

## RUNGE-KUTTA METHODS FOR PARTIAL DIFFERENTIAL EQUATIONS AND FRACTIONAL ORDERS OF CONVERGENCE

A. OSTERMANN AND M. ROCHE

**ABSTRACT.** We apply Runge-Kutta methods to linear partial differential equations of the form  $u_t(x, t) = \mathcal{L}(x, \partial) u(x, t) + f(x, t)$ . Under appropriate assumptions on the eigenvalues of the operator  $\mathcal{L}$  and the (generalized) Fourier coefficients of  $f$ , we give a sharp lower bound for the order of convergence of these methods. We further show that this order is, in general, fractional and that it depends on the  $L^r$ -norm used to estimate the global error. The analysis also applies to systems arising from spatial discretization of partial differential equations by finite differences or finite element techniques. Numerical examples illustrate the results.

### 1. INTRODUCTION

In this paper we study the order behavior of Runge-Kutta methods applied to certain classes of partial differential equations. As the order of the method will play an essential role throughout this paper, we start by summarizing some basic results related to this concept.

Consider the initial value problem (ODE)

$$(1.1) \quad y' = f(t, y), \quad y(t_0) = y_0,$$

and a so-called one-step method for its numerical solution. Starting from the initial value  $y_0$ , such a method constructs an approximation, say  $y_1$ , to the exact solution  $y(t_0 + h)$  for some step size  $h$  (for Runge-Kutta methods, see (2.2) below). The local error LE is defined as the difference between numerical and exact solution after one step

$$(1.2) \quad \text{LE} = y_1 - y(t_0 + h).$$

A method is said to have (*classical*) order  $p$  if

$$(1.3) \quad \text{LE} = \mathcal{O}(h^{p+1}) \quad \text{for } h \rightarrow 0.$$

Suppose that  $f$  in (1.1) as well as the numerical method satisfy a Lipschitz condition. Then the difference between the exact and numerical solution is seen

Received by the editor July 9, 1990 and, in revised form, June 27, 1991 and November 20, 1991.

1991 *Mathematics Subject Classification.* Primary 65L05, 65M20.

*Key words and phrases.* Runge-Kutta methods, method of lines, partial differential equations.

This work has been supported by the "Fonds national suisse de la recherche scientifique" (grants 20-26244.89 and 8220-025958).

to be  $\mathcal{O}(h^p)$ , uniformly on bounded intervals for sufficiently small  $h$ . In this case we call the method *convergent of order  $p$* . Estimates (1.3) are obtained by expanding the local error as a Taylor series in  $h$ , which implies that  $p$  is an integer. For stiff problems the behavior  $\text{LE} \approx Ch^{p+1}$  is observed for very small values of  $h$  only. This is due to the Lipschitz constant of (1.1), which is large in the presence of stiffness and which is involved in the estimate (1.3). Uniform convergence results can however be obtained for certain subclasses of (1.1) containing problems of arbitrary stiffness. This is the case for stiff differential equations satisfying a one-sided Lipschitz condition ( $B$ -convergence, see [3, 4, 8]) or for singularly perturbed problems [9].

The aim of the present paper is to give sharp orders of convergence for implicit Runge-Kutta methods applied to certain classes of partial differential equations. As differential operators are unbounded, equations of this type can be considered as infinitely stiff. Convergence results for such equations were derived in [1, 2, 6, 7, and 12]. Our approach, however, differs significantly and allows us to prove uniform convergence of order  $\mathcal{O}(h^\alpha)$ , where  $\alpha$  is not necessarily an integer. Order results for Rosenbrock-type methods can be obtained by similar techniques. The authors elaborated this in [15]. The order of multistep methods applied to *nonlinear* parabolic problems has been investigated by Le Roux [13] and more recently by Lubich [14]. Order results for explicit Runge-Kutta methods are given in [18].

A short overview of the present paper is as follows:

In §2 we apply Runge-Kutta methods to linear partial differential equations (PDE) and summarize some basic properties of these methods. Section 3 contains the main result of the paper. Its proof will be given in §4. We will show that the order of Runge-Kutta methods, applied to the PDE (2.1) is, in general, fractional and  $q+2$  at least ( $q$  denoting the stage order of the method). In §5, we prove convergence in sequence spaces, which leads to a deeper understanding why fractional orders occur. In addition, we generalize to nonhomogeneous boundary conditions. In §6 we will give a nice geometrical interpretation of fractional orders as superposition phenomena and thus reinterpret the results of §§3 and 5. We finally discuss the implications of our results to ODE systems coming from semidiscretization of parabolic differential equations.

## 2. RUNGE-KUTTA METHODS FOR LINEAR PDE'S

We consider the following linear partial differential equation

$$(2.1) \quad \begin{aligned} u_t(x, t) &= \mathcal{L}(x, \partial)u(x, t) + f(x, t), & x \in \Omega, \quad 0 \leq t \leq T, \\ u(x, 0) &= u_0(x), & x \in \Omega, \end{aligned}$$

with *homogeneous* boundary conditions. Here,  $\Omega$  is an open and bounded subset of  $\mathbf{R}^d$  with sufficiently smooth boundary  $\partial\Omega$ , and  $\mathcal{L}(x, \partial)$  denotes a differential operator, densely defined in  $L^2(\Omega)$  with spectrum contained in  $\{z \in \mathbf{C} ; \text{Re } z \leq 0\}$ . In order to put (2.1) into the more general framework of an abstract initial value problem in  $L^2(\Omega)$ , we consider the (unbounded) linear operator  $\mathcal{L} : L^2(\Omega) \mapsto L^2(\Omega)$

$$\mathcal{L}a = \mathcal{L}(\cdot, \partial)a \quad \text{for } a \in D(\mathcal{L})$$

with its domain

$$D(\mathcal{L}) = \{a \in L^2(\Omega) ; \mathcal{L}(\cdot, \partial)a \in L^2(\Omega) \text{ and } \mathcal{B}a = 0\}.$$

The derivative  $\mathcal{L}(\cdot, \partial)a$  as well as the boundary condition  $\mathcal{B}a$  are understood in the distributional sense. For example, for an elliptic differential operator of order  $2m$  with homogeneous Dirichlet boundary conditions we have (in the standard notation of Sobolev spaces)

$$D(\mathcal{L}) = H^{2m}(\Omega) \cap H_0^m(\Omega).$$

Considering  $u$  and  $f$  as functions of  $t$  with values in the Hilbert space  $L^2(\Omega)$ , equation (2.1) may be rewritten as

$$(2.1') \quad \begin{aligned} u'(t) &= \mathcal{L}u(t) + f(t), & 0 \leq t \leq T, \\ u(0) &= u_0. \end{aligned}$$

The unknown function  $u(t)$  will be approximated for  $t = t_n := nh$  by a Runge-Kutta method, step by step through the recursion

$$(2.2a) \quad u_{n+1} = u_n + h \sum_{i=1}^s b_i (\mathcal{L}U_i^n + f(t_n + c_i h)),$$

where the internal stages  $U_i^n$  ( $i = 1, \dots, s$ ) are defined by

$$(2.2b) \quad U_i^n = u_n + h \sum_{j=1}^s a_{ij} (\mathcal{L}U_j^n + f(t_n + c_j h)).$$

Here,  $h > 0$  is the step size,  $s$  the number of stages, and  $b_i, a_{ij}, c_i$  the (real) coefficients of the Runge-Kutta method. For notational convenience, we introduce some well-known abbreviations:

$$(2.3) \quad \begin{aligned} b^T &= (b_1, \dots, b_s), & c^k &= (c_1^k, \dots, c_s^k)^T, \\ A &= \begin{pmatrix} a_{11} & \cdots & a_{1s} \\ \vdots & & \vdots \\ a_{s1} & \cdots & a_{ss} \end{pmatrix}, & \mathbf{1} &= (1, \dots, 1)^T \in \mathbf{R}^s. \end{aligned}$$

The following conditions (simplifying assumptions, see [5, p. 214] or [10, p. 203]) on the Runge-Kutta coefficients play an important role throughout the paper,

$$(2.4) \quad C(q): \quad A c^{k-1} = \frac{1}{k} c^k, \quad k = 1, \dots, q.$$

The highest possible value  $q$  in (2.4) is called the *stage order* of the considered Runge-Kutta method. Condition  $C(q)$  says that the quadrature formula with weights  $a_{i1}, \dots, a_{is}$  is of order  $q$  in the interval  $[0, c_i]$ . Note that collocation methods with  $s$  stages satisfy  $C(s)$ . The stability function

$$(2.5) \quad R(z) = 1 + z b^T (I - zA)^{-1} \mathbf{1}$$

associated with the Runge-Kutta method is a rational function of  $z$ . A Runge-Kutta method is called *A-stable*, if its stability function satisfies  $|R(z)| \leq 1$  in the negative complex half-plane  $\mathbf{C}^- = \{z \in \mathbf{C}; \operatorname{Re} z \leq 0\}$ .

Further, consider the following rational function, depending on the coefficients of the Runge-Kutta method and on the stage order  $q$ ,

$$(2.6) \quad W_k(z) = \frac{b^T (I - zA)^{-1} (c^k - kA c^{k-1})}{1 - R(z)} \quad \text{for } k \geq 1.$$

A similar function plays an important role for obtaining  $B$ -convergence results, see [3]. It is evident that the condition  $C(q)$  implies  $W_k(z) \equiv 0$  for  $1 \leq k \leq q$ . In the formulation of the convergence results we shall refer to these conditions

$$(2.7) \quad W_k(z) \equiv 0 \quad \text{for } 1 \leq k \leq q.$$

For many important Runge-Kutta methods, such as Gauss, Radau, and Lobatto methods,  $q$  given by (2.7) is just the stage order (2.4). In general, however, condition (2.7) is weaker than  $C(q)$ , take for example the SDIRK methods treated in [19].

Note that for  $q + 1 \leq k \leq p - 1$  the function  $W_k(z)$  can be rewritten as

$$(2.8) \quad W_k(z) = z^{p-k-1} \widetilde{W}_k(z), \quad \text{with } \widetilde{W}_k(0) \neq 0.$$

This follows from the expansion  $(I - zA)^{-1} = I + zA + \dots$  around  $z = 0$ , the order conditions

$$b^T A^l c^k - k b^T A^{l+1} c^{k-1} = 0, \quad 0 \leq l \leq p - k - 1, \quad 1 \leq k \leq p - 1,$$

and from  $R(z) = 1 + z + \mathcal{O}(z^2)$  for small  $z$  (for  $p \geq 1$ ). Most convergence results of the paper are based on the following assumptions (cf. [4] and [7] for similar concepts):

$$(2.9a) \quad I - zA \quad \text{is regular in } \mathbf{C}^-,$$

$$(2.9b) \quad b^T z(I - zA)^{-1} \quad \text{is bounded in } \mathbf{C}^-,$$

$$(2.9c) \quad W_k(z) \quad \text{is bounded in } \mathbf{C}^- \text{ for } q + 1 \leq k \leq p - 1.$$

Note that condition (2.9) is satisfied for many well-known Runge-Kutta methods, such as the implicit midpoint rule, the trapezoidal rule, or RadauIIA and LobattoIIIC methods. For a differential operator  $\mathcal{L}(x, \vartheta)$  with spectrum contained in  $\{z \in \mathbf{C}; |\arg(z)| \geq \pi - \vartheta\}$  for a certain  $\vartheta$  with  $0 \leq \vartheta < \pi/2$ , condition (2.9) can be weakened to

$$(2.10a) \quad I - zA \quad \text{is regular in } S_\vartheta = \{z \in \mathbf{C}; |\arg(z)| \geq \pi - \vartheta\},$$

$$(2.10b) \quad b^T z(I - zA)^{-1} \quad \text{is bounded in } S_\vartheta,$$

$$(2.10c) \quad W_k(z) \quad \text{is bounded in } S_\vartheta \text{ for } q + 1 \leq k \leq p - 1.$$

Note that Gauss methods satisfy (2.10) but not (2.9) for  $s \geq 3$ . A condition similar to (2.9) has been pointed out already by Brenner et al. (formula (2.6) in [1]).

For a detailed description of Runge-Kutta methods applied to ordinary differential equations, we refer to standard textbooks [5, 10, 11]. The application to partial differential equations from the ODE viewpoint is treated in [17]. A more abstract analysis can be found in [1].

### 3. MODEL PROBLEM AND ORDER RESULTS FOR RUNGE-KUTTA METHODS

Our analysis relies heavily on an operational calculus involving fractional powers of  $\mathcal{L}$ . One possible setting for this is the theory of analytic semigroups. Such an approach was used in [15]. Here we will remain in a more classical framework based on spectral properties of the operator  $\mathcal{L}$ . Our point of view is a slight generalization of a standard assumption in spectral theory, namely  $-\mathcal{L}$  is selfadjoint, positive definite and has a compact inverse.

We therefore consider the class of partial differential equations given by (2.1) where we assume that

(3.1a)  $\mathcal{L}(x, \partial)$  has a pure point spectrum  $\{\lambda_1, \lambda_2, \lambda_3, \dots\}$  with  $\text{Re } \lambda_k \leq 0$

and that the eigenfunctions  $\varphi_k$  satisfy the following properties:<sup>1</sup>

(i) They form a basis of  $L^2(\Omega)$ , so that any  $\psi \in L^2(\Omega)$  can be expressed by the (generalized) Fourier series

$$(3.1b) \quad \psi = \sum_{k=1}^{\infty} \psi_k \varphi_k \quad \text{in } L^2(\Omega);$$

(ii) The mapping

$$(3.1c) \quad \begin{cases} L^2(\Omega) \rightarrow l^2 \\ \psi = \sum \psi_k \varphi_k \mapsto (\psi_k) \end{cases} \quad \text{is a homeomorphism}$$

(i.e., one-to-one and continuous in both directions). By  $l^2$  we denote, as usual, the Hilbert space of sequences  $(\psi_k)_{k \geq 1}$  for which  $\sum |\psi_k|^2 < \infty$ .

Note that (3.1c) implies the existence of two positive constants  $C_1$  and  $C_2$  such that every  $\psi = \sum_{k=1}^{\infty} \psi_k \varphi_k$  in  $L^2(\Omega)$  satisfies

$$(3.2) \quad C_1 \left( \sum_{k=1}^{\infty} |\psi_k|^2 \right)^{1/2} \leq \|\psi\|_{L^2(\Omega)} \leq C_2 \left( \sum_{k=1}^{\infty} |\psi_k|^2 \right)^{1/2}.$$

Assumption (3.1) will mainly be used to handle functions of the operator  $\mathcal{L}$ . Given a single-valued function  $g(z)$ , we may define the operator  $g(\mathcal{L})$  by defining it on the spectrum, i.e.,

$$(3.3a) \quad g(\mathcal{L})\psi = \sum_{k=1}^{\infty} g(\lambda_k) \psi_k \varphi_k \quad \text{for all } \psi \in D(g(\mathcal{L}))$$

with its domain

$$(3.3b) \quad D(g(\mathcal{L})) = \left\{ \psi = \sum_{k=1}^{\infty} \psi_k \varphi_k \in L^2(\Omega) ; (g(\lambda_k) \psi_k) \in l^2 \right\}.$$

If  $g(z)$  is bounded in a sector  $S$  containing the eigenvalues of  $\mathcal{L}$ , then  $D(g(\mathcal{L})) = L^2(\Omega)$  and

$$\|g(\mathcal{L})\psi\| = \left\| \sum_{k=1}^{\infty} g(\lambda_k) \psi_k \varphi_k \right\| \leq \frac{C_2}{C_1} \|\psi\| \sup_{z \in S} |g(z)|$$

by (3.2). Thus, we have proved the following lemma which will be of great use later:

**Lemma 1.** *Let  $g(z)$  be bounded in a sector containing the eigenvalues of  $\mathcal{L}$  and the origin. If  $\mathcal{L}$  satisfies (3.1), then  $g(h\mathcal{L})$  is bounded, independently of  $h$ . □*

The above definition can easily be extended to define fractional powers of the operator  $\mathcal{L}$ :

<sup>1</sup>Some authors call such a set a Riesz basis, e.g. [22].

Let  $\nu$  be any real number. We cut the complex plane along the positive real axis, represent  $z$  uniquely by  $r \exp(i\rho)$  with  $0 \leq \rho < 2\pi$  and set  $z^\nu = r^\nu \exp(i\rho\nu)$ . Thus,  $z^\nu$  becomes single-valued, and the above definition (3.3) applies.

A canonical example for (3.1) is the one-dimensional selfadjoint operator (Sturm-Liouville eigenvalue problem)

$$(3.4) \quad \mathcal{L}(x, \partial) u = \frac{\partial}{\partial x} \left( a(x) \frac{\partial u}{\partial x} \right) - b(x)u.$$

The reader should keep this in mind as a typical problem of the form (2.1).

We further suppose that the solution  $u$  and the source function  $f$  of (2.1) satisfy the following regularity assumption:

$$(3.5) \quad \begin{aligned} u &\in C^{p+1}([0, T], L^2(\Omega)), \\ f &\in C^p([0, T], L^2(\Omega)). \end{aligned}$$

We now define a property  $\mathcal{P}(\nu)$ , depending on a real number  $\nu$ . It will serve to characterize the order of the Runge-Kutta method (2.2) in terms of the differential operator and the Fourier coefficients of the source function. Let  $\|\cdot\|_\nu$  denote the following weighted  $L^2$ -norm:

$$\|\psi\|_\nu = \left( \sum_{k=1}^{\infty} (|\lambda_k|^\nu |\psi_k|)^2 \right)^{1/2}.$$

Now, for  $\mathcal{L}$  and  $f$  of (2.1) satisfying (3.1), we say that  $\mathcal{P}(\nu)$  holds (for a real number  $\nu$ ) if and only if there exists a constant  $C$  such that for all  $t \in [0, T]$

$$(3.6a) \quad \mathcal{P}(\nu) : \quad \|f^{(j)}(t)\|_\nu \leq C, \quad 0 \leq j \leq p.$$

As the domain of  $\mathcal{L}^\nu$  is just the set of all  $\psi \in L^2(\Omega)$  such that  $\|\psi\|_\nu < \infty$ , we have the following equivalent characterization of  $\mathcal{P}(\nu)$  in terms of  $D(\mathcal{L}^\nu)$ :

$$(3.6a') \quad \mathcal{P}(\nu) : \quad f^{(j)}(t) \in D(\mathcal{L}^\nu) \quad \text{for all } t \in [0, T] \text{ and } 0 \leq j \leq p.$$

Note that (3.6a') implies the existence of a constant  $C'$  such that

$$(3.7) \quad \|\mathcal{L}u^{(j)}(t)\|_{\min(\nu, p-j)} \leq C', \quad 0 \leq j \leq p.$$

This formula, which will be very useful in the proof of Theorem 2 below, can be deduced from (3.6) as follows:

(i) It holds for  $\nu = 0$ , since by (2.1) and (3.5),  $\mathcal{L}u^{(j)}(t) = u^{(j+1)}(t) - f^{(j)}(t) \in L^2(\Omega)$ , for  $0 \leq j \leq p$ .

(ii) If  $\nu \in (0, 1]$ , then (3.7) holds for  $j = p$  by (i). For  $j \leq p - 1$  one uses  $u^{(j+1)}(t) \in D(\mathcal{L})$ , which yields

$$\mathcal{L}u^{(j)}(t) = u^{(j+1)}(t) - f^{(j)}(t) \in D(\mathcal{L}^\nu).$$

(iii) Similarly, if  $\nu \in (l, l + 1]$ , where  $l$  is a positive integer, (3.7) holds for  $j \geq p - l$  because it holds for  $0 \leq \nu \leq l$ , and for  $j \leq p - l - 1$  since  $u^{(j+1)}(t) \in D(\mathcal{L}^{l+1})$ .

We also refer to the supremum of all real numbers  $\nu$  with  $\mathcal{P}(\nu)$  and denote it by

$$(3.6b) \quad \bar{\nu} = \sup\{\nu \in \mathbf{R} ; \mathcal{P}(\nu) \text{ holds}\}.$$

As  $f^{(j)}(t) \in L^2(\Omega)$ , we always have  $\mathcal{P}(\nu)$  with  $\nu = 0$ . Further,  $\mathcal{P}(\nu)$  implies  $\mathcal{P}(\mu)$  as long as  $\mu \leq \nu$ . Note, however, that  $\mathcal{P}(\bar{\nu})$  in general does *not* hold. (This can be seen by the canonical example  $\lambda_k = -\pi^2 k^2$  and  $f_k = 1/k$ , which yields  $\bar{\nu} = 1/4$ .) The value of  $\nu$  strongly depends on the asymptotic behavior of the eigenvalues  $\lambda_k$  of the operator  $\mathcal{L}(x, \partial)$ . The asymptotic distribution of the eigenvalues is known for many classes of operators and domains, see, e.g., [20, §5.6.2]. For example, in the case of the Laplacian in the  $d$ -dimensional unit square, we have  $\lambda_k = \mathcal{O}(k^{2/d})$  for the natural ordering  $|\lambda_k| \leq |\lambda_l|$  if  $k \leq l$ .

Let us illustrate the fact that  $\nu$  also depends on the regularity of  $f$ . Consider, for example, the one-dimensional heat equation with homogeneous Dirichlet boundary conditions,

$$(3.8) \quad u_t = u_{xx} + a(x)g(t) \quad \text{in } \Omega = (0, 1).$$

Here,  $\varphi_k(x) = \sin(k\pi x)$ ,  $\lambda_k = -k^2\pi^2$ , and simple integration by part shows that if  $a$  is sufficiently differentiable and satisfies

$$(3.9) \quad a^{(2m)}(0) = a^{(2m)}(1) = 0, \quad m = 0, \dots, M - 1,$$

then we have  $\mathcal{P}(\nu)$  with  $\nu < \bar{\nu} = M + 1/4$ . These artificial boundary conditions for  $M \geq 1$  have been pointed out already in [17]. Similar conditions (formulated in terms of the domain of the operator  $\mathcal{L}$ ) were also noticed by Crouzeix [6]. For the formulation of the theorem, we introduce the following notation. Let  $E(h)$  be a function satisfying

$$(3.10) \quad E(h) = \mathcal{O}(h^\alpha), \quad 0 \leq \alpha < \bar{\alpha},$$

but eventually  $E(h) \neq \mathcal{O}(h^{\bar{\alpha}})$ . In this case, we write

$$E(h) = \mathcal{O}_h(h^{\bar{\alpha}}).$$

It is in general *not* possible to deduce from (3.10) a convergence rate of  $\mathcal{O}(h^{\bar{\alpha}})$ ; as an example, consider  $h \cdot \log h$ , which is  $\mathcal{O}(h^\delta)$  for all  $\delta < 1$  but not  $\mathcal{O}(h)$ .

We now give the main result of the paper concerning the convergence in the  $L^2$ -norm.

**Theorem 2.** *Consider the equation (2.1) satisfying (3.1) with solution and source function satisfying (3.5). Apply an  $A$ -stable Runge-Kutta method (2.2), which has classical order  $p$  and satisfies  $C(q)$  and (2.9). Then we have the following estimate for the global error ( $nh = \text{const}$ ):*

(i) *If  $p \leq q + 2$ , then*

$$(3.11a) \quad \|u(nh) - u_n\|_{L^2} = \mathcal{O}(h^p).$$

(ii) *Let  $\bar{\nu}$  be given by (3.6b). Then*

$$(3.11b) \quad \|u(nh) - u_n\|_{L^2} = \begin{cases} \mathcal{O}(h^p) & \text{if } q + 2 < p < q + 2 + \bar{\nu}, \\ \mathcal{O}_h(h^{q+2+\bar{\nu}}) & \text{if } p \geq q + 2 + \bar{\nu}. \end{cases}$$

Note that condition  $C(q)$  can be replaced by the weaker condition (2.7) without changing the theorem, see [15]. The second line in (3.11b) can of course be written as  $\mathcal{O}(h^{q+2+\nu})$  for a  $\nu < \bar{\nu}$ . From the numerical point of view, it is impossible to notice if the bound  $\nu = \bar{\nu}$  is really attained or not. We therefore use the term *order of convergence* for

$$(3.11') \quad \alpha_2 = \min(p, q + 2 + \bar{\nu})$$

and write

$$\|u(x, nh) - u_n(x)\|_{L^2} = \mathcal{O}_h(h^{\alpha_2}).$$

Formula (3.11) shows that an order reduction from  $p$  down to  $q + 2 + \bar{\nu}$  can occur, depending on the value of  $\bar{\nu}$ . As (3.6a) holds with  $\nu = 0$  for every function  $f$  satisfying (3.5), a lower bound for  $\alpha_2$  is given by

$$\alpha_2 \geq \min(p, q + 2).$$

This was first shown by Brenner et al. [1] for a model problem similar to (2.1). However, for many classes of PDE's (2.1), even fractional order can occur. For the important class of one-dimensional second-order parabolic differential equations (3.8) and regular functions  $a$  (e.g., differentiable), not vanishing on  $\partial\Omega$ , the value of  $\bar{\nu}$  in Theorem 2 is  $\bar{\nu} = 1/4$ , hence we get the order

$$(3.12) \quad \alpha_2 = \min(p, q + 9/4).$$

We now study the actual form of the function  $\mathcal{O}_h(\dots)$  in (3.11b). It depends strongly on the structure of the constant  $C$  in the basic condition (3.6a), i.e., how  $C$  depends on  $\nu$ . We will illustrate this by a simple example and consider

$$(3.13) \quad \|f^{(j)}(t)\|_{\nu} \leq \frac{\bar{C}}{(\bar{\nu} - \nu)^{1/2}}, \quad 0 \leq j \leq p,$$

where  $\bar{C}$  is a constant. Note that (3.13) holds for the one-dimensional heat equation (3.8) if  $a$  is sufficiently regular. We prove the following result:

**Theorem 3.** *Under the assumptions of Theorem 2 and the additional condition (3.13), if  $p \geq q + 2 + \bar{\nu}$ , we have*

$$\|u(nh) - u_n\|_{L^2} = \mathcal{O}(h^{q+2+\bar{\nu}} \sqrt{|\log h|}).$$

Similar order results can also be obtained for  $A(\vartheta)$ -stable methods, i.e., methods whose stability region contains the sector

$$(3.14) \quad S_{\vartheta} = \{z \in \mathbf{C}; |\arg(z)| \geq \pi - \vartheta\}.$$

**Theorem 4.** *Consider the equation (2.1) satisfying (3.1) with solution and source function satisfying (3.5). Suppose that the eigenvalues of the operator  $\mathcal{L}$  lie in the sector (3.14). Apply an  $A(\vartheta)$ -stable Runge-Kutta method (2.2), with classical order  $p$  and satisfying  $C(q)$  and (2.10). Then we have the following estimate for the global error ( $nh = \text{const}$ ):*

(i) *If  $p \leq q + 2$ , then*

$$\|u(nh) - u_n\|_{L^2} = \mathcal{O}(h^p).$$

(ii) *Let  $\bar{\nu}$  be given by (3.6b). Then*

$$\|u(nh) - u_n\|_{L^2} = \begin{cases} \mathcal{O}(h^p) & \text{if } q + 2 < p < q + 2 + \bar{\nu}, \\ \mathcal{O}_h(h^{q+2+\bar{\nu}}) & \text{if } p \geq q + 2 + \bar{\nu}. \quad \square \end{cases}$$

Of course, the theorem remains valid if condition  $C(q)$  is replaced by (2.7).

It is possible to obtain similar convergence results in the  $L^r$ -norm for  $r \neq 2$  by defining a real number  $\bar{\nu}_r$  as in (3.6a'), (3.6b), but with the domain  $D(\mathcal{L}^{\nu})$  taken with respect to  $L^r(\Omega)$ . The only additional assumption concerns a substitute for Lemma 1, i.e., an operational calculus in  $L^r$ . Such an operational



calculus is provided, for instance, within the theory of analytic semigroups.<sup>2</sup> Note that an immediate extension of Lemma 1 within our assumptions seems difficult, as its proof is based on property (3.1c) which does not hold except for  $r = 2$ . With these ingredients one gets

$$(3.15) \quad \|u(nh) - u_n\|_{L^r} = \begin{cases} \mathcal{O}(h^p) & \text{if } p < q + 2 + \bar{\nu}_r, \\ \mathcal{O}_h(h^{q+2+\bar{\nu}_r}) & \text{if } p \geq q + 2 + \bar{\nu}_r. \end{cases}$$

Further, if  $1 \leq r \leq 2$  and  $r'$  are conjugate indices ( $1 = 1/r + 1/r'$ ), we have, owing to the continuity of the embedding  $L^{r'}(\Omega) \subset L^r(\Omega)$ ,

$$(3.16) \quad \alpha_\infty \leq \alpha_{r'} \leq \alpha_2 \leq \alpha_r \leq \alpha_1$$

with  $\alpha_\infty := \inf\{\alpha_r ; 1 \leq r < \infty\}$  and

$$(3.17) \quad \alpha_r = \min(p, q + 2 + \bar{\nu}_r) \quad \text{for } 1 \leq r < \infty.$$

In the simple situation of the one-dimensional heat equation (3.8) with regular  $a(x)$  one has for  $1 < r < \infty$ , see [20, §4.3.3]<sup>3</sup>

$$(3.18) \quad \bar{\nu}_r = M + \frac{1}{2r} \quad \text{with } M \text{ given by (3.9).}$$

For the generic cases  $a(0) \neq 0$  or  $a(1) \neq 0$  this gives

$$(3.19) \quad \alpha_1 = \min(p, q + 5/2), \quad \alpha_\infty = \min(p, q + 2).$$

*Remark.* Theorem 2 (and also Theorems 3 and 4) remains valid for the local error ( $n = 1$ ) of Runge-Kutta methods applied to problem (2.1) with  $p$  replaced by  $p + 1$  in formula (3.11). This explains the asymptotic  $h^{3.25}$ -behavior in the  $L^2$ -norm of the local error of a two-stage DIRK method ( $p = 3, q = 1$ ) observed by Verwer [21, formulas (4.25a), (4.27)] on a problem of class (3.8) with  $\bar{\nu} = 1/4$ . Similarly, the example of [1, p. 13] can be explained by the equivalent formula (3.11) for the local error, with  $\bar{\nu} = 1/4, q = 1$ .

The proof of Theorems 2 and 3 will be given in §4. The proof of Theorem 4 is a straightforward extension of that of Theorem 2. Therefore, it will be omitted. The implications of the theorems to the case where the operator  $\mathcal{L}(x, \partial)$  in (2.1) is discretized in space (by standard finite differences or finite elements) will be discussed in §6. We will show there that the global error of Runge-Kutta methods (satisfying the conditions of Theorems 2, 3, or 4) applied to the discretized system behaves like

$$\begin{cases} h^{\alpha_2} & \text{for } h_0 \leq h \leq H, \\ h^p & \text{for } h \leq h_0 \end{cases}$$

with  $\alpha_2$  given by (3.11') and appropriate constants  $h_0$  and  $H$ .

#### 4. PROOFS OF THEOREM 2 AND THEOREM 3

*Proof of Theorem 2.* We insert the exact solution of (2.1') into (2.2), expand into Taylor series and use the simplifying assumptions  $C(q)$  of (2.4). This

<sup>2</sup>The case  $r = \infty$  requires some modifications. We do not elaborate on this point, cf. also (3.17) below.

<sup>3</sup>In our context, formula (3.18) remains valid for  $r = 1$  (take, for instance,  $a(x) \equiv 1$ ).

yields (here,  $t_n = nh$ )

$$(4.1) \quad u(t_{n+1}) = u(t_n) + h \sum_{i=1}^s b_i u'(t_n + c_i h) + \mathcal{O}(h^{p+1}),$$

$$u(t_n + c_i h) = u(t_n) + h \sum_{j=1}^s a_{ij} u'(t_n + c_j h) + \delta_i,$$

where the defect  $\Delta = (\delta_1, \dots, \delta_n)^T$  is given by

$$(4.2) \quad \Delta = \sum_{k=q+1}^p \frac{h^k}{k!} (c^k - kAc^{k-1})u^{(k)}(t_n) + \mathcal{O}(h^{p+1}).$$

Next we subtract (2.2) from (4.1) and denote the global error by  $e_n = u(nh) - u_n$ . Using the abbreviations  $K_i^n = u(t_n + c_i h) - U_i^n$ ,  $K^n = (K_1^n, \dots, K_s^n)^T$ , we get

$$(4.3) \quad \begin{aligned} e_{n+1} &= e_n + h(b^T \otimes \mathcal{L})K^n + \mathcal{O}(h^{p+1}), \\ (I \otimes \mathcal{F} - hA \otimes \mathcal{L})K^n &= \mathbf{1} \otimes e_n + \Delta. \end{aligned}$$

System (4.3) yields the following recursion formula for the global error:

$$(4.4) \quad e_{n+1} = R(h\mathcal{L})e_n + h(b^T \otimes \mathcal{L})(I \otimes \mathcal{F} - hA \otimes \mathcal{L})^{-1}\Delta + \mathcal{O}(h^{p+1}).$$

Inserting (4.2) into (4.4) and using (2.9b) and Lemma 1 gives

$$(4.5) \quad e_{n+1} = R(h\mathcal{L})e_n + (\mathcal{F} - R(h\mathcal{L})) \sum_{k=q+1}^p \frac{h^{k+1}}{k!} W_k(h\mathcal{L})\mathcal{L}u^{(k)}(t_n) + \mathcal{O}(h^{p+1}).$$

We first show that the term with  $k = p$  is  $\mathcal{O}(h^{p+1})$  and can thus be neglected. By (3.7) one has  $\mathcal{L}u^{(p)}(t) \in L^2(\Omega)$ . Therefore, it is sufficient to prove the boundedness of the operator  $(\mathcal{F} - R(h\mathcal{L}))W_p(h\mathcal{L})$ . But this follows easily from (2.9a,b) and Lemma 1. Thus, we have to consider instead of (4.5) the recursion

$$(4.6) \quad e_{n+1} = R(h\mathcal{L})e_n + (\mathcal{F} - R(h\mathcal{L})) \sum_{k=q+1}^{p-1} \frac{h^{k+1}}{k!} W_k(h\mathcal{L})\mathcal{L}u^{(k)}(t_n) + \mathcal{O}(h^{p+1}).$$

For  $p = q, q + 1$  formula (4.6) simply reads

$$e_{n+1} = R(h\mathcal{L})e_n + \mathcal{O}(h^{p+1}),$$

hence  $A$ -stability together with Lemma 1 gives  $e_n = \mathcal{O}(h^p)$ . This implies (3.11a) for  $p = q, q + 1$ .

For  $p \geq q + 2$ , we solve the recursion (4.6) and use  $e_0 = 0$  to obtain

$$(4.7) \quad \begin{aligned} e_n &= \sum_{k=q+1}^{p-1} \frac{h^{k+1}}{k!} W_k(h\mathcal{L}) \sum_{i=0}^{n-1} R(h\mathcal{L})^{n-i-1} (\mathcal{F} - R(h\mathcal{L}))\mathcal{L}u^{(k)}(t_i) \\ &\quad + \mathcal{O}(h^p), \end{aligned}$$

which, by regrouping the second sum, can be rewritten as

$$\begin{aligned} e_n &= \sum_{k=q+1}^{p-1} \frac{h^{k+1}}{k!} W_k(h\mathcal{L}) \left( \mathcal{L}u^{(k)}(t_{n-1}) - R(h\mathcal{L})^n \mathcal{L}u^{(k)}(t_0) \right. \\ &\quad \left. + \sum_{i=0}^{n-2} R^{n-i-1}(h\mathcal{L})\mathcal{L}(u^{(k)}(t_i) - u^{(k)}(t_{i+1})) \right) + \mathcal{O}(h^p). \end{aligned}$$

Using  $A$ -stability and Lemma 1, we can estimate the global error  $e_n$  in the  $L^2$ -norm by

$$(4.8) \quad \sum_{k=q+1}^{p-1} \frac{h^{k+1}}{k!} \left( \|W_k(h\mathcal{L})\mathcal{L}u^{(k)}(t_{n-1})\| + \|W_k(h\mathcal{L})\mathcal{L}u^{(k)}(t_0)\| \right. \\ \left. + \int_0^{t_{n-1}} \|W_k(h\mathcal{L})\mathcal{L}u^{(k+1)}(s)\| ds \right) + \mathcal{O}(h^p).$$

It thus remains to estimate terms of the form

$$(4.9) \quad h^k \|W_k(h\mathcal{L})h\mathcal{L}u^{(l)}(t)\|, \quad l = k, k + 1$$

for  $q + 1 \leq k \leq p - 1$ .

If  $p < q + 2 + \bar{\nu}$  and hence  $\nu > p - q - 2$  (for  $\nu < \bar{\nu}$  and sufficiently near to  $\bar{\nu}$ ), condition (3.7) implies

$$\mathcal{L}u^{(l)}(t) \in D(\mathcal{L}^{p-k-1}), \quad l = k, k + 1 \quad \text{and} \quad q + 1 \leq k \leq p - 1.$$

We rewrite (4.9) as

$$(4.10) \quad h^p \|W_k(h\mathcal{L})(h\mathcal{L})^{k+1-p}\mathcal{L}^{p-k}u^{(l)}(t)\|$$

and have to show that the operator  $W_k(h\mathcal{L})(h\mathcal{L})^{k+1-p}$  is bounded. This is a consequence of (2.8), (2.9), and Lemma 1.

Finally, for  $p \geq q + 2 + \bar{\nu}$  and hence  $\nu < p - q - 2$  (for  $\nu < \bar{\nu}$ ), condition (3.7) shows that

$$\mathcal{L}u^{(l)}(t) \in D(\mathcal{L}^{q+1+\nu-k}), \quad l = k, k + 1 \quad \text{and} \quad q + 1 \leq k \leq p - 1.$$

We distinguish two cases: If  $k \geq q + 1 + \bar{\nu}$ , then (4.9) can be bounded by

$$(4.11) \quad h^{k+1} \cdot \|W_k(h\mathcal{L})\| \cdot \|\mathcal{L}u^{(l)}(t)\|,$$

which is  $\mathcal{O}(h^{k+1})$  by Lemma 1. If  $k < q + 1 + \bar{\nu}$  we rewrite (4.9) as

$$(4.12) \quad h^{q+2+\nu} \|W_k(h\mathcal{L})(h\mathcal{L})^{k-q-1-\nu}\mathcal{L}^{q+2+\nu-k}u^{(l)}(t)\|$$

and use again (2.8), (2.9), and Lemma 1. This shows that (4.12) is of size  $\mathcal{O}(h^{q+2+\nu})$ , which completes the proof of Theorem 2.  $\square$

*Proof of Theorem 3.* The foregoing proof gives sharp results up to (4.8). Then, because of (3.13), terms of the form (4.12) can be bounded by

$$(4.13) \quad \bar{C} \frac{h^{q+2+\nu}}{(\bar{\nu} - \nu)^{1/2}} \quad \text{for all } \nu < \bar{\nu}.$$

Therefore, they are also bounded by the infimum taken over all  $\nu < \bar{\nu}$ , which is easily seen to be

$$\sqrt{2e} \bar{C} h^{q+2+\bar{\nu}} \sqrt{|\log h|}.$$

This completes the proof of Theorem 3.  $\square$

### 5. ADDITIONAL CONVERGENCE RESULTS

**5.1. Convergence in sequence spaces.** In view of the isomorphism (3.1c), the convergence results obtained in §3 can easily be translated to convergence results

in the sequence space  $l^2$ . Although it is straightforward, we elaborate this point, since the  $l^2$ -interpretation of convergence gives—in our opinion—much more insight why fractional order appears (see §6.1). Further, the  $l^2$  approach easily extends to  $l^r$ -norms with  $r \neq 2$ .

To start, we represent the solution  $u(x, t)$  of (2.1) by the (generalized) Fourier series

$$(5.1) \quad u(x, t) = \sum_{k=1}^{\infty} u_k(t)\varphi_k(x),$$

where  $\{\varphi_k\}$  is the basis of eigenfunctions of  $\mathcal{L}$  satisfying (3.1). Inserting (5.1) into (2.1) and comparing the coefficients of  $\varphi_k$ , we obtain an infinite sequence of ordinary differential equations

$$(5.2) \quad u'_k(t) = \lambda_k u_k(t) + f_k(t), \quad k \geq 1,$$

where  $f_k(t)$  is the Fourier coefficient of  $f$ . The initial value  $u_k(0)$  is the Fourier coefficient of the initial function  $u_0(x)$  of (2.1). Let  $U(t) = (u_1(t), u_2(t), \dots)$  denote the exact solution of (5.2). Applying  $n$ -times a Runge-Kutta method with step size  $h$  to (5.2) yields the numerical solution, which we denote by  $U_n$  and which approximates  $U(nh)$ . Because of (3.1c), the error of a Runge-Kutta method (which has to satisfy, of course, the assumptions of Theorem 2), applied to the decoupled system (5.2), has an asymptotic behavior given by (3.11). Thus, an alternative proof of Theorem 2 is the following:

Consider first the scalar equation

$$(5.3) \quad y'(t) = \lambda y(t) + g(t)$$

with some initial value  $y_0$ . Apply a Runge-Kutta method for its solution,

$$(5.4) \quad \begin{aligned} Y_i &= y_n + h \sum_{j=1}^s a_{ij}(\lambda Y_j + g(t_n + c_j h)), \\ y_{n+1} &= y_n + h \sum_{i=1}^s b_i(\lambda Y_i + g(t_n + c_i h)), \end{aligned}$$

and call  $E_n(\lambda, g, h)$  the global error for  $t = nh$ . Then the error of the whole system (5.2) is given by

$$(5.5) \quad \|U(nh) - U_n\|_{l^2} = \left( \sum_{k=1}^{\infty} |E_n(\lambda_k, f_k, h)|^2 \right)^{1/2}$$

with  $\lambda_k$  of (3.1b) and  $f_k$  of (5.2). It can be estimated as in the proof of Theorem 2.

We like to stress the fact that the above approach does not need any operational calculus and is therefore more elementary. The only ingredient needed is the following lemma whose proof is given by (4.1)–(4.7) with  $\mathcal{L}$  replaced by a complex number  $\lambda$ .

**Lemma 5.** *The global error  $E_n(\lambda, g, h)$  of the Runge-Kutta method (5.4) satisfying the assumptions of Theorem 2 is given by ( $nh = \text{const}$ )*

$$(5.6) \quad \begin{aligned} E_n(\lambda, g, h) &= y(nh) - y_n \\ &= \sum_{k=q+1}^{p-1} \frac{h^{k+1}}{k!} W_k(z) \sum_{i=0}^{n-1} R(z)^{n-i-1} (1 - R(z)) \lambda y^{(k)}(t_i) + \mathcal{O}(h^p), \end{aligned}$$

with  $z = h\lambda$  and  $W_k(z)$  of (2.6).  $\square$

To estimate the global error in  $l^r$ -norms, we consider again the decoupled system (5.2). As in §3, we will characterize the order of convergence with the help of a real number  $\mu$  and the property  $\widehat{\mathcal{P}}_r(\mu)$ ,

$$(5.7a) \quad \widehat{\mathcal{P}}_r(\mu) : \quad (|\lambda_k|^\mu f_k^{(j)}(t)) \in l^r, \quad 0 \leq j \leq p,$$

$$(5.7b) \quad \bar{\mu}_r = \sup\{\mu \in \mathbf{R} ; \widehat{\mathcal{P}}_r(\mu) \text{ holds}\}$$

with  $\lambda_k$  of (3.1b) and  $f_k(t)$  given by (5.2). Recall that  $l^r$  for  $r \geq 1$  is the Banach space of sequences  $(\psi_k)_{k \geq 1}$  satisfying  $\sum |\psi_k|^r < \infty$ . For simplicity we assume that

$$(5.8) \quad (u_k) \in C^{p+1}([0, T], l^1), \quad (f_k) \in C^p([0, T], l^1).$$

Then  $\widehat{\mathcal{P}}_r(\mu)$  holds with  $\mu = 0$  for all  $r \geq 1$ , and since the embeddings ( $r$  and  $r'$  denote conjugate indices, i.e.,  $1 = 1/r + 1/r'$ )

$$l^1 \subset l^r \subset l^2 \subset l^{r'} \subset l^\infty$$

are continuous, we have

$$0 \leq \bar{\mu}_1 \leq \bar{\mu}_r \leq \bar{\mu}_2 \leq \bar{\mu}_{r'} \leq \bar{\mu}_\infty.$$

Using the same techniques as in the  $l^2$ -case shows the following theorem:

**Theorem 6.** *Under the assumptions of Theorem 2 together with (5.7), (5.8) instead of (3.5), (3.6), the global error of a Runge-Kutta method, applied to the system (5.2), is given by ( $nh = \text{const}$ )*

$$\|U(nh) - U_n\|_{l^r} = \begin{cases} \mathcal{O}(h^p) & \text{if } p < q + 2 + \bar{\mu}_r, \\ \mathcal{O}_h(h^{q+2+\bar{\mu}_r}) & \text{if } p \geq q + 2 + \bar{\mu}_r. \quad \square \end{cases}$$

The convergence behavior is thus  $\mathcal{O}_h(h^{\hat{\alpha}_r})$  with

$$(5.9) \quad \hat{\alpha}_r = \min(p, q + 2 + \bar{\mu}_r).$$

A careful analysis of the proof (§4) shows that if the Runge-Kutta method satisfies additional conditions (cf. (5.15) below), then the above theorem can be extended to situations where

$$(|\lambda_k|^{-\gamma} f_k^{(j)}(t)) \in l^1 \quad \text{for some } 0 \leq \gamma \leq 1.$$

We omit details.

In the simple situation (3.8) of the one-dimensional Laplace operator and  $f_k(t) = a_k \cdot g(t) = \mathcal{O}(k^{-1})$  one easily deduces from (5.7), (5.9)

$$(5.10) \quad \hat{\alpha}_1 = \min(p, q + 2) \quad \text{and} \quad \hat{\alpha}_\infty = \min(p, q + 5/2),$$

which is conjugate to (3.19), i.e.,  $\hat{\alpha}_1 = \alpha_\infty$  and  $\hat{\alpha}_\infty = \alpha_1$ .

**5.2. Nonhomogeneous boundary conditions.** The case of nonhomogeneous boundary conditions is more difficult to investigate from a theoretical point of view. Numerical experiments, however, indicate that the order reduction can be more severe than the one predicted by Theorem 2. To illustrate this, we consider the PDE (3.8) with nonhomogeneous boundary conditions

$$u(0, t) = \Phi(t), \quad u(1, t) = \Theta(t).$$

Discretization of the Laplacian by standard 3-point finite differences leads to the ODE system

$$(5.11) \quad U' = L_N U + B_N(t) + F_N(t),$$

where  $L_N$  is the  $N \times N$  matrix

$$(5.12) \quad L_N = -(N+1)^2 \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix}$$

and  $B_N(t) = (N+1)^2(\Phi(t), 0, \dots, 0, \Theta(t))^T$ . Introducing the affine function

$$w(x, t) = x\Theta(t) + (1-x)\Phi(t)$$

allows us to rewrite (5.11) as

$$(5.13) \quad U' = L_N(U - W) + F_N(t),$$

where  $W = (w_1(t), \dots, w_N(t))^T$  with  $w_k(t) = w(\frac{k}{N+1}, t)$ . As the eigenvectors of (5.12) are orthogonal, system (5.13) can be decoupled by an orthogonal transformation  $Q$  into the diagonal system ( $V = QU$ )

$$(5.14) \quad V' = \Lambda(V - QW) + QF_N(t),$$

where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$  with the eigenvalues  $\lambda_k$  of (5.12). Note that Runge-Kutta methods are invariant under the transformation from (5.11) to (5.14). System (5.14) consists of  $N$  scalar differential equations of the Prothero-Robinson type (see [16]),

$$v'_k = \lambda_k(v_k - a_k(t)) + g_k(t), \quad 1 \leq k \leq N.$$

Instead of (2.10c) we consider the conditions

$$(5.15) \quad zW_k(z) \text{ is bounded in } \{z \in \mathbf{C}; |\arg(z)| \geq \pi - \vartheta\} \text{ for } q+1 \leq k \leq p.$$

Many well-known Runge-Kutta methods such as Radau IIA or Lobatto IIIC methods fulfill (5.15). An analysis similar to that made in §5 for the infinite system (5.2) is now possible, essentially with  $q+2$  replaced by  $q+1$ . Suppose that the Runge-Kutta method satisfies the assumptions of Theorem 4, but with (5.15) instead of (2.10c), and apply it to (5.11). Then its global error behaves in the Euclidian norm like  $C \cdot h^{\beta_2}$  (for  $h$  not too small) with

$$(5.16) \quad \beta_2 = \min(p, q+1 + \bar{\chi}).$$

Here,  $\chi$  (and  $\bar{\chi}$ ) is defined by (3.6) with  $f_k^{(j)}$  replaced by  $a_k$ . As  $\lambda_k \approx -\pi^2 k^2$

and

$$a_k(t) \approx \frac{(-1)^{k+1}}{k\pi} \Theta(t) + \frac{1}{k\pi} \Phi(t),$$

we have  $\bar{\chi} = 1/4$  for general  $\Theta$  and  $\Phi$  and hence (compare with (3.12))

$$(5.17) \quad \beta_2 = \min(p, q + 5/4).$$

Note that Gauss methods with an even number of stages do not satisfy (5.15). As a consequence, one loses another power of  $h$  when applying these methods to (5.11) and gets

$$(5.18) \quad \beta_2 = \min(p, q + 1/4)$$

instead of (5.17).

*Remark.* Formula (5.17) remains valid for the local error with  $p$  replaced by  $p + 1$ . This explains perfectly the asymptotic  $h^{2.25}$ -behavior in the  $L^2$ -norm of the local error of a two-stage DIRK method ( $p = 3, q = 1$ ) observed by Verwer [21, formulas (4.25b), (4.27)] on a problem of class (3.8) with nonhomogeneous boundary conditions. Similarly, formula (5.17) confirms the order of convergence observed by Verwer for the same class of methods [21, Table 4.1].

### 6. MORE ON ERROR BEHAVIOR

**6.1. Superposition.** The proof of Theorem 2 leads to a nice geometrical interpretation of the encountered fractional order. We consider again equation (5.3). The global error of a Runge-Kutta method, applied to it, has been derived in Lemma 5. For a fixed value of  $\lambda \in \{z \in \mathbb{C} ; \text{Re } z \leq 0\}$  with  $|\lambda|$  sufficiently large, this term exhibits two different  $h$ -behaviors:

(i) For  $|h\lambda|$  large, the function  $E_n(\lambda, g, h)$  in (5.6) is (see also (4.8))

$$(6.1) \quad E_n(\lambda, g, h) \approx F_1(\lambda)h^{q+2-\omega},$$

where the integer  $\omega$  depends on the method and is for  $p \geq q + 2$  determined by the behavior of the rational function  $W_{q+1}(z)$  at infinity, i.e.,  $W_{q+1}(z) = \mathcal{O}(z^{-\omega})$  for  $z \rightarrow \infty$ . Radau IIA methods ( $s \geq 3$ ), for example, have  $\omega = 2$ . This can be seen from  $b_i = a_{si}$ , which implies  $c_s = 1$  and  $b^T A^{-1}(c^{q+1} - qAc^q) = 0$ .

(ii) For  $h \rightarrow 0$ , (2.8) shows that  $W_k(z) = \mathcal{O}(h^{p-k-1})$ , and therefore  $E_n$  behaves like

$$(6.2) \quad E_n(\lambda, g, h) \approx F_2(\lambda)h^p.$$

The constants  $F_1, F_2$  depend on the method and on  $\lambda$ , but not on  $h$ . Plotted in double-logarithmic scale, the function  $E_n(\lambda, g, h)$  consists thus essentially of two segments with slopes  $q + 2 - \omega$  and  $p$ , respectively.

Considering now a sequence of equations (5.3) with different values of  $\lambda$  gives a sequence of self-similar curves  $E_n(\lambda, g, h)$  which are, however, displaced and therefore superpose each other. The global error of the whole system (5.2) for different values of  $h$  is thus dominated by curves associated with different components of (5.2). This is particularly evident in the  $l^\infty$ -norm, where the error of the whole system is just the envelope of the set of individual curves.

Let us illustrate this phenomenon with a picture. We consider a system of type (5.2),

$$(6.3) \quad u'_k = -10^k u_k + t^5, \quad k = 1, \dots, 4,$$

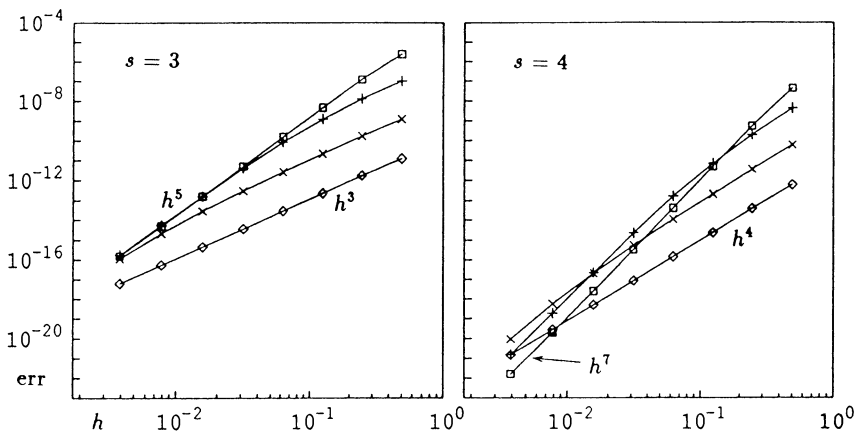


FIGURE 1. Global error of Radau IIA methods on (6.3) at  $t = 1$

with initial values on the smooth solution. In Figure 1 we have plotted the global error of the four components of (6.3) in dependence of  $h$  for the three- and four-stage Radau IIA methods. As  $\bar{\mu}_r = 0$  for all  $r$ , one observes the superposed orders 5 and 6 as the slope of the envelope of the four curves.

**6.2. Order reduction for semidiscretized PDE's.** Full discretization of a problem (2.1) gives rise to two types of errors: a space truncation error due to the discretization of the space variables  $x$  by finite elements or finite differences, and a time truncation error from the numerical integration of the resulting ODE by a Runge-Kutta method. The following analysis treats only the time truncation error.

Spatial discretization (by standard finite elements or finite differences) of the partial differential equation (2.1) leads to the ODE system

$$(6.4) \quad U' = L_N U + F_N(t),$$

where  $L_N$  is a constant  $N \times N$  matrix whose eigenvalues tend to the  $N$  first eigenvalues of the operator  $L(x, \partial)$  when  $N \rightarrow \infty$ . The global error of Runge-Kutta methods is governed by a similar superposition as described above. But as there are just  $N$  components, the superposition takes place only for  $h$  sufficiently large. For  $h \rightarrow 0$ , we observe, of course, classical order of convergence. There exists thus an  $h$ -zone where the order result of Theorem 2 applies. This zone becomes arbitrarily large when  $N$  tends to infinity.

We illustrate this superposition with a numerical example. Consider the PDE (3.8) with  $a(x) = x$ ,  $g(t) = e^{-t}$  and homogeneous Dirichlet boundary conditions. We discretize by standard 3-point finite differences. In this case the matrix  $L_N$  is given by (5.12) and the parameter  $\bar{\nu}$  of Theorem 2 is  $\bar{\nu} = 1/4$ . As the eigenvectors of (5.12) are orthogonal, system (6.4) can be decoupled by an orthogonal transformation into the diagonal system

$$(6.5) \quad V' = \Lambda V + \tilde{F}(t),$$

where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$  with the eigenvalues  $\lambda_k$  of (5.12). This decoupling is the perfect numerical analogue of the continuous decoupling of (3.8) with the eigenfunctions  $\varphi_k(x) = \sin(k\pi x)$ .



TABLE 1. Observed orders of convergence

$n$	$\alpha_{1,ob}$	$\alpha_{2,ob}$	$\alpha_{\infty,ob}$	$\hat{\alpha}_{1,ob}$	$\hat{\alpha}_{2,ob}$	$\hat{\alpha}_{\infty,ob}$
2	6.45	6.36	6.19	6.05	6.36	6.50
4	6.61	6.37	6.11	6.05	6.37	6.79
8	6.57	6.32	6.06	6.04	6.32	6.43
16	6.53	6.29	6.04	6.03	6.29	6.53
32	6.52	6.28	6.04	6.02	6.28	6.53
64	6.53	6.27	6.04	6.02	6.27	6.51
128	6.55	6.28	6.03	6.03	6.28	6.53

We denote by  $U_n$  and  $V_n$  the numerical solution of the four-stage Radau IIA method ( $p = 7, q = 4$ ) applied  $n$ -times with step size  $h = 1/n$  to the systems (6.4) and (6.5), respectively, with initial values on the smooth solution.  $U_n$  and  $V_n$  are approximations to the exact solutions  $U(t)$  and  $V(t)$  at  $t = 1$ . We computed the global errors  $e_n = U_n - U(1)$  and  $\hat{e}_n = V_n - V(1)$  in the three norms  $\|\cdot\|_1, \|\cdot\|_2$ , and  $\|\cdot\|_\infty$ . The observed orders of convergence are obtained through the formulas

$$\alpha_{i,ob} = \log_2(\|e_n\|_i / \|e_{2n}\|_i), \quad i = 1, 2, \infty,$$

and

$$\hat{\alpha}_{i,ob} = \log_2(\|\hat{e}_n\|_i / \|\hat{e}_{2n}\|_i), \quad i = 1, 2, \infty.$$

We display these values in Table 1 for  $N = 50$  and  $n = 2, 4, 8, \dots, 128$ .

Table 1 nicely shows that the observed orders of convergence correspond to the theoretical values given in Theorems 2 and 6, in particular, the order reduction  $\alpha_2 = 6.25$  predicted by (3.12) and  $\hat{\alpha}_1 = 6, \hat{\alpha}_\infty = 6.5$  predicted by (5.10), can be observed. Notice also the identities  $\alpha_1 = \hat{\alpha}_\infty$  and  $\alpha_\infty = \hat{\alpha}_1$ , which follow from (3.19) and (5.10).

ACKNOWLEDGMENT

We would like to thank M. Crouzeix, E. Hairer, G. Wanner, J.P. Kauthen, and P.A. Raviart for several observations which improved the presentation of the paper considerably. Further, we thank an anonymous referee for valuable suggestions.

BIBLIOGRAPHY

1. P. Brenner, M. Crouzeix, and V. Thomée, *Single step methods for inhomogeneous linear differential equations in Banach space*, RAIRO Anal. Numér. **16** (1982), 5–26.
2. P. Brenner, V. Thomée, and L. B. Wahlbin, *Besov spaces and applications to difference methods for initial value problems*, Lecture Notes in Math., vol. 434, Springer-Verlag, Berlin-Heidelberg, 1975.
3. K. Burrage and W. H. Hundsdorfer, *The order of B-convergence of algebraically stable Runge-Kutta methods*, BIT **27** (1987), 62–71.
4. K. Burrage, W. H. Hundsdorfer, and J. G. Verwer, *A study of B-convergence of Runge-Kutta methods*, Computing **36** (1986), 17–34.

5. J. C. Butcher, *The numerical analysis of ordinary differential equations: Runge-Kutta and general-linear methods*, Wiley, Chichester, 1987.
6. M. Crouzeix, *Sur l'approximation des équations différentielles opérationnelles linéaires par des méthodes de Runge-Kutta*, Thèse d'Etat, Université Paris VI, 1975.
7. M. Crouzeix and P. A. Raviart, *Méthodes de Runge-Kutta*, Unpublished Lecture Notes, Université de Rennes, 1980.
8. R. Frank, J. Schneid, and C. W. Ueberhuber, *The concept of B-convergence*, SIAM J. Numer. Anal. **18** (1981), 753–780.
9. E. Hairer, Ch. Lubich, and M. Roche, *Error of Runge-Kutta methods for stiff problems studied via differential algebraic equations*, BIT **28** (1988), 678–700.
10. E. Hairer, S. P. Nørsett, and G. Wanner, *Solving ordinary differential equations I. Nonstiff problems*, Springer-Verlag, Berlin-Heidelberg, 1987.
11. E. Hairer and G. Wanner, *Solving ordinary differential equations II, Stiff and differential-algebraic problems*, Springer-Verlag, Berlin-Heidelberg, 1991.
12. M. N. Le Roux, *Semidiscretization in time for parabolic problems*, Math. Comp. **33** (1979), 919–931.
13. ———, *Méthodes multipas pour des équations paraboliques non linéaires*, Numer. Math. **35** (1980), 143–162.
14. Ch. Lubich, *On the convergence of multistep methods for nonlinear stiff differential equations*, Numer. Math. **58** (1991), 839–853, and Corrigendum, *ibid.* **61** (1992), 277–279.
15. A. Ostermann and M. Roche, *Rosenbrock methods for partial differential equations and fractional orders of convergence*, Submitted for publication.
16. A. Prothero and A. Robinson, *On the stability and accuracy of one-step methods for solving stiff systems of ordinary differential equations*, Math. Comp. **28** (1974), 145–162.
17. J. M. Sanz-Serna and J. G. Verwer, *Stability and convergence at the PDE/Stiff ODE interface*, Appl. Numer. Math. **5** (1989), 117–132.
18. J. M. Sanz-Serna, J. G. Verwer, and W. H. Hundsdorfer, *Convergence and order reduction of Runge-Kutta schemes applied to evolutionary problems in partial differential equations*, Numer. Math. **50** (1987), 405–418.
19. S. Scholz, *Order barriers for the B-convergence of SDIRK methods*, Preprint, TU Dresden, 1987.
20. H. Triebel, *Interpolation theory, function spaces, differential operators*, North-Holland, Amsterdam, 1978.
21. J. G. Verwer, *Convergence and order reduction of diagonally implicit Runge-Kutta schemes in the method of lines*, Numerical Analysis (D. F. Griffiths and G. A. Watson, eds.), Pitman Research Notes in Math., vol. 140, 1986, pp. 220–237.
22. R. M. Young, *An introduction to nonharmonic Fourier series*, Academic Press, New York, 1980.

UNIVERSITÄT INNSBRUCK, INSTITUT FÜR MATHEMATIK UND GEOMETRIE, TECHNIKERSTRASSE  
13, A-6020 INNSBRUCK, AUSTRIA  
E-mail address: alex@mat0.uibk.ac.at

CRAY RESEARCH (SUISSE), CENTRE DE CALCUL EPFL, BÂTIMENT MA, CH-1015 LAUSANNE,  
SWITZERLAND  
E-mail address: roche@craysun1.epfl.ch