix} \left[ \begin{pmatrix} 1/2 \\ 1/2 \end{pmatrix} + \mathfrak{C}^{(1)} \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \mathfrak{C}^{(2)}\vec{\epsilon} \right] = 0. \end{aligned}$$

Since  $L(0) = L_{-1} + L_0 + L_1 = 0$ , the terms with  $\mathfrak{C}^{(2)}$  negate each other. Simplifying, we obtain the system for  $\mathfrak{C}^{(1)} = \text{diag} \{ \mathfrak{C}_L^{(1)}, \mathfrak{C}_R^{(1)} \}$ :

$$\begin{cases} -\mathfrak{C}_L^{(1)} + \mathfrak{C}_R^{(1)} - \frac{1}{2} = 0; \\ \frac{4}{3}\mathfrak{C}_L^{(1)} - \frac{4}{3}\mathfrak{C}_R^{(1)} + \frac{2}{3} = 0. \end{cases}$$

The solution of this system is

$$\mathfrak{C}^{(1)} = \begin{pmatrix} 0 & 0 \\ 0 & 1/2 \end{pmatrix} + c_1 I, \quad (11.7)$$

and the matrix  $\mathfrak{C}^{(2)}$  is arbitrary. For  $c_1 = 0$  and  $\mathfrak{C}^{(2)} = 0$  we obtain the map  $\Pi_h^{(1,2)}$  prescribed by  $(\Pi_h f)_{j,L}^{(1,2)} = f(jh)$ ,  $(\Pi_h^{(1,2)} f)_{j,R} = f(jh) + (h/2)f'(jh)$ . Note that  $\Pi_h^{(1,2)} f = \bar{\Pi}_h f + O(h^2)$  for  $f$  smooth enough.

Following the algorithm write the system

$$\begin{cases} (\epsilon_1(x, \Pi_1^{(1,3)}))_0 = 0; \\ (\epsilon_1(x^2/2, \Pi_1^{(1,3)}))_0 = 0; \\ (\epsilon_1(x^3/6, \Pi_1^{(1,3)}))_0 = 0. \end{cases}$$

for  $\mathfrak{C}^{(1)}, \mathfrak{C}^{(2)}, \mathfrak{C}^{(3)}$ . The first equation is valid for each  $\mathfrak{C}^{(m)}$ . The second one is equivalent to (11.7). By Lemma 10.7 we can drop the identity term in the expression (11.7) for  $\mathfrak{C}^{(1)}$  by putting  $c_1 = 0$ . The third one is

$$\begin{aligned} & -\mathfrak{C}^{(2)} \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \begin{pmatrix} 0 & -1 \\ 0 & 1/6 \end{pmatrix} \left[ \begin{pmatrix} -1/6 \\ -1/6 \end{pmatrix} + \mathfrak{C}^{(1)} \begin{pmatrix} 1/2 \\ 1/2 \end{pmatrix} + \mathfrak{C}^{(2)} \begin{pmatrix} -1 \\ -1 \end{pmatrix} + \mathfrak{C}^{(3)}\vec{\epsilon} \right] + \\ & + \begin{pmatrix} 0 & 1 \\ -4/3 & 0 \end{pmatrix} \mathfrak{C}^{(3)}\vec{\epsilon} + \begin{pmatrix} 0 & 0 \\ 4/3 & -1/6 \end{pmatrix} \left[ \begin{pmatrix} 1/6 \\ 1/6 \end{pmatrix} + \mathfrak{C}^{(1)} \begin{pmatrix} 1/2 \\ 1/2 \end{pmatrix} + \mathfrak{C}^{(2)} \begin{pmatrix} 1 \\ 1 \end{pmatrix} + \mathfrak{C}^{(3)}\vec{\epsilon} \right] = 0. \end{aligned}$$

Since  $L(0) = 0$ , the terms with  $\mathfrak{C}^{(3)}$  negate each other. Substituting  $\mathfrak{C}_L^{(1)} = 0$ ,  $\mathfrak{C}_R^{(1)} = 1/2$  and simplifying,

$$\begin{cases} -\mathfrak{C}_L^{(2)} + \mathfrak{C}_R^{(2)} - \frac{1}{12} = 0; \\ \frac{4}{3}\mathfrak{C}_L^{(2)} - \frac{4}{3}\mathfrak{C}_R^{(2)} + \frac{1}{6} = 0. \end{cases}$$

This system is inconsistent. This proves that the values  $P = 1$  and  $Q = 2$  are optimal.

950

Note that the use of Algorithm 4 would give the wrong result  $P' = Q' = 1$ . Indeed, following Algorithm 4 we write

$$f^2 = -(\epsilon_1(x^2/2, \Pi_1))_0 \neq 0$$

and  $f^2 \notin \text{Im}L(0)$  since  $\text{Im}L(0) = \{0\}$ .

#### 11.4. Family of schemes with unlimited growth of $\mathfrak{C}$

In Theorem 1 for a scheme with the truncation error of order  $P_A$  and the formal order of accuracy  $P = P_A + 1$  we proved the existence of  $\Pi_h^{(P,P)}$  that gives the truncation error of order  $P$  with an estimate for  $\mathfrak{C}^{(m)}$ ,  $|\mathbf{m}| = P$ . Theorem 3 does not give an analogous estimate on  $\mathfrak{C}^{(m)}$ . In this section we show that this estimate is not possible without additional assumptions.

955

Consider the family of schemes parametrized by  $\gamma > 0$  with 3 DOFs per cell:

$$h \frac{du_j}{dt} + \left( -\frac{1}{12}u_{j+2} + \frac{2}{3}u_{j+1} - \frac{2}{3}u_{j-1} + \frac{1}{12}u_{j-2} \right) + M(u_{j+1} - 2u_j + u_{j-1}) - \gamma Gu_j = 0,$$

where

$$G = \begin{pmatrix} 0 & 1 & -1 \\ -1 & 0 & 1 \\ 1 & -1 & 0 \end{pmatrix}, \quad H = \begin{pmatrix} 0 & 0 & -1 \\ 0 & 0 & 1 \\ 1 & -1 & 0 \end{pmatrix}, \quad M = G + \gamma H.$$

It is of the form (2.8) with the coefficients

$$Z_0 = I, \quad L_0 = -\gamma G - 2M, \quad L_{\pm 1} = M \pm 2I/3, \quad L_{\pm 2} = \mp I/12,$$

where  $I$  is the identity matrix. We will use the mapping  $\Pi_h$  given by  $(\Pi_h f)_{\eta,1} = (\Pi_h f)_{\eta,2} = (\Pi_h f)_{\eta,3} = f(\eta h)$  and the Euclidean norm on  $\mathbb{C}^3$ .

The matrix  $A(\phi)$  for this scheme takes the form

$$\begin{aligned} A(\phi) &= i\phi I - L(\phi) = \\ &= i\phi I - \left( -\frac{I}{12}e^{2i\phi} + (M + 2I/3)e^{i\phi} - \gamma G - 2M + (M - 2I/3)e^{-i\phi} + \frac{I}{12}e^{-2i\phi} \right) = \\ &= i \left( \phi - \frac{4}{3}\sin\phi + \frac{1}{6}\sin 2\phi \right) I + 4M \sin^2(\phi/2) + \gamma G = \\ &= if(\phi)I + \begin{pmatrix} 0 & \gamma + g & -\gamma - g - \gamma g \\ -\gamma - g & 0 & \gamma + g + \gamma g \\ \gamma + g + \gamma g & -\gamma - g - \gamma g & 0 \end{pmatrix}, \end{aligned}$$

where

$$g = g(\phi) = 4 \sin^2(\phi/2), \quad f(\phi) = \phi - \frac{4}{3}\sin\phi + \frac{1}{6}\sin 2\phi = O(\phi^5).$$

Since  $A^*(\phi) = -A(\phi)$ , the scheme is stable with  $K = 1$ . We have

$$v(\phi, \Pi_h) = (\Pi_1 e^{i\phi x})_0 = \vec{\epsilon} = (1, 1, 1)^T,$$

$$\hat{\epsilon}(\phi, \Pi_h) = A(\phi)v(\phi, \Pi_h) = if(\phi) \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + \gamma g(\phi) \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix},$$

so by Lemma 6.10 the optimal value of the order of the truncation error is  $P_A = 1$ .

Introduce the function

$$w(\phi) = \frac{1}{\gamma + g} \begin{pmatrix} \gamma + g + \gamma g \\ \gamma + g + \gamma g \\ \gamma + g \end{pmatrix} = \vec{\epsilon} + \frac{\gamma g(\phi)}{\gamma + g(\phi)} \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}.$$

By construction  $\gamma + g(\phi) \geq \gamma > 0$ . There holds

$$\|w(\phi) - v(\phi, \Pi_h)\| \leq \sqrt{2}g(\phi) \leq \sqrt{2}\phi^2, \quad \|w(\phi)\| \leq \sqrt{3} + \sqrt{2}\phi^2.$$

Since  $A(\phi)w(\phi) \equiv if(\phi)w(\phi)$ , we have

$$\|A(\phi)w(\phi)\| \leq c|\phi|^5(\sqrt{3} + \sqrt{2}\phi^2)$$

960 with  $c$  independent of  $\gamma$  and  $\phi$ . By Lemma 8.1 the scheme possesses the 2nd order of accuracy and the 4th order of the long-time simulation, and the constants in the estimate (2.17) are uniform with respect to  $\gamma$ .

Now we put  $\mathcal{P}_h = \Pi_h$  and find the values of  $\mathfrak{C}^{(m)}$  providing the 4th order of the truncation error in the sense of  $\Pi_h^{(2,4)}$ . We have  $L(0) = -\gamma G$  and

$$-(\epsilon_1(x^2/2, \Pi_1))_0 = -M\vec{\epsilon} = -\gamma(-1, 1, 0)^T.$$

By Lemma 10.4 the diagonal matrix  $\mathfrak{C}^{(2)}$  should satisfy  $-\gamma G\mathfrak{C}^{(2)}\vec{\epsilon} = -\gamma(-1, 1, 0)^T$ , thus we get

$$\mathfrak{C}^{(2)} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} + c_2 I,$$

where  $c_2 \in \mathbb{R}$  is arbitrary.

For any  $\mathfrak{C}^{(2)}$  we have

$$\left( \epsilon_1 \left( \frac{x^3}{6}, \Pi_1^{(2,2)} \right) \right)_0 = \left( \epsilon_1 \left( \frac{x^3}{6}, \Pi_1 \right) \right)_0 - Z_0 \mathfrak{C}^{(2)} (\Pi_1 1)_0 + \sum_{\eta=-2}^2 L_\eta \mathfrak{C}^{(2)} (\Pi_1 x)_\eta = \left( -Z_0 + \sum_{\eta=-2}^2 \eta L_\eta \right) \mathfrak{C}^{(2)} \vec{\epsilon} = 0$$

and the system  $-\gamma G\mathfrak{C}^{(3)}\vec{\epsilon} = 0$  yields  $\mathfrak{C}^{(3)} = c_3 I$ ,  $c_3 \in \mathbb{R}$ .

Further,

$$\begin{aligned} \left( \epsilon_1 \left( \frac{x^4}{24}, \Pi_1^{(2,3)} \right) \right)_0 &= \left( \epsilon_1 \left( \frac{x^4}{24}, \Pi_1 \right) \right)_0 - Z_0 \mathfrak{C}^{(3)} (\Pi_1 1)_0 + \sum_{\eta=-2}^2 L_\eta \mathfrak{C}^{(3)} (\Pi_1 x)_\eta + \sum_{\eta=-2}^2 L_\eta \mathfrak{C}^{(2)} \left( \Pi_1 \frac{x^2}{2} \right)_\eta = \\ &= \frac{1}{12} M\vec{\epsilon} + \left( -Z_0 + \sum_{\eta=-2}^2 \eta L_\eta \right) \mathfrak{C}^{(3)} \vec{\epsilon} + \left( \sum_{\eta=-2}^2 \frac{\eta^2}{2} L_\eta \right) \mathfrak{C}^{(2)} \vec{\epsilon} = \frac{1}{12} M\vec{\epsilon} + M\mathfrak{C}^{(2)} \vec{\epsilon} = \\ &= \frac{1}{12} \gamma \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} + (1 + \gamma) \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} + c_2 \gamma \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} = \left( 1 + \gamma \left( \frac{13}{12} + c_2 \right) \right) \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}. \end{aligned}$$

The solution of  $-\gamma G\mathfrak{C}^{(4)}\vec{\epsilon} = -(\epsilon_1(x^4/24, \Pi_1^{(2,3)}))_0$  is

$$\mathfrak{C}^{(4)} = \left( \frac{1}{\gamma} + \frac{13}{12} + c_2 \right) \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} + c_4 I.$$

965 Looking on the expressions for  $\mathfrak{C}^{(2)}$  and  $\mathfrak{C}^{(4)}$  we see that we cannot choose  $c_2$ ,  $c_3$ , and  $c_4$  such that  $\mathfrak{C}^{(2)}$  and  $\mathfrak{C}^{(4)}$  are bounded uniformly in  $\gamma$  simultaneously.

Note that since  $A(0) = -L(0) = \gamma G$ , for each  $\gamma > 0$  the scheme is “simple”. For  $\gamma = 0$  the scheme possesses the 4th order of accuracy and the 4th order in the long-time simulation, but it is no more “simple”:  $A(0) = 0$ .

## 12. The bad

The crucial difference between the one-dimensional and multidimensional cases can be seen from the proof of Lemma 9.20. Consider a scheme possessing the formal order of accuracy  $P$  and the long-time simulation order  $Q$ . The function  $\mathcal{W}^{(P,Q)}(\phi)$  given by (8.5) can be componentwise represented as a ratio of the two analytical functions. This ratio is bounded as  $\phi \rightarrow 0$  (otherwise the scheme does not possess the orders  $P$  and  $Q$ ). In 1D case (and this remains true for quasi-1D case) this ratio is holomorphic. Then the Taylor expansion for  $\mathcal{W}^{(P,Q)}(\phi)$  by Lemma 6.7 can be transformed to the coefficients of a local map  $\Pi_h^{(P,Q)}$  that gives the  $Q$ -th order of the truncation error.

In the multidimensional case the ratio of two holomorphic functions can be bounded but not holomorphic. The example is  $f(\phi_x, \phi_y) = \phi_x \phi_y / (\phi_x^2 + \phi_x \phi_y + \phi_y^2)$ . Once the function  $\mathcal{W}^{(P,Q)}(\phi)$  behaves like this, the auxiliary mapping becomes non-local, see Section 12.2. Besides, the behavior of the solution error can become surprising. Particularly, the optimal value of the long-time simulation order can be non-integer, see Section 12.4.

Let  $P \in \mathbb{N}$  be the optimal value of the order of accuracy for the scheme (2.8). If one prescribes  $e \in \mathring{\Omega}$ , this leads to the 1D formulation. By Proposition 8.7 for each  $e \in \mathring{\Omega}$  the  $Q$ -th order of the long-time simulation order on  $H_{per,e}^q(\mathbb{R}^d)$  (with  $q$  big enough) is equivalent to the existence of  $\Pi_{h,e}^{(P,Q)}$  providing the  $Q$ -th order of the truncation error on the direction  $e$ . Here we have the following fork of the possible cases.

1. For some  $e \in \mathring{\Omega}$  there exist no  $\mathfrak{C}_e^{(P)}, \dots, \mathfrak{C}_e^{(Q)}$  such that the mapping  $\Pi_{h,e}^{(P,Q)}$  given by (2.12) provides the  $Q$ -th order of the truncation error on  $e$ . Then the scheme does not possess the orders  $P$  and  $Q$ . This may hold even if the resulting system for  $\mathfrak{C}_e^{(P)}, \dots, \mathfrak{C}_e^{(Q)}$  is consistent for almost all  $e \in \Omega$ , see an example in Section 12.1.
2. For each  $e \in \mathring{\Omega}$  there exist  $\mathfrak{C}_e^{(m)}, m = P, \dots, Q$ , such that  $\Pi_{h,e}^{(P,Q)}$  provides the  $Q$ -th order of accuracy on  $e$  and  $\sup_{e \in \mathring{\Omega}} \|\mathfrak{C}_e^{(m)}\| < \infty, m = P, \dots, Q$ . Then by Corollary 6.19 the scheme possesses the formal order of accuracy  $P$  and long-time simulation order  $Q$ . See an example in Section 12.2.
3. For each  $e \in \mathring{\Omega}$  there exist  $\mathfrak{C}_e^{(P)}, \dots, \mathfrak{C}_e^{(Q)}$  such that  $\Pi_{h,e}^{(P,Q)}$  provides the  $Q$ -th order of accuracy on  $e$  but the previous condition does not hold. Let  $P, Q'$  be optimal values of the formal order of accuracy and of the long-time simulation order. Then the following cases are possible.
  - $Q' \geq Q$ . See Section 12.3 for example.
  - $P \leq Q' < Q$ , and  $Q'$  may be not an integer. See Section 12.4 for example.

Throughout this section we consider the Cauchy problem for the transport equation (2.6), (2.7) in  $\mathbb{R}^2$  with  $\omega = 0$ . We put  $\mathbf{a}_1 = (1, 0)^T, \mathbf{a}_2 = (0, 1)^T$ , so  $T = I$ . On  $\mathbb{C}^{M^0}$  the Euclidean norm will be implied. All the schemes considered in this section are artificial and not intended to represent any scheme used in practice. They have skew-Hermitian matrix  $A(\phi)$  and thus are stable with  $K = 1$ . Note that  $\omega = 0$  is not a limitation; to each scheme considered in this section for the case of nonzero  $\omega$  we can add the central difference approximation of  $\omega \cdot \nabla u$  of the order high enough. This will keep  $A(\phi)$  skew-Hermitian and thus preserve the stability.

### 12.1. Minus-special direction with respect to the long-time simulation accuracy

The following scheme has the following properties: on  $C_{per,e}^3(\mathbb{R}^2)$ , it possesses the formal order of accuracy  $P = 1$  and the long-time simulation order  $Q = 1$  if  $e$  is aligned with the vertical axis, and  $P = 1, Q = 2$  otherwise. I. e. the long-time simulation order degrades if the wave vector is aligned with the vertical axis.

Put  $M^0 = \{L, R\}$ . Consider the scheme

$$h \frac{du_{j,k}}{dt} + E(u_{j+1,k} - u_{j-1,k}) + W(u_{j,k+1} - 2u_{j,k} + u_{j,k-1}) = 0,$$

$$E = \begin{pmatrix} -1/2 & 1/2 \\ 1/2 & -1/2 \end{pmatrix}, \quad W = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}.$$

It is of the form (2.8) with the coefficients  $Z_{0,0} = I$ ,

$$L_{0,0} = -2W, \quad L_{\pm 1,0} = \pm E, \quad L_{0,\pm 1} = W,$$

Unspecified coefficients are zero. The mapping  $\Pi_h$  is defined by  $(\Pi_h f)_{\eta,L} = (\Pi_h f)_{\eta,R} = f(\eta_1 h, \eta_2 h)$ , so  $v(\phi, \Pi_h) = \vec{e} = (1, 1)^T$ . The matrix  $A(\phi) = -L(\phi)$  is

$$A(\phi_1, \phi_2) = \begin{pmatrix} i \sin \phi_1 & -i \sin \phi_1 - 4 \sin^2(\phi_2/2) \\ -i \sin \phi_1 + 4 \sin^2(\phi_2/2) & i \sin \phi_1 \end{pmatrix}. \quad (12.1)$$

We have  $A(0) = 0$  and  $A(\phi)\vec{e} = 4 \sin^2(\phi_2/2)(-1, 1)^T$ , thus the optimal values of the truncation error and of the formal order of accuracy are  $P_A = P = 1$ .

We look for the auxiliary mapping  $\Pi_{h,e}^{(1,2)}$  of the form (2.13):

$$\Pi_{h,e}^{(1,2)} f = \Pi_h f + h \mathfrak{C}_e^{(1)} \Pi_h \frac{\partial f}{\partial e} + h^2 \mathfrak{C}_e^{(2)} \Pi_h \frac{\partial^2 f}{\partial e^2}$$

that provides the second order of the truncation error on planar waves aligned with  $e$ . Since  $L(0) = 0$ , the truncation error on quadratic polynomials does not depend on  $\mathfrak{C}_e^{(2)}$ , so it is enough to consider mappings of the form

$$\Pi_{h,e}^{(1,1)} f = \Pi_h f + h \mathfrak{C}_e \Pi_h \frac{\partial f}{\partial e}. \quad (12.2)$$

On the linear function we have  $(\epsilon_1(e \cdot r, \Pi_{1,e}^{(1,1)}))_0 = 0$  for each  $\mathfrak{C}_e$  because  $L(0) = 0$ . For the quadratic function the truncation error is

$$\left( \epsilon_1 \left( \frac{(e \cdot r)^2}{2}, \Pi_{1,e}^{(1,1)} \right) \right)_0 = \sum_{\eta} L_{\eta} \left( \frac{(e \cdot \eta)^2}{2} \vec{e} + (e \cdot \eta) \mathfrak{C}_e \vec{e} \right).$$

Equating to zero and substituting the values of  $L_{\eta}$  we get  $e_y^2 W \vec{e} + 2e_x E \mathfrak{C}_e \vec{e} = 0$ . Substituting the values for  $W$  and  $E$  we get

$$e_x \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} (\mathfrak{C}_e \vec{e}) + e_y^2 \begin{pmatrix} -1 \\ 1 \end{pmatrix} = 0.$$

For  $e_x \neq 0$  the system is consistent and has the solution

$$\mathfrak{C}_e = -\frac{e_y^2}{e_x} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + c_e I, \quad c_e \in \mathbb{R},$$

while for  $e_x = 0$  the system is inconsistent. Thus for each  $e \in \mathring{\Omega}$  with  $e_x \neq 0$  the scheme possesses the first order of accuracy and the second order in the long-time simulation  $C_{per,e}^3(\mathbb{R}^2)$ . The following estimate is valid:

$$\|\varepsilon_h(t, e^{i\alpha \cdot r}, \Pi_h)\| \leq \min \left\{ 2 \frac{\alpha_y^2}{|\alpha_x|} h + C_2 |\alpha|^3 h^2 t, C_3 |\alpha|^2 h t \right\}.$$

Let us write the explicit expression for the solution error. Denote  $a = \sin \phi_1$ ,  $b = 4 \sin^2(\phi_2/2)$ ,  $v = a + bi$ . Then (12.1) rewrites as

$$A(\phi) = \begin{pmatrix} ia & -ia - b \\ -ia + b & ia \end{pmatrix} = S \Lambda S^{-1},$$

where

$$S = \frac{1}{2} \begin{pmatrix} -1 & 1 \\ v/|v| & \bar{v}/|v| \end{pmatrix}, \quad \Lambda = \begin{pmatrix} ia + i|v| & 0 \\ 0 & ia - i|v| \end{pmatrix}, \quad S^{-1} = \begin{pmatrix} -1 & \bar{v}/|v| \\ 1 & v/|v| \end{pmatrix}.$$

Then

$$\hat{\varepsilon}(\phi, \nu, \Pi_h) = (e^{\nu A(\phi)} - I) \vec{e} = S \begin{pmatrix} (\exp(i\nu(a + |v|)) - 1)(\bar{v}/|v| - 1) \\ (\exp(i\nu(a - |v|)) - 1)(\bar{v}/|v| + 1) \end{pmatrix}.$$

Since  $A(0) = 0$ , the function  $\hat{\varepsilon}(\alpha h, t/h, \Pi_h)$  is holomorphic in  $h$  (see Section 9.1). Thus it can be expressed using the Taylor series in  $h$ . Omitting the calculations, we get

$$S^{-1}(\alpha h)\hat{\varepsilon}(\alpha h, t/h, \Pi_h) = i \frac{\alpha_y^2}{4\alpha_x} (1 - \exp(2it\alpha_x))h \begin{pmatrix} 1 \\ 0 \end{pmatrix} + O(h^2), \quad \alpha_x > 0;$$

$$S^{-1}(\alpha h)\hat{\varepsilon}(\alpha h, t/h, \Pi_h) = i \frac{\alpha_y^2}{4\alpha_x} (1 - \exp(2it\alpha_x))h \begin{pmatrix} 0 \\ -1 \end{pmatrix} + O(h^2), \quad \alpha_x < 0.$$

Thus for each  $\alpha_x \neq 0$  the term with  $h$  is a bounded function of the time. In contrast, for  $\alpha_x = 0$  we have

$$\exp(\nu A(\phi)) = \begin{pmatrix} \cos(\nu b) & -\sin(\nu b) \\ \sin(\nu b) & \cos(\nu b) \end{pmatrix}, \quad \hat{\varepsilon}(\alpha h, t/h, \Pi_h) = \begin{pmatrix} \cos(\nu b) - 1 - \sin(\nu b) \\ \cos(\nu b) - 1 + \sin(\nu b) \end{pmatrix},$$

where  $\nu b = 4t \sin^2(\alpha_y h/2)/h$ . It is clear that there is no estimate of the form  $O(h + h^2 t)$  for  $\hat{\varepsilon}(\alpha h, t/h, \Pi_h)$ .

### 12.2. The case of nonlocal auxiliary mapping

1010 In this section we present a scheme that possesses the first order of accuracy and the second order in the long-time simulation, however, there is no local mapping of the form (2.12) that gives the second order of the truncation error.

Consider the scheme with three degrees of freedom per cell given by

$$h \frac{du_{j,k}}{dt} + E(u_{j+1,k} - u_{j-1,k}) + F(u_{j,k+1} - u_{j,k-1}) + W(u_{j+1,k+1} - u_{j+1,k-1} - u_{j-1,k+1} + u_{j-1,k-1}) = 0,$$

$$E = \begin{pmatrix} 1 & 0 & -1 \\ 0 & -1 & 1 \\ -1 & 1 & 0 \end{pmatrix}, \quad F = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 0 & 1 \\ 0 & 1 & -1 \end{pmatrix}, \quad W = \begin{pmatrix} 0 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & -1 & 0 \end{pmatrix}$$

and let  $\Pi_h$  be given by  $(\Pi_h f)_{\eta, \xi} = f(\eta_1 h, \eta_2 h)$  for each  $\xi$ . This scheme is of the form (2.8) with  $Z_{(0,0)} = I$ ,

$$\begin{aligned} L_{(-1,-1)} &= W, & L_{(0,-1)} &= -F, & L_{(1,-1)} &= -W, \\ L_{(-1,0)} &= -E, & L_{(0,0)} &= 0, & L_{(1,0)} &= E, \\ L_{(-1,1)} &= -W, & L_{(0,1)} &= F, & L_{(1,1)} &= W, \end{aligned}$$

(all unspecified coefficients are zero).

We start with the analysis of the truncation error. We have  $(\epsilon_1(f, \Pi_1))_0 = (L\Pi_1 f)_0 = \sum_{\eta} L_{\eta}(\Pi_1 f)_{\eta}$  and  $(\Pi_1 \mathbf{r}^m)_{\eta} = \eta^m \vec{\mathbf{e}}$ , where  $\vec{\mathbf{e}} = (1, 1, 1)^T$ ,

$$\begin{aligned} \sum_{\eta} L_{\eta} &= 0, & \sum_{\eta} \eta_1 L_{\eta} &= 2E, & \sum_{\eta} \eta_2 L_{\eta} &= 2F, \\ \sum_{\eta} \eta_1^2 L_{\eta} &= 0, & \sum_{\eta} \eta_1 \eta_2 L_{\eta} &= 4W, & \sum_{\eta} \eta_2^2 L_{\eta} &= 0. \end{aligned}$$

Since  $E\vec{\mathbf{e}} = F\vec{\mathbf{e}} = 0$ , the scheme possesses the first order of the truncation error; since  $W\vec{\mathbf{e}} \neq 0$ , it does not possess the second order of the truncation error. Since  $L(0) = 0$ , by Theorem 1 the optimal order of accuracy is equal to  $P = 1$ .

Now we consider this scheme on  $C_{per,e}^3(\mathbb{R}^2)$  for some  $e \in \hat{\Omega}$ . We will look for the auxiliary mapping  $\Pi_{h,e}^{(1,1)}$  given by (12.2) providing the second order of the truncation error on  $e$ . Since the scheme possesses the first order of the truncation error, we need to consider the function  $f(\mathbf{r}) = (e \cdot \mathbf{r})^2/2$  only. Write

$$\left( \epsilon_1(f, \Pi_{1,e}^{(1,1)}) \right)_0 = \sum_{\eta} L_{\eta} \left( \frac{(e \cdot \eta)^2}{2} \vec{\mathbf{e}} + (e \cdot \eta) \mathfrak{C}_e \vec{\mathbf{e}} \right) = 4e_x e_y W \vec{\mathbf{e}} + 2(e_x E + e_y F) \mathfrak{C}_e \vec{\mathbf{e}}.$$

Substituting the expressions for  $E$ ,  $F$ , and  $W$  and equating to zero we get

$$\begin{pmatrix} e_x + e_y & -e_y & -e_x \\ -e_y & -e_x & e_x + e_y \\ -e_x & e_x + e_y & -e_y \end{pmatrix} \mathfrak{C}_e \vec{\mathbf{e}} = -2e_x e_y \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}.$$

The solution of this system is

$$\mathfrak{C}_e = \frac{2e_x e_y}{e_x^2 + e_x e_y + e_y^2} \begin{pmatrix} -e_x & 0 & 0 \\ 0 & e_y & 0 \\ 0 & 0 & 0 \end{pmatrix} + c_e I.$$

1015 We can put  $c_e = 0$ . Since  $\mathfrak{C}_e$  is bounded in  $e \in \mathring{\Omega}$ , by Corollary 6.19 the scheme possesses the first order of accuracy and the second order in the long-time simulation on  $C_{per}^3(\mathbb{R}^2)$ .

Now show that there exists no mapping  $\Pi_h^{(1,2)}$  of the form (2.12) providing the second order of the truncation error. Indeed, since  $L(0) = 0$ , the truncation error in the sense of  $\Pi_h^{(1,2)}$  on quadratic polynomials does not depend on  $\mathfrak{C}^{(m)}$  for each  $m$  such that  $|m| = 2$ . Thus we need to check mappings  $\Pi_h^{(1,1)}$  of the form

$$\left( \Pi_h^{(1,1)} f \right)_\eta = (\Pi_h f)_\eta + h \mathfrak{C}^{(x)} \left( \Pi_h \frac{\partial f}{\partial x} \right)_\eta + h \mathfrak{C}^{(y)} \left( \Pi_h \frac{\partial f}{\partial y} \right)_\eta.$$

1020 Considering  $f = x^2/2$  we get  $\mathfrak{C}^{(x)} \equiv \mathfrak{C}_e$  with  $e = (1, 0)^T$ . Considering  $f = y^2/2$  we get  $\mathfrak{C}^{(y)} \equiv \mathfrak{C}_e$  with  $e = (0, 1)^T$ . Thus we get  $\mathfrak{C}^{(x)} = \kappa_x I$ ,  $\mathfrak{C}^{(y)} = \kappa_y I$ ,  $\kappa_x, \kappa_y \in \mathbb{R}$ . It remains to see that the truncation error on  $f = xy$  is equal to  $4Wh$  whatever  $\kappa_x$  and  $\kappa_y$ . Thus the mapping  $\Pi_h^{(1,2)}$  of the form (2.12) providing the second order of the truncation error does not exist.

### 12.3. Fake special direction

1025 Let a scheme of the form (2.8) possess the formal order of accuracy  $P \in \mathbb{N}$  and the long-time simulation order  $Q \in \mathbb{N}$ ,  $Q > P$ . By Lemma 8.7 for each  $e \in \mathring{\Omega}$  there exist coefficients  $\mathfrak{C}_e^{(m)}$ ,  $m = P, \dots, Q$  such that the scheme possesses the truncation error of order  $Q$  in the sense of  $\Pi_{h,e}^{(P,Q)}$  given by (2.13) on  $e$ . But it is generally impossible to define these coefficients such that they satisfy  $\max_{m=P,\dots,Q} \sup_{e \in \mathring{\Omega}} \|\mathfrak{C}_e^{(m)}\| < \infty$ . To show this, we adapt the example of Section 11.4.

Consider the scheme

$$h \frac{du_{j,k}}{dt} + \sum_{\eta_1=-2}^2 \sum_{\eta_2=-2}^2 L_{\eta_1, \eta_2} u_{j+\eta_1, k+\eta_2} = 0$$

with

$$L_{\eta_1, \eta_2} = \left( -c_{\eta_2}^{(2)} + c_{\eta_1}^{(4)} + c_{\eta_2}^{(4)} \right) W + \left( c_{\eta_1}^{(2)} c_{\eta_2}^{(2)} + c_{\eta_2}^{(4)} \right) G,$$

$$W = \begin{pmatrix} 0 & 1 & -1 \\ -1 & 0 & 1 \\ 1 & -1 & 0 \end{pmatrix}, \quad G = \begin{pmatrix} 0 & 0 & -1 \\ 0 & 0 & 1 \\ 1 & -1 & 0 \end{pmatrix},$$

and  $c_\eta^{(m)}$  be the coefficients of  $(5-m)$ -th order finite difference approximations of the  $m$ -th derivative on the uniform mesh with the unit step:

$$c_0^{(2)} = -5/2; \quad c_{\pm 1}^{(2)} = 4/3; \quad c_{\pm 2}^{(2)} = -1/12;$$

$$c_0^{(4)} = 6; \quad c_{\pm 1}^{(4)} = -4; \quad c_{\pm 2}^{(4)} = 1.$$

It is of the form (2.8) with  $Z_0 = I$ ,  $Z_\eta = 0$  for  $\eta \neq 0$  and  $L_\eta$  given above. Define  $\Pi_h$  is by  $(\Pi_h f)_{\eta, \xi} = f(\eta_1 h, \eta_2 h)$  for each  $\xi$ , so  $v(\phi, \Pi_h) = \vec{e} = (1, 1, 1)^T$ . The matrix  $L(\phi)$  is

$$L(\phi) = (F^{(2)}(\phi_2) + F^{(4)}(\phi_1) + F^{(4)}(\phi_2))W + (F^{(2)}(\phi_1)F^{(2)}(\phi_2) + F^{(4)}(\phi_2))G,$$

where

$$F^{(m)}(\psi) = (-i)^m \sum_{\eta} c_\eta^{(m)} \exp(i\eta\psi) = \psi^m + O(|\psi|^{m+1}). \quad (12.3)$$



Then

$$A(\phi) = -L(\phi) = \begin{pmatrix} 0 & -x & x+y \\ x & 0 & -x-y \\ -x-y & x+y & 0 \end{pmatrix},$$

where

$$x \equiv x(\phi) = F^{(2)}(\phi_2) + F^{(4)}(\phi_2) + F^{(4)}(\phi_1), \quad y \equiv y(\phi) = F^{(2)}(\phi_1)F^{(2)}(\phi_2) + F^{(4)}(\phi_2). \quad (12.4)$$

For  $\phi$  in a neighborhood of  $\phi = 0$  introduce the vector

$$w(\phi) = \frac{1}{x(\phi)} \begin{pmatrix} x(\phi) + y(\phi) \\ x(\phi) + y(\phi) \\ x(\phi) \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + \frac{y(\phi)}{x(\phi)} \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}, \quad \phi \neq 0; \quad w(0) = 0.$$

Obviously,  $A(\phi)w(\phi) \equiv 0$ . Using (12.3) and  $v(\phi, \Pi_h) = (1, 1, 1)^T$  in a neighborhood of  $\phi = 0$  we get

$$\|w(\phi) - v(\phi, \Pi_h)\| = \left| \frac{y(\phi)}{x(\phi)} \right| \left\| \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \right\| = \sqrt{2} \left| \frac{F^{(2)}(\phi_1)F^{(2)}(\phi_2) + F^{(4)}(\phi_2)}{F^{(2)}(\phi_2) + F^{(4)}(\phi_2) + F^{(4)}(\phi_1)} \right| \leq 2|\phi|^2.$$

By Lemma 8.1 the scheme possesses the second order of accuracy and the infinite order in the long-time simulation. These values are not optimal ( $\|\hat{e}(\phi, \Pi_h)\| \sim \sqrt{2}|y(\phi)| = O(|\phi|^4)$ , so the optimal value of the solution error is equal to 3) but the scheme may be easily modified by adding a disjoint spurious component that will decrease the formal order of accuracy.

**Proposition 12.1.** Let  $\{\mathfrak{C}_e^{(m)}, m = 2, \dots, 9, e \in \mathring{\Omega}\}$  be a set of diagonal matrices such that for each  $e \in \mathring{\Omega}$  the scheme possesses the truncation error of order 9 on the direction  $e$  in the sense of  $\Pi_{h,e}^{(2,9)}$  given by (2.13). Then

$$\max_{m=2,\dots,9} \sup_{e \in \mathring{\Omega}} \|\mathfrak{C}_e^{(m)}\| = \infty.$$

*Proof.* Assume the converse, then for each  $e = (e_1, e_2) \in \mathring{\Omega}$  for the function

$$w_e(\psi) = v(\psi e, \Pi_{h,e}^{(2,9)}) = \left( I + \sum_{m=2}^9 \psi^m \mathfrak{C}_e^{(m)} \right) \vec{e}$$

in a neighborhood of  $\psi = 0$  there holds

$$\|A(e\psi)w_e(\psi)\| \leq c_e |\psi|^{10}, \quad \|v(e\psi, \Pi_h) - w_e(\psi)\| \leq c_e |\psi|^2, \quad (12.5)$$

$$\max_{m=2,\dots,9} \sup_{e \in \mathring{\Omega}} \left\| \frac{d^m w_e}{d\psi^m}(0) \right\| < \infty. \quad (12.6)$$

The eigenvalues of  $A(\phi)$  are 0 and  $\pm i(x^2 + 2(x+y)^2)^{1/2}$  where  $x$  and  $y$  are given by (12.4). The last two eigenvalues satisfy  $|\lambda| \geq |\phi|^4/2$  in a neighborhood of  $\phi = 0$ . Thus the inequality  $\|A(e\psi)w_e(\psi)\| = O(|\psi|^{10})$  is possible only if there exists  $\gamma_e(\psi) > 0$  such that  $w_e(\psi) = \gamma_e(\psi)w(e\psi) + O(|\psi|^6)$ . Taking  $\psi = 0$  we get  $\gamma_e(0) = 1$ . This means that for  $\phi = \psi e$  we have

$$w_e(\psi) = \gamma_e(\psi) \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + \frac{\gamma_e(\psi)y(e\psi)}{x(e\psi)} \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} + O(|\psi|^6). \quad (12.7)$$

By assumption  $d^n w_e/d\psi^n(0)$  is bounded in  $e$  for  $n = 0, \dots, 4$ . Taking the last component of (12.7) we see that this also holds for the derivatives of  $\gamma_e(\psi)$ ; thus this holds for the derivatives of  $\gamma_e(\psi)y(e\psi)/x(e\psi)$  and  $y(e\psi)/x(e\psi)$ . But the Taylor expansion

$$\frac{y(\psi e_1, \psi e_2)}{x(\psi e_1, \psi e_2)} \approx \frac{\psi^4 e_2^2}{\psi^2 e_2^2 + \psi^4 (1 - 2e_1^2 e_2^2)} = \psi^2 \frac{1}{1 + \psi^2 (e_2^{-2} - 2e_1^2)} = \psi^2 - \psi^4 (e_2^{-2} - 2e_1^2) + O(\psi^6)$$

shows that the 4-th derivative by  $\psi$  of this expression is of order  $e_2^{-2}$ . This contradiction proves the proposition.  $\square$

Note that there exist no mapping  $\Pi_h^{(2,9)}$  that gives the 9-th order of the truncation error. Assuming the converse we get that  $v(\phi, \Pi_h^{(2,9)})$  is a holomorphic function satisfying  $\|v(\phi, \Pi_h^{(2,9)}) - v(\phi, \Pi_h)\| = O(|\phi|^2)$  and  $A(\phi)v(\phi, \Pi_h^{(2,9)}) = O(|\phi|^{10})$ . Thus  $w_e(\psi) = v(e\psi, \Pi_h^{(2,9)})$  satisfies (12.5) and (12.6), which leads to the same contradiction as above.

#### 12.4. Strange special direction

In this section we consider a family of schemes with the following properties. The optimal value of the order of accuracy is equal to  $P = 1$ . For some  $q \in \mathbb{N}$  and each  $e \in \hat{\Omega}$  there exist diagonal matrices  $\mathfrak{C}_e^{(m)}$ ,  $m = 1, \dots, q$ , such that the scheme possesses the truncation error of order  $q$  in the sense of  $\Pi_{h,e}^{(1,q)}$  on the direction  $e$ . However, the scheme does not possess the formal order  $P$  and the long-time simulation order  $q$ . For instance, the optimal value  $Q$  of the long-time simulation order may be equal to  $q - 1$  or  $q - 1/2$ .

Prior to construct the schemes we prove the following lemma. Denote  $\mathbb{R}_+^d = \{(x_1, \dots, x_d) : x_1, \dots, x_d \geq 0\}$  and for  $\mathbf{x}, \mathbf{y} \in \mathbb{R}_+^d$  put  $\mathbf{x}^{\mathbf{y}} = x_1^{y_1} \cdot \dots \cdot x_d^{y_d}$ .

**Lemma 12.2.** *Let  $\mathbf{m}, \mathbf{n}_j \in \mathbb{R}_+^d$ ,  $j = 1, \dots, N$ . The following statements are equivalent.*

1. *The ratio of  $\phi^{\mathbf{m}}$  and  $\sum_j \phi^{\mathbf{n}_j}$  is bounded on  $[0, 1]^d$ ;*
2. *There exist  $\delta_1, \dots, \delta_N \geq 0$  with  $\sum \delta_j = 1$  such that  $\sum \delta_j \mathbf{n}_j \leq \mathbf{m}$ .*

*Proof.* Assume the second statement. Then for each  $\phi \in [0, 1]^d$  there holds

$$\phi^{\mathbf{m}} \leq \phi^{\sum \delta_j \mathbf{n}_j} = \prod_{j=1}^N (\phi^{\mathbf{n}_j})^{\delta_j} \leq \sum_{j=1}^N \phi^{\mathbf{n}_j}.$$

The last inequality is due to the Young inequality. Thus we get the first statement.

Assume the first statement. Let  $\mathbf{q} \in \mathbb{R}_+^d$ . Put  $\phi_1 = \varepsilon^{q_1}, \dots, \phi_d = \varepsilon^{q_d}$  and send  $\varepsilon$  to zero. Since  $\phi^{\mathbf{m}} = \varepsilon^{\mathbf{q} \cdot \mathbf{m}}$  and  $\sum_j \phi^{\mathbf{n}_j} = \sum_j \varepsilon^{\mathbf{q} \cdot \mathbf{n}_j}$  we get

$$\mathbf{q} \cdot \mathbf{m} \geq \min_j \{\mathbf{q} \cdot \mathbf{n}_j\}. \quad (12.8)$$

Assume that the second statement does not hold. Let  $M = \text{Conv}\{\mathbf{n}_j\} + \mathbb{R}_+^d$  where the plus sign means the Minkowski addition. Obviously,  $M$  is convex and  $\mathbf{m} \notin M$ . By the hyperplane separation theorem there exists  $\mathbf{q} \in \mathbb{R}^d$  such that there holds  $\mathbf{q} \cdot \mathbf{a} > C > \mathbf{q} \cdot \mathbf{m}$  for each  $\mathbf{a} \in M$ . Therefore

$$\min_j \{\mathbf{q} \cdot \mathbf{n}_j\} > C > \mathbf{q} \cdot \mathbf{m}. \quad (12.9)$$

It remains to see that  $\mathbf{q} \in \mathbb{R}_+^d$ . Assume the converse, i. e.  $q_k < 0$  for some  $k$ . Let  $\mathbf{a} \in M$ . Consider the sequence

$$\mathbf{a}_l = \mathbf{a} + l(0, \dots, 1, \dots, 0)^T,$$

where unit stays in the  $k$ -th position. This vector belongs to  $M$  by construction, but  $\mathbf{q} \cdot \mathbf{a}_l \rightarrow -\infty$  and thus can't be greater than  $C$  for each  $l$ . This contradiction proves that  $\mathbf{q} \in \mathbb{R}_+^d$ . So (12.9) contradicts (12.8), which proves the lemma.  $\square$

Now we construct the schemes. Put  $M^0 = \{1, 2, 3, 4\}$  and consider the schemes of the form (2.8) with the coefficients  $Z_0 = I$ ,

$$L_{\boldsymbol{\eta}} = c_{\eta_1}^{(\alpha)} c_{\eta_2}^{(\beta)} W + c_{\eta_1}^{(\gamma)} c_{\eta_2}^{(\delta)} E + c_{\eta_1}^{(1)} R,$$

where  $\alpha, \beta, \gamma$  are odd natural numbers and  $\delta$  is an even natural number,

$$W = \begin{pmatrix} 0 & 1 & -1 & 0 \\ -1 & 0 & 1 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \quad E = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \quad R = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix},$$

and  $c_\eta^{(m)}$  are the coefficients of high-order finite difference approximations for the  $m$ -th derivative on the uniform mesh with the unit step such that  $c_{-\eta}^{(m)} = (-1)^m c_\eta^{(m)}$ . The order of the approximation should be not less than  $\max\{\alpha, \beta, \gamma, \delta\} + 1 - m$ . For example, for  $\alpha = 5, \beta = \gamma = 1, \delta = 8$  we can put

$$\begin{aligned} c_{\pm 1}^{(1)} &= \pm 4/5; & c_{\pm 2}^{(1)} &= \mp 1/5; & c_{\pm 3}^{(1)} &= \pm 4/105; & c_{\pm 4}^{(1)} &= \mp 1/280; \\ c_{\pm 1}^{(5)} &= \pm 29/6; & c_{\pm 2}^{(5)} &= \mp 13/3; & c_{\pm 3}^{(5)} &= \pm 3/2; & c_{\pm 4}^{(5)} &= \mp 1/6; \\ c_0^{(8)} &= 70; & c_{\pm 1}^{(8)} &= -56; & c_{\pm 2}^{(8)} &= 28; & c_{\pm 3}^{(8)} &= -8; & c_{\pm 4}^{(8)} &= 1 \end{aligned}$$

(unspecified coefficients are zeros). The mapping  $\Pi_h$  is defined by

$$(\Pi_h f)_{\eta, \xi} = f(\eta h), \quad \xi = 1, \dots, 3; \quad (\Pi_h f)_{\eta, 4} = h \frac{\partial f}{\partial x}(\eta h),$$

so  $v(\phi, \Pi_h) = (1, 1, 1, i\phi_x)^T$ . The matrix  $A(\phi) = -L(\phi)$  is

$$A(\phi) = f(\phi)W - ig(\phi)E - iF^{(1)}(\phi_x)R, \quad (12.10)$$

where

$$f(\phi) = i^{\alpha+\beta+2} F^{(\alpha)}(\phi_x) F^{(\beta)}(\phi_y), \quad g(\phi) = i^{\gamma+\delta-1} F^{(\gamma)}(\phi_x) F^{(\delta)}(\phi_y), \quad (12.11)$$

and

$$F^{(m)}(\psi) = (-i)^m \sum_{\eta} c_\eta^{(m)} \exp(i\eta\psi) = \psi^m + O(|\psi|^{\max\{\alpha, \beta, \gamma, \delta\}+1}).$$

Since  $(A(\phi))^* = -A(\phi)$  the scheme is stable. Note that  $i^{\alpha+\beta+2} = \pm 1$  and  $i^{\gamma+\delta-1} = \pm 1$ . By construction  $F^{(m)}(\psi)$ ,  $f(\phi)$ , and  $g(\phi)$  are real-valued for real-valued  $\psi$  and  $\phi$ .

Determine the order of the truncation error first. We have

$$\hat{\epsilon}(\phi, \Pi_h) = A(\phi)v(\phi, \Pi_h) = (-ig(\phi), 0, ig(\phi), \phi_x F^{(1)}(\phi_x))^T.$$

By assumption  $\gamma + \delta \geq 3$ , then  $\hat{\epsilon}(\phi, \Pi_h) \leq c|\phi|^2$  and the scheme possesses the truncation error of the order  $P_A = 1$ . Since the last component is of the second order as  $\phi_x \rightarrow 0$ , the value  $P_A = 1$  is optimal. Since  $A(0) = 0$ , by Theorem 1 the optimal value of the order of accuracy is  $P = 1$ .

The optimal value of the long-time simulation order will be obtained with the use of Theorem 8.5. Since  $A^*(\phi) = -A(\phi)$ , the functional  $\mathcal{F}(\phi)$  defined by (8.9) can be rewritten as

$$\mathcal{F}(\phi) = |\phi|^{-2P} (\hat{\epsilon}(\phi, \Pi_h))^* \left( A^*(\phi)A(\phi) + |\phi|^{2(Q+1-P)} \right)^{-1} \hat{\epsilon}(\phi, \Pi_h).$$

Denote by  $\mathcal{A}(\phi_x, \phi_y)$  the restriction of  $A(\phi_x, \phi_y)$  to the first three components. Taking into account that  $P = 1$  we get

$$\mathcal{F}(\phi) = |\phi|^{-2} (g(\phi))^2 (1, 0, -1) (\mathcal{A}^*(\phi)\mathcal{A}(\phi) + |\phi|^{2Q})^{-1} \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} + \frac{|\phi_x|^2}{|\phi|^2} \frac{|F^{(1)}(\phi_x)|^2}{|F^{(1)}(\phi_x)|^2 + |\phi|^{2Q}}.$$

The last term (corresponding to the last component of the solution error) is bounded for each  $Q$ . Now write

$$\mathcal{A} \equiv \mathcal{A}(\phi) = \begin{pmatrix} -ig & f & -f \\ -f & 0 & f \\ f & -f & ig \end{pmatrix}, \quad \mathcal{A}^* \mathcal{A} = -\mathcal{A}^2 = \begin{pmatrix} g^2 + 2f^2 & ifg - f^2 & -f^2 \\ -igf - f^2 & 2f^2 & -igf - f^2 \\ -f^2 & igf - f^2 & g^2 + 2f^2 \end{pmatrix},$$

where  $f \equiv f(\phi)$  and  $g \equiv g(\phi)$  are defined by (12.11). Now see that  $(1, 0, -1)^T$  is a right eigenvector of  $\mathcal{A}^* \mathcal{A}$  corresponding to the eigenvalue  $3f^2 + g^2$ . Thus

$$\mathcal{F}(\phi) = 2g^2 |\phi|^{-2} (3f^2 + g^2 + |\phi|^{2Q})^{-1} + < \text{bounded term} >.$$

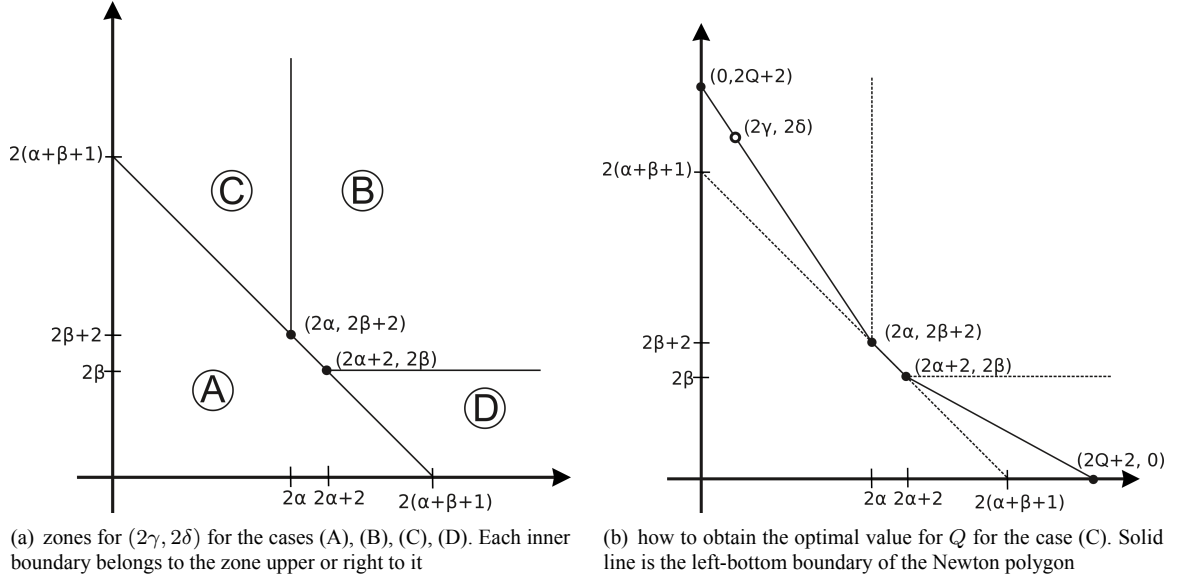


Figure 4: Illustration to the example of Section 12.4

The optimal value of the long-time simulation order is the highest value  $Q$  such that  $\mathcal{F}(\phi)$  is bounded at  $\phi = 0$ , i. e. there exist  $c > 0$  and a neighborhood of  $\phi = 0$  such that

$$\exists c : (g(\phi))^2 \leq c|\phi|^2(3(f(\phi))^2 + (g(\phi))^2 + |\phi|^{2Q}).$$

Now we substitute the expressions (12.11) for  $f(\phi)$  and  $g(\phi)$ . Obviously, replacing a function  $F^{(m)}(\phi_j)$  by  $\phi_j^m$  does not affect the condition to check. The multiplier 3 can be omitted also. Thus we need to check that

$$\exists c : \phi_x^{2\gamma} \phi_y^{2\delta} \leq c \left[ (\phi_x^2 + \phi_y^2)(\phi_x^{2\alpha} \phi_y^{2\beta} + \phi_x^{2\gamma} \phi_y^{2\delta}) + (\phi_x^2 + \phi_y^2)^{Q+1} \right].$$

The expression  $(\phi_x^2 + \phi_y^2)^{Q+1}$  is equivalent to  $|\phi_x|^{2Q+2} + |\phi_y|^{2Q+2}$ . The term  $(\phi_x^2 + \phi_y^2)\phi_x^{2\gamma} \phi_y^{2\delta}$  does not affect the condition to check. Thus we get

$$\exists c : \phi_x^{2\gamma} \phi_y^{2\delta} \leq c \left[ \phi_x^{2\alpha+2} \phi_y^{2\beta} + \phi_x^{2\alpha} \phi_y^{2\beta+2} + \phi_x^{2Q+2} + \phi_y^{2Q+2} \right]. \quad (12.12)$$

We need the biggest value of  $Q$  for which this condition is satisfied. Use Lemma 12.2 to find it. At the left-hand side we have  $\phi^m$ , where  $m = (2\gamma, 2\delta)$ . At the right-hand side we have  $\sum_{j=1}^4 \phi^{n_j}$ , where

$$n_1 = (2\alpha + 2, 2\beta), \quad n_2 = (2\alpha, 2\beta + 2), \quad n_3 = (2Q + 2, 0), \quad n_4 = (0, 2 + 2Q).$$

We have four possible cases here.

- (A)  $\gamma + \delta < \alpha + \beta + 1$ . Then  $Q = \gamma + \delta - 1$ , i. e. the long-time simulation order is defined by the truncation error for the first three components of the solution.
- (B)  $\gamma \geq \alpha$  and  $\delta \geq \beta$ , and  $\gamma + \delta > \alpha + \beta$ . Then  $Q = \infty$ .
- (C)  $\gamma < \alpha$ ,  $\gamma + \delta \geq \alpha + \beta + 1$ . Then  $Q = \beta + \alpha(\delta - \beta - 1)/(\alpha - \gamma)$ .
- (D)  $\delta < \beta$ ,  $\gamma + \delta \geq \alpha + \beta + 1$ . Then  $Q = \alpha + \beta(\gamma - \alpha - 1)/(\beta - \delta)$ .

In Fig. 4(a) for a fixed pair  $(\alpha, \beta)$  the zones for  $(2\gamma, 2\delta)$  corresponding to the cases (A), (B), (C), and (D) are plotted. In Fig. 4(b) for the case (C) we illustrate the location of  $m$  and  $n_j$ . For example, for  $\alpha = 5$ ,  $\beta = \gamma = 1$ ,  $\delta = 8$  the optimal value for  $Q$  can be found from the condition that the points  $(10, 4)$ ,  $(2, 16)$ , and  $(0, 2Q + 2)$  can be

1060 connected by a straight line. Thus we get  $Q = 17/2$ . The solution error possesses the estimate  $O(h + th^{17/2})$ , and the values  $P = 1$ ,  $Q = 17/2$  are optimal.

Compare this result with the accuracy on  $C_{per,e}^1(\mathbb{R}^2)$  for a fixed  $e$ . If  $e = (\pm 1, 0)$  or  $e = (0, \pm 1)$  the first three components of the solution error are zero. Otherwise we can write  $\phi_x = e_1\psi$  and  $\phi_y = e_2\psi$ , and the criterion (12.12) reduces to

$$\exists c : \quad \psi^{2(\gamma+\delta)} \leq c \left[ \psi^{2(\alpha+\beta+1)} + \psi^{2Q+2} \right].$$

So in the case (A) we still have  $Q = \gamma + \delta - 1$ . But in the cases (B), (C), (D) we have  $Q = \infty$ , i. e. the scheme possesses the error estimate

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \leq C_1(e)h\|\nabla v_0\| \quad (12.13)$$

on  $C_{per,e}^1(\mathbb{R}^2)$  for each  $e \in \hat{\Omega}$ .

1065 This example also reveals the difference between the long-time simulation accuracy in the strong sense and in the weak sense. To show this, we need an estimate for the multiplier  $C_1(e)$  in (12.13). Theorem 8.5 does not provide this estimate (one can obtain an estimate for  $C_1(e)$  this way, but this is difficult due to the need of explicit control of the neighborhood of  $\phi = 0$  where intermediate estimates are valid). We will use the direct spectral analysis instead.

Now we put  $\alpha = 3$ ,  $\beta = \gamma = 1$ ,  $\delta = 8$ .

The equation for the last component of the solution is separate from others, so we will consider only the first three components. The eigenvalues of  $\mathcal{A}(\phi)$  are  $\pm ir$  and 0 where  $r = \sqrt{3f^2 + g^2}$  with  $f \equiv f(\phi)$  and  $g \equiv g(\phi)$  given by (12.11). If  $f = g = 0$ , then  $\mathcal{A} = 0$  and the error is exactly zero; below assume  $r \neq 0$ . We don't need to provide an explicit expression of the eigenvectors; it is enough to see that

$$A(\phi) = S(\phi)M(\phi)S^{-1}(\phi),$$

where

$$S(\phi) = \frac{1}{r} \begin{pmatrix} -ig & r & f & 0 \\ -f & 0 & f+ig & 0 \\ f & -r & f & 0 \\ 0 & 0 & 0 & r \end{pmatrix}, \quad M(\phi) = \begin{pmatrix} -f-ig & 2r & 0 & 0 \\ -f(2f+ig)/r & f+ig & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -iF^{(1)}(\phi_x) \end{pmatrix},$$

$$S^{-1}(\phi) = \frac{1}{r^2} \begin{pmatrix} r(f+ig) & -2fr & r(f+ig) & 0 \\ f(2f+ig) & -f(f+ig) & ifg-g^2-f^2 & 0 \\ rf & r(f-ig) & rf & 0 \\ 0 & 0 & 0 & r^2 \end{pmatrix}.$$

Obviously,  $\|S(\phi)\| \leq \tilde{c}$ ,  $\|S^{-1}(\phi)\| \leq \tilde{c}$ , where  $\tilde{c}$  does not depend on  $\phi$ , and by stability for each  $\nu \geq 0$  there holds  $\|\exp(\nu M(\phi))\| \leq \|S(\phi)\| \|\exp(\nu A(\phi))\| \|S(\phi)\| \leq \tilde{c}^2$ . Then

$$\hat{\varepsilon}(\phi, \nu, \Pi_h) = \left( e^{\nu A(\phi)} - I \right) v(\phi, \Pi_h) = S(\phi) \left( e^{\nu M(\phi)} - I \right) S^{-1}(\phi) v(\phi, \Pi_h).$$

Denote by  $\|\cdot\|_{[1,2]}$  and by  $\|\cdot\|_{[1,2,3]}$  the seminorms on  $\mathbb{C}^4$  defined as  $\|(a, b, c, d)^T\|_{[1,2]} = \sqrt{|a|^2 + |b|^2}$ , and  $\|(a, b, c, d)^T\|_{[1,2,3]} = \sqrt{|a|^2 + |b|^2 + |c|^2}$ . Then

$$S^{-1}(\phi)v(\phi, \Pi_h) = \begin{pmatrix} 2ig/r \\ (ifg-g^2)/r^2 \\ (3f-ig)/r \\ i\phi_x \end{pmatrix}, \quad \|S^{-1}(\phi)v(\phi, \Pi_h)\|_{[1,2]} = \sqrt{4r^2 + |g-if|^2} \frac{|g|}{r^2} \leq \sqrt{5} \frac{|g|}{r},$$

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\|_{[1,2,3]} \leq \|S(\phi)\| \left( 1 + \sup_{\nu \geq 0} \|e^{\nu M(\phi)}\|_{[1,2]} \right) \|S^{-1}(\phi)v(\phi, \Pi_h)\|_{[1,2]} \leq \sqrt{5}(\tilde{c} + \tilde{c}^3) \frac{|g|}{r}.$$

Neglecting high order terms we can write  $f \approx i^{\alpha+\beta+2} \phi_x^\alpha \phi_y^\beta$  and  $g \approx i^{\gamma+\delta-1} \phi_x^\gamma \phi_y^\delta$  and, for the chosen parameters,  $f \approx -\phi_x^3 \phi_y$  and  $g \approx \phi_x \phi_y^8$ . Thus

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\|_{[1,2,3]} \leq c \frac{|\phi_y|^7}{|\phi_x|^2 + |\phi_y|^7} = c \frac{|\phi_y|^{7/2}}{|\phi_x|} \frac{|\phi_x| |\phi_y|^{7/2}}{|\phi_x|^2 + |\phi_y|^7} \leq c \frac{|\phi_y|}{|\phi_x|} |\phi|^{5/2}. \quad (12.14)$$

Let  $v_0 \in H_{per}^3(\mathbb{R}^2)$ ,  $N_0$  be its period. Then  $w = \Pi_h v_0$  has the period  $N_0/h$ , so its Fourier series  $F[w](\phi)$  may be nonzero only for  $\phi = 2\pi \mathbf{k}h/N_0$ ,  $\mathbf{k} \in (\mathbb{Z} \cap [-N_0/(2h), N_0/(2h)])^2$ , see (6.1). For such  $\phi$ , we have either  $\phi_x = 0$  or

$$\frac{|\phi_y|}{|\phi_x|} = \frac{|k_2|}{|k_1|} \leq |k_2| = \frac{N_0}{2\pi h} |\phi_y|.$$

If  $\phi_x \neq 0$ , (12.14) yields

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\|_{[1,2,3]} \leq c \frac{N_0}{2\pi h} |\phi|^{7/2} \leq c \frac{N_0}{\sqrt{2\pi h}} |\phi|^3. \quad (12.15)$$

If  $\phi_x = 0$  then  $A(\phi) = 0$  and (12.15) also holds. Adding the disjointed component, we finally obtain

$$\|\hat{\varepsilon}(\phi, \nu, \Pi_h)\| \leq c N_0 h^{-1} |\phi|^3 + 2|\phi|.$$

Note that  $H_{per}^3(\mathbb{R}^2) \subset\subset C_{per}(\mathbb{R}^2)$ , so  $\Pi_h$  is well defined on  $H_{per}^3(\mathbb{R}^2)$  and is a bounded homogeneous mapping of  $H_{per}^3(\mathbb{R}^2)$  to  $V_{per}$ . Let  $F_h$  be the operator that takes  $v_0 \in H_{per}^3(\mathbb{R}^2)$  to  $F_h v_0 = \varepsilon_h(t, v_0, \Pi_h)$ . By Lemma 6.9 we have  $\|F_h\|_{(3,h)} \leq 2\|\Pi_h\|_{(3,h)}$ . Thus the conditions of Lemma 6.14 are satisfied with  $p(|\alpha|) = c N_0 h^2 |\alpha|^3 + 2h|\alpha|$ . For each  $v_0 \in H_{per}^3(\mathbb{R}^2)$  with period  $N_0$  Lemma 6.14 yields

$$\|\varepsilon_h(t, v_0, \Pi_h)\| \equiv \|F_h v_0\| \leq \tilde{C}(h\|\nabla v_0\| + N_0 h^2 \|\nabla^3 v_0\|).$$

1070 This means that the scheme possesses the first formal order of accuracy and the infinite long-time simulation order on  $H_{per}^3(\mathbb{R}^2)$  in sense of  $\Pi_h$  in the weak sense.

Note that the set of values  $\alpha = 3$ ,  $\beta = \gamma = 1$ ,  $\delta = 8$  falls in the case (C) and the optimal values of the formal order of accuracy and of the long-time simulation order in the strong sense are  $P = 1$ ,  $Q = 10$ .

### 13. Order of accuracy depending on the transport velocity

1075 Up to this point, we studied properties of a fixed scheme (2.8) for the transport equation (2.6) with a fixed the transport velocity  $\omega$  in (2.6). In this section we assume that the scheme coefficients linearly depend on  $\omega = (\omega_x, \omega_y)$ , which is constant in space and time.

Consider the following scheme with three degrees of freedom per cell:

$$h \frac{du_{j,k}}{dt} + \omega_x \frac{u_{j+1,k} - u_{j-1,k}}{2} + \omega_y \frac{u_{j,k+1} - u_{j,k-1}}{2} + \omega_x F(u_{j,k+1} - 2u_{j,k} + u_{j,k-1}) + \omega_y G u_{j,k} = 0,$$

$$F = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 0 & -1 \\ -1 & 1 & 0 \end{pmatrix}, \quad G = \begin{pmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

It is of the form (2.8) with  $L_\eta = L_\eta^{(x)} \omega_x + L_\eta^{(y)} \omega_y$ , where

$$Z_{0,0} = I, L_{\pm 1,0}^{(x)} = \pm I/2, L_{0,\pm 1}^{(x)} = F, L_{0,0}^{(x)} = -2F, L_{0,\pm 1}^{(y)} = \pm I/2, L_{0,0}^{(y)} = G,$$

and other coefficients are zero. Let  $\Pi_h$  be given by

$$(\Pi_h f)_{\eta,1} = (\Pi_h f)_{\eta,2} = -h^2 \frac{\partial^2 f}{\partial x^2}(\eta h), \quad (\Pi_h f)_{\eta,3} = f(\eta h).$$

Then  $v(\phi, \Pi_h) = (\phi_x^2, \phi_y^2, 1)^T$  and  $A(\phi) = i\phi \cdot \omega - L(\phi)$  is

$$A(\phi) = i\omega_x(\phi_x - \sin \phi_x)I + i\omega_y(\phi_y - \sin \phi_y)I + \begin{pmatrix} 0 & \omega_y & \omega_x g(\phi_y) \\ -\omega_y & 0 & -\omega_x g(\phi_y) \\ -\omega_x g(\phi_y) & \omega_x g(\phi_y) & 0 \end{pmatrix},$$

where  $g(\phi_y) = 4 \sin^2(\phi_y/2)$ . Since  $(A(\phi))^* = -A(\phi)$ , the scheme is stable with the stability constant  $K = 1$  for each  $\omega$ . We have

$$\hat{\epsilon}(\phi, \Pi_h) = i[\omega_x(\phi_x - \sin \phi_x) + \omega_y(\phi_y - \sin \phi_y)] \begin{pmatrix} \phi_x^2 \\ \phi_y^2 \\ 1 \end{pmatrix} + [\omega_y \phi_x^2 + 4\omega_x \sin^2(\phi_y/2)] \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix}.$$

Thus for  $\omega \neq 0$  we have  $\hat{\epsilon}(\phi, \Pi_h) = O(|\phi|^2)$  and the optimal value of the truncation error is  $P_A = 1$ .

Now we use Theorem 1 to obtain the order of the solution error. We have

$$\begin{aligned} (\Pi_1 x^2/2)_{(\eta_x, \eta_y)} &= (-1, -1, \eta_x^2/2)^T, \\ (\Pi_1 xy)_{(\eta_x, \eta_y)} &= (0, 0, \eta_x \eta_y)^T, \\ (\Pi_1 y^2/2)_{(\eta_x, \eta_y)} &= (0, 0, \eta_y^2/2)^T, \end{aligned}$$

so

$$\begin{aligned} (\epsilon_1(x^2/2, \Pi_h))_0 &= -\omega_y G(-1, -1, 0)^T = \omega_y(-1, 1, 0)^T, \\ (\epsilon_1(xy, \Pi_h))_0 &= 0, \\ (\epsilon_1(y^2/2, \Pi_h))_0 &= \omega_x F(0, 0, 1)^T = \omega_x(-1, 1, 0)^T. \end{aligned}$$

Note that  $L(0) = \omega_y G$ . If  $\omega_y \neq 0$ , then all the vectors  $(\epsilon_1(x^2/2, \Pi_h))_0$ ,  $(\epsilon_1(xy, \Pi_h))_0$ ,  $(\epsilon_1(y^2/2, \Pi_h))_0$  belong to the image of  $L(0)$  and thus by Theorem 1 the scheme possesses the second order of accuracy. For  $\omega_y = 0$ ,  $\omega_x \neq 0$  the vector  $(\epsilon_1(y^2/2, \Pi_h))_0$  does not belong to  $\text{Im } L(0) = \{0\}$ , thus the scheme does not possess the second order of accuracy.

## 14. Conclusion

The Lax – Ryabeknii theorem states that a stable scheme possessing the truncation error  $O(h^{P_A})$  provides the solution error with an estimate  $O(h^{P_A}t)$ . For the classical finite-difference schemes these estimates are optimal, however, for some schemes they can be very rough. The method of auxiliary mapping is a powerful tool to obtain better estimates, which was successfully used for the DG scheme by several authors. However, up to now this method has been applied on ad hoc basis.

We consider a linear  $L_2$ -stable scheme with several degrees of freedom per cell on a uniform mesh and a local mapping  $\Pi_h$  to the mesh space. Based on the method of auxiliary mapping we present a unified approach to find the largest possible  $P$  and, for this  $P$ , the largest possible  $Q$  such that the scheme possesses a solution error estimate of the form  $O(h^P + th^Q)$ . The outline of our results is below.

To analyze the accuracy of the scheme one first finds the optimal order  $P_A$  of the truncation error. The optimal value  $P$  of the formal order of accuracy can be either  $P_A$  or  $P_A + 1$ . To check which of these two cases holds, one considers auxiliary mappings  $\Pi_h^{(P_A+1, P_A+1)}$  of the form (2.12) with undetermined coefficients  $\mathfrak{C}^{(m)}$ . The condition that the truncation error is of order  $P_A + 1$  in the sense of  $\Pi_h^{(P_A+1, P_A+1)}$  forms linear systems for  $\mathfrak{C}^{(m)}$ ,  $|m| = P_A + 1$ . The scheme possesses the order of accuracy  $P = P_A + 1$  if and only if each of these systems is consistent, see Algorithm 2 for details. Below we assume that  $P > 0$ .

For a quasi-1D or a “simple” scheme (see Definitions 6 and 7), the optimal value of  $Q$  in an estimate of the form  $O(h^P + th^Q)$  is an integer (or it can be infinity, if the transport velocity is zero). In order to check whether the scheme possesses this estimate, one considers auxiliary mappings  $\Pi_h^{(P, Q)}$  of the form (2.12) with undetermined coefficients  $\mathfrak{C}^{(m)}$ . The condition that the truncation error is of order  $Q$  in the sense of  $\Pi_h^{(P, Q)}$  forms a linear system for  $\mathfrak{C}^{(m)}$ ,  $P \leq |m| \leq Q$ . If this system is consistent, then the scheme possesses a solution error estimate of the form

$O(h^P + th^Q)$ , otherwise it does not. See Algorithm 3 for details. For “simple” schemes the diagonal matrices  $\mathfrak{C}^{(m)}$  can be found successively, and thus Algorithm 3 simplifies to Algorithm 4.

If a scheme is neither quasi-1D nor simple, the situation may be tricky. The definitions of the long-time simulation order in the weak and strong sense, namely Definition 3 and Definition 4, are not equivalent to each other. The optimal value of the long-time simulation order  $Q$  (with the formal order  $P$ ) in the sense of Definition 4 is generally *not* a natural number. If the scheme possesses the formal order of accuracy  $P$  and the long-time simulation order  $Q$  in the sense of Definition 4, then there exists a homogeneous mapping  $\tilde{\Pi}_h$  such that  $\|\tilde{\Pi}_h f - \Pi_h f\| = O(h^P)$  and the scheme possesses the truncation error of order  $Q$  in the sense of  $\tilde{\Pi}_h$ . However, even if  $Q$  is a natural number, there can be no such local mapping and thus Algorithm 3 generally fails. A general criterion of the  $Q$ -th order in the long-time simulation is given by Theorem 8.5, however we do not have an algorithm to check it.

In the general case one can specify a unit vector  $e$  and restrict the analysis to solutions of the form  $f(t, e \cdot r)$  and write a system for diagonal matrices  $\mathfrak{C}_e^{(m)}$ ,  $m = P, \dots, Q$ , such that the scheme possesses the long-time simulation order  $Q$  in the sense of the mapping  $\Pi_e^{(P,Q)}$  of the form (2.13). If these matrices are bounded over  $e$  on the unit sphere then the scheme possesses the long-time simulation order  $Q$ . If for some  $e \in \hat{\Omega}$  these matrices do not exist then the scheme does not possess the long-time simulation order  $Q$ . If these matrices exist but there is no way to specify them boundedly over the unit sphere then this tells us nothing.

All cases and methods mentioned above were demonstrated on the artificial examples of numerical schemes.

## A. Appendix

**Lemma A.1.** *Let  $n \in \mathbb{N} \cup \{0\}$  and  $w \in H_{per}^n(\mathbb{R}^d)$ . Then the values  $\|\nabla^n w\|^2$  defined by (2.2) and (2.5) coincide.*

*Proof.* Let  $w \in H_{per}^n(\mathbb{R}^d)$  have the Fourier series (2.4). For a multiindex  $\mathbf{n} = (n_1, \dots, n_d)$ ,  $|\mathbf{n}| = n$ , the Fourier coefficients for the function  $D^{\mathbf{n}}w$  are equal to  $(i\alpha_1)^{n_1} \dots (i\alpha_d)^{n_d} w_{\alpha}$ . Hence, by Parseval’s identity

$$\|D^{\mathbf{n}}w\|^2 = \sum_{\alpha \in \mathcal{A}} |\alpha_1|^{2n_1} \dots |\alpha_d|^{2n_d} |w_{\alpha}|^2.$$

Multiply this by  $n!/|\mathbf{n}|!$  and take the sum over all multiindexes  $\mathbf{n}$  such that  $|\mathbf{n}| = n$ . We have

$$\begin{aligned} \|\nabla^n w\|^2 &= \sum_{\alpha \in \mathcal{A}} |w_{\alpha}|^2 \sum_{\substack{n_1 \geq 0, \dots, n_d \geq 0, \\ n_1 + \dots + n_d = n}} \frac{n!}{|\mathbf{n}|!} |\alpha_1|^{2n_1} \dots |\alpha_d|^{2n_d} = \\ &= \sum_{\alpha \in \mathcal{A}} |w_{\alpha}|^2 (\alpha_1^2 + \dots + \alpha_d^2)^n = \sum_{\alpha \in \mathcal{A}} |w_{\alpha}|^2 |\alpha|^{2n}. \end{aligned}$$

The proof of the lemma is complete.  $\square$

**Lemma A.2.** *For  $e \in \hat{\Omega}$  the spaces  $H_{per,e}^q(\mathbb{R}^d)$  and  $C_{per,e}^q(\mathbb{R}^d)$  are infinite dimensional for each  $q \in \mathbb{N} \cup \{0\}$ .*

*Proof.* By definition there exists  $\lambda > 0$  such that  $\lambda e \cdot \mathbf{a}_j = n_j/d_j$ , where  $n_j \in \mathbb{Z}$  and  $d_j \in \mathbb{N}$ , for all  $j = 1, \dots, d$ . Consider  $m \in \mathbb{Z}$ , the functions  $g(x) = \exp(2\pi i m \lambda x)$  and  $f(\mathbf{r}) = g(e \cdot \mathbf{r})$ . Then we have

$$f(\mathbf{r} + N_0 \mathbf{a}_j) = g(e \cdot \mathbf{r} + e \cdot \mathbf{a}_j N_0) = g(e \cdot \mathbf{r}) \exp(2\pi i m \lambda e \cdot \mathbf{a}_j N_0) = f(\mathbf{r}) \exp(2\pi i m \lambda e \cdot \mathbf{a}_j N_0).$$

Assuming  $N_0 = \prod_j d_j$  we get  $f(\mathbf{r} + N_0 \mathbf{a}_j) = f(\mathbf{r})$  for each  $j = 1, \dots, d$ , then  $N_0$  is a period of  $f$ . Hence,  $f \in C_{per,e}^q(\mathbb{R}^d) \subseteq H_{per,e}^q(\mathbb{R}^d)$ .  $\square$

**Lemma A.3.** *Let  $G$  be a bounded domain in  $\mathbb{R}^d$ . For each function  $f \in L_{2,per}(\mathbb{R}^d)$  with period  $N_0$  and each  $h$  such that  $1/h \in \mathbb{N}$  there holds*

$$Y(f) := \frac{1}{(N_0/h)^d} \sum_{\boldsymbol{\eta} \in \{0, \dots, N_0/h-1\}^d} \int_G |f(h\mathbf{r} + hT\boldsymbol{\eta})|^2 d\mathbf{r} \leq C \|f\|^2 \quad (\text{A.1})$$

with  $C > 0$  independent of  $f$ ,  $N_0$ ,  $h$ .



*Proof.* Let  $\square$  be the parallelepiped generated by vectors  $N_0 \mathbf{a}_1, \dots, N_0 \mathbf{a}_d$  and  $\square_1$  be the parallelepiped generated by the vectors  $\mathbf{a}_1, \dots, \mathbf{a}_d$ . If  $G$  is a translation of  $\square_1$  then we have

$$\frac{1}{(N_0/h)^d} \sum_{\boldsymbol{\eta} \in \{0, \dots, N_0/h-1\}^d} \int_G |f(h\mathbf{r} + hT\boldsymbol{\eta})|^2 d\mathbf{r} = \frac{1}{N_0^d} \int_{\square} |f(\mathbf{r})|^2 d\mathbf{r} = |\square_1| \cdot \|f\|^2.$$

Since  $G$  is bounded,  $G$  belongs to the union of  $m \in \mathbb{N}$  translations of  $\square_1$ . Then (A.1) holds with  $C = m|\square_1|$ .  $\square$

**Lemma A.4.** A local mapping with kernel  $\mu \in (W_2^q(G))^*$  is a bounded homogeneous mapping of  $H_{per}^q(\mathbb{R}^d)$  to  $V_{per}$ .

*Proof.* We need to show that for a local mapping  $\Pi_h$  with  $\mu \in (W_2^q(G))^*$  the norm  $\|\Pi_h\| = \sup \|\Pi_h f\| / \|f\|_{(q,h)}$  is bounded. First equip  $W_2^q(G)$  with the norm

$$\|f\|_{W_2^q}^2 = \sum_{|\mathbf{m}| \leq q} \frac{|\mathbf{m}|!}{\mathbf{m}!} \int_G |D^{\mathbf{m}} f(\mathbf{r})|^2 d\mathbf{r}.$$

Let  $f \in H_{per}^q(\mathbb{R}^d)$  have a period  $N_0$ . Put  $N = N_0/h$ . By the norm equivalence on the finite-dimensional space  $\mathbb{C}^{M^0}$  we get

$$\|\Pi_h f\|^2 = \frac{1}{N^d} \sum_{\boldsymbol{\eta} \in \{0, \dots, N-1\}^d} \|(\Pi_h f)_{\boldsymbol{\eta}}\|^2 \leq C \frac{1}{N^d} \sum_{\boldsymbol{\eta} \in \{0, \dots, N-1\}^d} \sum_{\xi \in M^0} |(\Pi_h f)_{\boldsymbol{\eta}, \xi}|^2.$$

The value of  $(\Pi_h f)_{\boldsymbol{\eta}, \xi}$  results from the application of the functional  $\mu_{\xi} \in (W_2^q(G))^*$  to the function  $g(\mathbf{r}) = f(h(\mathbf{r} + T\boldsymbol{\eta}))$  so  $|(\Pi_h f)_{\boldsymbol{\eta}, \xi}| \leq \|\mu_{\xi}\| \|g\|_{W_2^q(G)}$ . Thus

$$\begin{aligned} \|\Pi_h f\|^2 &\leq C|M^0| \left( \max_{\xi \in M^0} \|\mu_{\xi}\| \right)^2 \frac{1}{N^d} \sum_{\boldsymbol{\eta} \in \{0, \dots, N-1\}^d} \sum_{|\mathbf{m}| \leq q} h^{2|\mathbf{m}|} \frac{|\mathbf{m}|!}{\mathbf{m}!} \int_G |(D^{\mathbf{m}} f)(h(\mathbf{r} + T\boldsymbol{\eta}))|^2 d\mathbf{r} = \\ &= C|M^0| \left( \max_{\xi \in M^0} \|\mu_{\xi}\| \right)^2 \sum_{|\mathbf{m}| \leq q} h^{2|\mathbf{m}|} \frac{|\mathbf{m}|!}{\mathbf{m}!} Y(D^{\mathbf{m}} f), \end{aligned}$$

where  $Y$  is given by (A.1). By Lemma A.3 we get

$$\|\Pi_h f\|^2 \leq \tilde{C} \sum_{|\mathbf{m}| \leq q} h^{2|\mathbf{m}|} \frac{|\mathbf{m}|!}{\mathbf{m}!} \|D^{\mathbf{m}} f\|^2 = \tilde{C} \|f\|_{(q,h)}^2,$$

1130 thus  $\|\Pi_h\|_{(q,h)} \leq \tilde{C} < \infty$ .  $\square$

**Lemma A.5.** Let  $\{(e_1 \cdot \mathbf{r})^m, \dots, (e_J \cdot \mathbf{r})^m\}$  be a complete system in the space of homogeneous polynomials of order  $m$ . Then  $\{(e_1 \cdot \mathbf{r})^n, \dots, (e_J \cdot \mathbf{r})^n\}$  is a complete system in the space of homogeneous polynomials of order  $n \leq m$ .

*Proof.* For  $m = n$  this is obvious. If  $m > n$  choose any  $\mathbf{e}_0 \in \Omega$ . Let  $u$  be a homogeneous polynomial of order  $n$ . Then there exists a homogeneous polynomial  $U$  of order  $m$  such that  $u = (\mathbf{e}_0 \cdot \nabla)^{m-n} U$ . For instance, one can construct  $U$  using an  $(m-n)$ -th primitive of  $u(t(\mathbf{e}_0 \cdot \mathbf{r})\mathbf{e}_0 + \mathbf{r} - (\mathbf{e}_0 \cdot \mathbf{r})\mathbf{e}_0)$ . By assumption  $U = \sum_{j=1}^J U_j (e_j \cdot \mathbf{r})^m$ , thus

$$u(\mathbf{r}) = \sum_{j=1}^J U_j (\mathbf{e}_0 \cdot \nabla)^{m-n} (e_j \cdot \mathbf{r})^m = \frac{m!}{n!} \sum_{j=1}^J U_j (\mathbf{e}_0 \cdot \mathbf{e}_j)^{m-n} (e_j \cdot \mathbf{r})^n.$$

Thus  $u \in \text{span}\{(e_j \cdot \mathbf{r})^n, j = 1, \dots, J\}$ .  $\square$

**Lemma A.6.** Let  $n, d \in \mathbb{N}$ ,  $\Omega_d$  be the unit sphere in  $\mathbb{R}^d$  and  $\check{C}_d$  be an open cone in  $\mathbb{R}^d$ . Then there exists a set of vectors  $\{e_k \in \check{\Omega}_d \cap \check{C}_d, k = 1, \dots, C_{n+d-1}^{d-1}\}$  such that  $\{(e_k \cdot \mathbf{r})^n\}$  form a basis in the space of homogeneous polynomials of order  $n$  of  $d$  variables.

*Proof.* First prove that there exists a set of vectors  $\{e_k \in \Omega_d \cap \check{C}_d, k = 1, \dots, C_{n+d-1}^{d-1}\}$  such that  $\{(e_k \cdot \mathbf{r})^n\}$  form a basis in the space of homogeneous polynomials of order  $n$  of  $d$  variables.

Assume without loss that  $(1, 0, \dots, 0) \in \check{C}_d$ .

The proof is by induction in  $d$ . For  $d = 1$  the statement is obvious. Let us prove this for  $d = 2$ . Let  $k \in \{0, \dots, n\}$ . Let  $\alpha = (\alpha_0, \dots, \alpha_n)$  be a set of pairwise different numbers. Find  $\gamma = (\gamma_0, \dots, \gamma_n)$  such that

$$x^k y^{n-k} = \sum_{j=0}^n \gamma_j (x + \alpha_j y)^n.$$

This is equivalent to the system of equations

$$\sum_{j=0}^n \gamma_j \alpha_j^l = \frac{\delta_{l, n-k}}{C_n^l}, \quad l = 0, \dots, n.$$

This is a system of linear equations for  $\gamma_j$  with the Vandermonde determinant and for any set of pairwise different  $\alpha_j$  this system has a unique solution. Thus for  $d = 2$  and any system of vectors  $e_j = (1, \alpha_j), j = 0, \dots, n$  with pairwise different  $\alpha_j$  the polynomials  $(e_j \cdot \mathbf{r})^n$  form a basis in the space of homogeneous polynomials of order  $n$ . Clearly we can choose  $\alpha$  so that the vectors  $(1, \alpha_j)$  belong to  $\check{C}_d$ .

Now let the statement hold for the space dimension  $d \geq 2$ , we will prove it for  $d + 1$ . Let  $\{e_j \in \Omega_d, j = 1, \dots, C_{n+d-1}^{d-1}\}$ , be the set of vectors given by the induction assumption. Let  $\{g_j, j = 1, \dots, d + 1\}$  be the standard basis in  $\mathbb{R}^{d+1}$ . For  $\mathbf{r} \in \mathbb{R}^{d+1}$  denote  $\mathbf{r} = (\mathbf{r}', r_{d+1}), \mathbf{r}' \in \mathbb{R}^d$ .

Using the argument for  $d = 2$ , given any set of pairwise different numbers  $\alpha_j, j = 0, \dots, n$ , we can find  $\gamma_{k,j}, k, j = 0, \dots, n$ , such that

$$(\mathbf{r}' \cdot \mathbf{e}')^k (r_{d+1})^{n-k} = \sum_{j=0}^n \gamma_{k,j} ((\mathbf{r}' \cdot \mathbf{e}') + \alpha_j r_{d+1})^n, \quad \mathbf{e}', \mathbf{r}' \in \mathbb{R}^d, \quad r_{d+1} \in \mathbb{R}.$$

By the induction assumption and Lemma A.5, for any multiindex  $\mathbf{k} = (k_1, \dots, k_d)$  with  $|\mathbf{k}| = k$  we can find coefficients  $\beta_{\mathbf{k},l}, l = 1, \dots, J$  where  $J = C_{n+d-1}^{d-1}$ , such that

$$(\mathbf{r}')^{\mathbf{k}} = \sum_{l=1}^J \beta_{\mathbf{k},l} (\mathbf{r}' \cdot \mathbf{e}_l)^k$$

and thus

$$(\mathbf{r}')^{\mathbf{k}} (r_{d+1})^{n-k} = \sum_{l=1}^J \beta_{\mathbf{k},l} \sum_{j=0}^n \gamma_{k,j} ((\mathbf{r}' \cdot \mathbf{e}_l) + \alpha_j r_{d+1})^n.$$

So for any set of pairwise different numbers  $\{\alpha_j\}_{j=0}^n$ , for the system of vectors

$$\mathbf{e}_{l,j} = (1 + \alpha_j^2)^{-1/2} (\mathbf{e}_l, \alpha_j) \in \mathbb{R}^{d+1}, \quad l = 1, \dots, C_{n+d-1}^{d-1}, \quad j = 0, \dots, n,$$

the set of polynomials  $(\mathbf{e}_{l,j} \cdot \mathbf{r})^n$  forms a complete system in the space of homogeneous polynomials of order  $n$  in  $d + 1$  variables. By construction all these vectors have the form  $t(1, \alpha_{i_1}, \dots, \alpha_{i_d}), i_k \in \{0, \dots, n\}, t > 0$ . Choosing sufficiently small  $\alpha$  we can guarantee that these vectors belong to  $\Omega_{d+1} \cap \check{C}_d$ . Extracting a basis from this complete system we prove the induction statement.

Now let us show that we can find the required system of vectors in  $\check{\Omega}_d \cap \check{C}_d$ . Let  $\mathbf{e}_l \in \Omega_d \cap \check{C}_d, l = 1, \dots, C_{n+d-1}^{d-1}$  be the output of the argument above. It is clear that any set of vectors sufficiently close to a basis also form a basis.

Given  $\varepsilon > 0$ , for each  $\mathbf{e}_l$  we find a vector  $\mathbf{g}_l \in \mathbb{Q}^d$  such that  $\|U^{-1}\mathbf{e}_l - \mathbf{g}_l\| < \varepsilon$ . Then  $\|\mathbf{e}_l - U\mathbf{g}_l\| \leq \varepsilon\|U\|$ , thus  $|||U\mathbf{g}_l|| - 1| = |||\mathbf{g}_l|| - \|\mathbf{e}_l\|| \leq \varepsilon\|U\|$  and

$$\left\| \mathbf{e}_l - \frac{U\mathbf{g}_l}{\|U\mathbf{g}_l\|} \right\| \leq \|\mathbf{e}_l - U\mathbf{g}_l\| + \left\| U\mathbf{g}_l - \frac{U\mathbf{g}_l}{\|U\mathbf{g}_l\|} \right\| = \|\mathbf{e}_l - U\mathbf{g}_l\| + |||U\mathbf{g}_l|| - 1| \leq 2\varepsilon\|U\|.$$

Thus for sufficiently small  $\varepsilon > 0$  the system  $U\mathbf{g}_l/\|U\mathbf{g}_l\|$  is the required set of vectors in  $\tilde{\Omega} \cap \check{C}_d$ .  $\square$

## Acknowledgements

The research is funded by the Russian Science Foundation, project 22-11-00199.

## References

- 1155 [1] Cockburn B., Luskin M., Shu C.-W., Süli E., Enhanced accuracy by post-processing for finite element methods for hyperbolic equations, *Mathematics of Computation* 72 (2003) 577–606.
- [2] Vichnevetsky R., Bowles J. B., *Fourier analysis of numerical approximations of hyperbolic equations*, SIAM, 1982.
- [3] Lowrie R., Compact higher-order numerical methods for hyperbolic conservation laws, Ph.D. thesis, The University of Michigan (1996).
- 1160 [4] Hu F., Hussaini M., Rasetarinera P., An Analysis of the Discontinuous Galerkin Method for Wave Propagation Problems, *J. Comput. Phys.* 151 (1999) 921–946.
- [5] Zhang M., Shu C.-W., An analysis of and a comparison between the discontinuous Galerkin and the spectral finite volume methods, *Computers and Fluids* 34 (2005) 581–592.
- 1165 [6] Guo W., Zhong X., Qui C.-M., Superconvergence of discontinuous Galerkin and local discontinuous Galerkin methods: Eigen-structure analysis based on Fourier approach, *J. Comput. Phys.* 235 (2013) 458–485.
- [7] Balan A., May G., Schöberl J., A stable high-order Spectral Difference method for hyperbolic conservation laws on triangular elements, *J. Comput. Phys.* 231 (2012) 2359–2375.
- 1170 [8] Vanharen J., Puigt G., Vasseur X., Boussuge J.-F., Sagaut P., Revisiting the spectral analysis for high-order spectral discontinuous methods, *J. Comput. Phys.* 337 (2017) 379–402.
- [9] Huynh H. T., A flux reconstruction approach to high-order schemes including discontinuous Galerkin methods, AIAA paper No 2007-4079.
- [10] Vincent P. E., Castonguay P., Jameson A., Insights from von Neumann analysis of high-order flux reconstruction schemes, *J. Comput. Phys.* 230 (2011) 8134–8154.
- 1175 [11] Cheng Y., Shu C.-W., Superconvergence and time evolution of discontinuous Galerkin finite element solutions, *J. Comput. Phys.* 227 (22) (2008) 9612–9627.
- [12] Cheng Y., Shu C.-W., Superconvergence of discontinuous Galerkin and local discontinuous Galerkin schemes for linear hyperbolic and convection-diffusion equations in one space dimension, *SIAM Journal on Numerical Analysis* 47 (2010) 4044–4072.
- 1180 [13] Zhang Q., Shu C.-W., Stability analysis and a priori error estimates of the third order explicit Runge–Kutta discontinuous Galerkin method for scalar conservation laws, *SIAM Journal on Numerical Analysis* 48 (3) (2010) 1038–1063.
- [14] Yang Y., Shu C.-W., Analysis of optimal superconvergence of discontinuous Galerkin method for linear hyperbolic equations, *SIAM Journal on Numerical Analysis* 50 (6) (2012) 3110–3133.

- 1185 [15] Cao W., Zhang Z., Zou Q., Superconvergence of discontinuous Galerkin methods for linear hyperbolic equations, *SIAM Journal on Numerical Analysis* 52 (5) (2014) 2555–2573.
- [16] Bouche D., Ghidaglia J.-M., Pascal F., Error estimate and the geometric corrector for the upwind finite volume method applied to the linear advection equation, *SIAM Journal on Numerical Analysis* 43 (2006) 557–603.
- 1190 [17] Cao W., Shu C.-W., Yang Y., Zhang Z., Superconvergence of discontinuous Galerkin methods for two-dimensional hyperbolic equations, *SIAM Journal on Numerical Analysis* 53 (4) (2015) 1651–1671.
- [18] Cao W., Shu C.-W., Zhang Z., Superconvergence of discontinuous Galerkin methods for 1-D linear hyperbolic equations with degenerate variable coefficients, *ESAIM: M2AN* 51 (6) (2017) 2213–2235.
- [19] Kučera V., Shu C.-W., On the time growth of the error of the DG method for advective problems, *IMA Journal of Numerical Analysis* 39 (2019) 687–712.
- 1195 [20] Liu Y., Shu C.-W., Zhang M., Optimal error estimates of the semidiscrete discontinuous Galerkin methods for two dimensional hyperbolic equations on Cartesian meshes using  $P_k$  elements, *ESAIM: M2AN* 54 (2020) 705–726.
- [21] Liu Y., Shu C.-W., Zhang M., Superconvergence of energy-conserving discontinuous Galerkin methods for linear hyperbolic equations, *Communications on Applied Mathematics and Computation* 1 (2019) 101–116.
- 1200 [22] Xu Y., Meng X., Shu C.-W., Zhang Q., Superconvergence analysis of the Runge–Kutta discontinuous Galerkin methods for a linear hyperbolic equation, *Journal of Scientific Computing* 84 (2020) article number 23.
- [23] Bouche D., Ghidaglia J.-M., Pascal F., Error estimate for the upwind finite volume method for the nonlinear scalar conservation law, *Journal of Computational and Applied Mathematics* 235 (2011) 5394–5410.
- [24] Work C. D., Katz A. J., Aspects of the flux correction method for solving the Navier – Stokes equations on unstructured meshes, *AIAA paper No. 2015-0834*.
- 1205 [25] Kolmogorov A. N., On inequalities between the upper bounds of the successive derivatives of an arbitrary function on an infinite integral, *Amer. Math. Soc. Translations, Ser. 1* 2 (1962) 233–243.
- [26] Kato T., *Perturbation theory for linear operators*, *Grund. math. Wiss., B. 132*, Springer, 1966.
- [27] Kreiss H. O., Über Matrizen die beschränkte Halbgruppen erzeugen, *Mathematica Scandinavica* (1959) 71–80.
- 1210 [28] Bakhvalov P. A., Surnachev M. D., Transformation to a block-diagonal form of matrices generating bounded semigroups, *Linear Algebra and Its Applications* 603 (2020) 275–288.
- [29] Bakhvalov P. A., Surnachev M. D., On analytical families of matrices generating bounded semigroups, *Numerical Analysis and Applications* 14 (2021) 1–12.
- [30] Baumgartel H., *Analytic perturbation theory for matrices and operators*, Basel; Boston; Stuttgart: Birkhäuser Verlag, 1985.
- 1215 [31] Shu C.-W., *Discontinuous Galerkin method for time-dependent problems: survey and recent developments*, Springer, Cham, 2014.