

## IMPLICIT RUNGE-KUTTA METHODS FOR HIGHER INDEX DIFFERENTIAL-ALGEBRAIC SYSTEMS

E. HAIRER and L. JAY

Section de Mathématiques, Université de Genève  
Case postale 240, CH-1211 Genève 24, Switzerland

### ABSTRACT

This article considers the numerical treatment of differential-algebraic systems by implicit Runge-Kutta methods. The perturbation index of a problem is discussed and its relation to the numerical solution is explained. Optimal convergence results of implicit Runge-Kutta methods for problems of index 1, 2, and 3 in Hessenberg form are then surveyed and completed. Their importance in the study of convergence for singular perturbation problems is shown and some comments on the numerical treatment of stiff Hamiltonian systems are given.

### 1. Introduction

The subject of this paper is the numerical treatment of nonlinear differential-algebraic equations (DAEs) of the form

$$0 = F(u, v) \quad (1)$$

where  $u$  and  $v$  are of the same dimension. The matrix  $\partial F/\partial v$  may be singular but is assumed to have constant rank. An important special case is the situation where the components can be separated in differential and algebraic parts as follows

$$y' = f(y, z), \quad 0 = g(y, z). \quad (2)$$

Differential-algebraic equations arise in a variety of applications, e.g., constrained mechanical systems, robotics, simulation of electrical networks and control engineering. They are also obtained as the limit of singular perturbation problems.

There are several ways for solving numerically the above problem. All of them have their own advantages:

- *Index reduction with projection.* Differentiate analytically the algebraic constraints and do some algebraic manipulations until an ordinary differential equation (ODE) is obtained. This ODE can then be solved by any ODE method (explicit or implicit, one-step or multistep). In order to avoid a "drift off" from the original algebraic constraints, it is recommended to combine this approach with certain projections onto the manifolds where the exact solution lies.

- *Direct approach.* Embed the original problem into a singular perturbation problem (e.g.,  $y' = f(y, z)$ ,  $\epsilon z' = g(y, z)$  for Eq. 2), apply formally an ODE method and

consider the limit  $\epsilon \rightarrow 0$ . This approach is restricted to implicit methods whose stability function is bounded by one at infinity. However, it provides much insight into the numerical solution of stiff and singular perturbation problems.

- *Special methods.* They are adapted to the particular problem. Usually one applies some explicit ODE method to the differential part of the DAE and some nonlinear equation solver to the algebraic part.

In this article we restrict ourselves to the direct approach combined with the use of implicit Runge-Kutta methods. For further results on the numerical treatment of DAEs we refer the reader to the monographs of Griepentrog and März<sup>5</sup>, Brenan, Campbell and Petzold<sup>2</sup>, and Hairer, Lubich and Roche<sup>7</sup>.

Consider the DAE of Eq. 1 and assume that consistent initial values are given ( $u_0$  is consistent if a function  $u(t)$  exists which satisfies Eq. 1 and  $u(t_0) = u_0$ ). For an  $s$ -stage implicit Runge-Kutta method with coefficients  $a_j, b_j$ , the direct approach yields<sup>8,12</sup>

$$U_i = u_0 + h \sum_{j=1}^s a_{ij} U_j^i, \quad 0 = F(U_i, U_i^i),$$

$$u_1 = u_0 + h \sum_{i=1}^s b_i U_i^1. \tag{3}$$

The first line of Eq. 3 represents a nonlinear system for  $U_i, U_i^i, i = 1, \dots, s$ , and  $u_1$  is the numerical approximation for the solution at  $t_0 + h$ . Due to the wide variety of problems included in the formulation of Eq. 1, there is no hope for a unified convergence theory. There exist perfectly meaningful DAEs which cause difficulties to every numerical method. Fortunately, the problems arising in practice have some additional structures and permit a successful application of the above method.

This paper is organized in the following way. We begin with a classification of DAEs (perturbation index) which measures how strong the problem is ill-conditioned. For several important problems in Hessenberg form (of index 1, 2, and 3) we then present optimal convergence results for implicit Runge-Kutta methods. These are used to give some insight into the numerical solution of singular perturbation problems. As an example, the numerical treatment of stiff Hamiltonian systems is discussed.

**2. Influence of Perturbations**

For ordinary differential equations  $u' = f(u)$  it follows from the lemma of Gronwall that the difference between the exact solution  $u(t)$  and a perturbed solution  $\tilde{u}(t)$  with defect  $\delta(t) := \tilde{u}'(t) - f(\tilde{u}(t))$  can be estimated as

$$\|\tilde{u}(t) - u(t)\| \leq e^{L(t-t_0)} (\|\tilde{u}(t_0) - u(t_0)\| + \max_{t_0 \leq \tau \leq t} \|\int_{t_0}^{\tau} \delta(\tau) d\tau\|). \tag{4}$$

Working on a bounded interval ( $t_0 \leq t \leq T$ ) this means that small perturbations in the data of the problem cause small perturbations in the solution. For differential-algebraic equations the situation may be completely different.

*Index 1 problems.* We consider the DAE of Eq. 2 together with a perturbed version

$$y' = f(y, z), \quad \tilde{y}' = f(\tilde{y}, \tilde{z}) + \delta_1(t),$$

$$0 = g(y, z), \quad 0 = g(\tilde{y}, \tilde{z}) + \delta_2(t), \tag{5}$$

and assume that

$$g_z \text{ is invertible} \tag{6}$$

in a neighbourhood of the solution ( $g_z$  denotes the derivative of  $g$  with respect to  $z$ ). Using the implicit function theorem we can solve the algebraic relations in Eq. 5 for  $z$  (resp.  $\tilde{z}$ ). Inserted into the differential equation of Eq. 5 this yields an ODE for  $y$  (resp.  $\tilde{y}$ ) and the estimate of Eq. 4 can be used to obtain

$$\|\tilde{y}(t) - y(t)\| \leq C (\|\tilde{y}(t_0) - y(t_0)\| + \int_{t_0}^t \|g_z(s)\| ds + \max_{t_0 \leq s \leq t} \|\int_{t_0}^s \delta_1(\tau) d\tau\|),$$

$$\|\tilde{z}(t) - z(t)\| \leq C (\|\tilde{y}(t) - y(t)\| + \|\delta_2(t)\|).$$

Again, the problem is in general well-conditioned.

*Index 2 problems.* An example where the condition of Eq. 6 is violated is

$$y' = f(y, z), \quad \tilde{y}' = f(\tilde{y}, \tilde{z}) + \delta(t),$$

$$0 = g(y), \quad 0 = g(\tilde{y}) + \theta(t). \tag{7}$$

It represents a typical control problem where  $z$  acts as a control variable and forces the solution  $y$  to stay on the manifold defined by  $0 = g(y)$ . The essential idea is to differentiate the algebraic constraints with respect to  $t$ . This yields  $0 = g_y(y)f(y, z)$  and a similar relation for the perturbed system. If we assume that

$$g_y f_z \text{ is invertible} \tag{8}$$

in a neighbourhood of the solution, the differentiated constraint can be solved for  $z$ . We thus obtain, as for the index 1 example, the estimates<sup>7</sup>

$$\|\tilde{y}(t) - y(t)\| \leq C (\|\tilde{y}(t_0) - y(t_0)\| + \int_{t_0}^t (\|\delta(s)\| + \|\theta'(s)\|) ds),$$

$$\|\tilde{z}(t) - z(t)\| \leq C (\|\tilde{y}(t_0) - y(t_0)\| + \max_{t_0 \leq s \leq t} \|\delta(s)\| + \max_{t_0 \leq \tau \leq t} \|\theta'(\tau)\|). \tag{9}$$

Although the estimate for the  $y$ -component can be slightly improved<sup>1</sup>, the dependence on  $\theta'(t)$  cannot be suppressed in general. Therefore, small perturbations in Eq. 8 may lead to large perturbations in the solution.

*Index 3 problems.* The equations of motion of constrained mechanical systems can be written in the form

$$y' = f(y, z), \quad \tilde{y}' = f(\tilde{y}, \tilde{z}) + \delta(t),$$

$$\begin{aligned} z' &= k(y, z, u), & z' &= k(\bar{y}, \bar{z}, \bar{u}) + \mu(t), \\ 0 &= g(y), & 0 &= g(\bar{y}) + \theta(t), \end{aligned} \tag{10}$$

where typically  $k$  depends linearly on  $u$ . If we differentiate twice the algebraic constraint, the assumption

$$g_y f z k_u \text{ is invertible} \tag{11}$$

allows to express  $u$  in terms of  $y, z$ , and estimates for  $\bar{y}(t) - y(t), \dots$  can be obtained. These estimates will depend on  $\theta''(t), \theta'(t)$ , and on  $\delta'(t)$ , so that this problem is even more ill-conditioned than the previous one. However, in some important situations (e.g.,  $f(y, z) = f_0(y) + f_1(y)z, k(y, z, u) = k_0(y, z) + k_1(y)u$ ) the differences  $\bar{y}(t) - y(t), \bar{z}(t) - z(t)$  (but not  $\bar{u}(t) - u(t)$ ) are independent of  $\theta''(t)$  and  $\delta'(t)$ .

These examples motivate the following definition of the index.

**Definition.** Eq. 1 has *perturbation index*  $m$  along a solution  $u(t)$  on  $[t_0, T]$ , if  $m$  is the smallest integer such that, for all functions  $\bar{u}(t)$  having a defect

$$F(\bar{u}(t), \bar{u}(t)) = \delta(t),$$

there exists on  $[t_0, T]$  an estimate

$$\|\bar{u}(t) - u(t)\| \leq C(\|\bar{u}(t_0) - u(t_0)\| + \max_{t_0 \leq s \leq t} \|\delta(s)\| + \dots + \max_{t_0 \leq s \leq t} \|\delta^{(m-1)}(s)\|)$$

whenever the expression on the right-hand side is sufficiently small.

It is also of interest to study the influence of perturbations in the Runge-Kutta equations (Eq. 3) on the numerical solution. Let us explain this at the example of the index 2 problem (Eq. 7). The internal stages of the Runge-Kutta method satisfy

$$\begin{aligned} Y_{ni} &= y_n + h \sum_{j=1}^s a_{ij} Y_{nj}, & Y'_{ni} &= f(Y_{ni}, Z_{ni}), \\ Z_{ni} &= z_n + h \sum_{j=1}^s a_{ij} Z'_{nj}, & 0 &= g(Y_{ni}), \end{aligned}$$

and the numerical approximation at  $t_{n+1} = t_0 + (n+1)h$  is given by

$$y_{n+1} = y_n + h \sum_{i=1}^s b_i Y'_{ni}, \quad z_{n+1} = z_n + h \sum_{i=1}^s b_i Z'_{ni}.$$

If we eliminate the variables  $Y'_{ni}, Z'_{ni}$ , we obtain the equivalent formulas

$$Y_{ni} = y_n + h \sum_{j=1}^s a_{ij} f(Y_{nj}, Z_{nj}), \quad 0 = g(Y_{ni}), \tag{12a}$$

$$y_{n+1} = y_n + h \sum_{i=1}^s b_i f(Y_{ni}, Z_{ni}), \quad z_{n+1} = \rho z_n + \sum_{i,j=1}^s b_i \omega_{ij} Z_{nj}, \tag{12b}$$

where

$$(\omega_{ij}) = (a_{ij})^{-1} \quad \text{and} \quad \rho = 1 - \sum_{i,j=1}^s b_i \omega_{ij} \tag{13}$$

(of course, we have assumed that  $A = (a_{ij})$  is invertible;  $\rho$  is the value of the stability function at infinity). Eq. 12a represents a nonlinear system for  $Y_{ni}, Z_{ni}, i = 1, \dots, s$ , and the numerical solution after one step is then explicitly given by Eq. 12b.

We next compare Eq. 12 to the perturbed Runge-Kutta method

$$\bar{Y}_{ni} = \bar{y}_n + h \sum_{j=1}^s a_{ij} f(\bar{Y}_{nj}, \bar{Z}_{nj}) + h \delta_{ni}, \quad 0 = g(\bar{Y}_{ni}) + \theta_{ni}, \tag{14a}$$

$$\bar{y}_{n+1} = \bar{y}_n + h \sum_{i=1}^s b_i f(\bar{Y}_{ni}, \bar{Z}_{ni}) + h \delta_{n,s+1}, \quad \bar{z}_{n+1} = \rho \bar{z}_n + \sum_{i,j=1}^s b_i \omega_{ij} \bar{Z}_{nj} \tag{14b}$$

with the aim of estimating the differences  $\Delta y_n = \bar{y}_n - y_n$  and  $\Delta z_n = \bar{z}_n - z_n$ . Since the nonlinear systems do not depend on the initial values of the  $z$ -component, it is possible to derive the estimates for  $\Delta y_n$  independently of those for  $\Delta z_n$ .

Subtracting Eq. 12 from Eq. 14 we obtain by linearization that<sup>7</sup>

$$\Delta y_{n+1} = P_n \Delta y_n + \rho Q_n \Delta y_n + \mathcal{O}(h \|\Delta y_n\| + h \delta_n + \theta_n) \tag{15}$$

where

$$P_n = (f_z(g_y f z)^{-1} g_y)(y_n, z_n), \quad Q_n = I - P_n$$

are suitable projectors and

$$\delta_n = \max_{i=1, \dots, s+1} \|\delta_{ni}\|, \quad \theta_n = \max_{i=1, \dots, s} \|\theta_{ni}\|$$

(we have implicitly used the fact that  $h, \delta_n, \theta_n$  and  $\Delta y_n$  are sufficiently small and that  $y_n$  is sufficiently close to consistent initial values, so that the numerical solution exists). We next assume that  $|\rho| \leq 1$  and deduce from Eq. 15 that

$$\|\Delta y_N\| \leq C(\|P_0 \Delta y_0\| + (|\rho|^N + h) \|Q_0 \Delta y_0\| + h \sum_{n=0}^{N-1} (\delta_n + \frac{\theta_n}{h})) \tag{16a}$$

for  $Nh \leq \text{Const}$ . Similarly one gets for the  $z$ -component

$$\|\Delta z_N\| \leq C(\|P_0 \Delta y_0\| + (|\rho|^N + h) \|Q_0 \Delta y_0\| + \max_{n=0, \dots, N-1} \delta_n + \max_{n=0, \dots, N-1} \frac{\theta_n}{h}). \tag{16b}$$

Eq. 16 is the numerical analogue of Eq. 9. We observe that the derivative in Eq. 9 becomes a division by the stepsize  $h$  in Eq. 16.

These estimates can be exploited in several ways. If one replaces  $\bar{y}_n, \bar{z}_n$  by values on the exact solution, one can obtain convergence results of the method. Furthermore, it is possible to interpret the perturbations as round-off errors or errors in the iterative solution of the nonlinear systems. Eq. 16 shows how such errors can affect the numerical solution. The ill-conditioning of the problem is

reflected by the factor  $1/h$  in front of the perturbations  $\theta_n$ . Therefore, one has to take care of the above-mentioned errors when the stepsize is very small.

Estimates similar to Eq. 16 can be derived also for the index 1 and index 3 problems of the beginning of this section. However, similar results for general DAEs, having perturbation index  $m$ , are not yet known.

### 3. Convergence Results

The above investigations make it clear that there is no unified convergence theory of Runge-Kutta methods for general DAEs, and that one has to treat separately all different types of problems. The first results (for linear inhomogeneous DAEs with constant coefficients, but of arbitrary index) have been obtained by Petzold<sup>12</sup>. In this section we collect convergence results for nonlinear semi-explicit problems of index 1, 2, and 3. For nearly all Runge-Kutta methods the estimates are optimal in the sense that the exponent of  $h$  cannot be improved without restricting the class of problems.

For the presentation of the convergence results we need the following abbreviations:

$$B(p) : \sum_{i=1}^p b_i c_i^{k-1} = \frac{1}{k} \quad \text{for } k = 1, \dots, p,$$

$$C(q) : \sum_{j=1}^q a_j c_j^{k-1} = \frac{c_k^k}{k} \quad \text{for } i = 1, \dots, s; \quad k = 1, \dots, q;$$

$$D(r) : \sum_{i=1}^r b_i c_i^{k-1} a_{ij} = \frac{b_i}{k} (1 - c_i^k) \quad \text{for } j = 1, \dots, s; \quad k = 1, \dots, r;$$

$$(S) : \quad a_{si} = b_i \quad \text{for } i = 1, \dots, s.$$

The assumptions  $B(p)$ ,  $C(q)$ ,  $D(r)$  have been introduced by Butcher<sup>4</sup> and play a crucial role in the construction of implicit Runge-Kutta methods<sup>4,8,9</sup>. The condition (S) means that the method is "stiffly accurate". Throughout this section we shall assume that the Runge-Kutta matrix  $A = (a_{ij})$  is invertible and we shall denote by  $\rho$  the value defined in Eq. 13. The integer  $q$  of condition  $C(q)$  is called the "stage order". We always assume that  $p \geq q$ .

**Theorem (index 1).** Consider the index 1 problem (Eq. 5, Eq. 6) and assume that the initial values are consistent. If the Runge-Kutta method satisfies  $B(p)$ ,  $C(q)$ ,  $D(r)$ , and  $|\rho| \leq 1$ , then the global error satisfies for  $hn \leq \text{Const}$

$$y(t_n) - y_n = \mathcal{O}(h^\eta), \quad z(t_n) - z_n = \mathcal{O}(h^\zeta),$$

where  $\eta = \min(p, 2q + 2, q + r + 1)$  and

$$\zeta = \begin{cases} \eta, & \text{if (S) holds,} \\ \min(p, q + 1) & \text{if } -1 \leq \rho < 1, \\ \min(p - 1, q) & \text{if } \rho = 1. \end{cases}$$

The proof of this theorem is based on the following idea: due to Eq. 6, the algebraic constraint  $0 = g(y, z)$  can formally be written as  $z = G(y)$ . Inserted into the differential equation of Eq. 5 we get the ODE  $y' = f(y, G(y))$ . Convergence for the  $y$ -component now follows from the fact that the numerical solution, obtained from Eq. 3, is identical to the approximation of the Runge-Kutta method applied to  $y' = f(y, G(y))$ . Hence the results from the ODE theory can be applied<sup>7,8</sup>.

The above estimates are valid for a constant stepsize application of the method. For variable stepsizes the same estimates hold with  $h = \max_n h_n$  (with the exception of the case  $\rho = -1$  where the results become those of the case  $\rho = +1$ ). The order reduction in the  $z$ -component can be avoided, if one computes the approximation  $\tilde{z}_n$  from  $g(y_n, \tilde{z}_n) = 0$ . Similar remarks can be made also for the subsequent theorems.

**Theorem (index 2).** Consider the index 2 problem (Eq. 7, Eq. 8) and assume that the initial values are consistent. If the Runge-Kutta method satisfies  $B(p)$ ,  $C(q)$ ,  $D(r)$ , and  $|\rho| \leq 1$ , then the global error satisfies for  $hn \leq \text{Const}$

$$y(t_n) - y_n = \mathcal{O}(h^\eta), \quad z(t_n) - z_n = \mathcal{O}(h^\zeta),$$

where

$$\eta = \begin{cases} \min(p, 2q, q + r + 1) & \text{if (S) holds,} \\ \min(p, q + 1) & \text{if } -1 \leq \rho < 1, \\ q & \text{if } \rho = 1, \end{cases}$$

$$\zeta = \begin{cases} q & \text{if } |\rho| < 1, \\ q - 1 & \text{if } \rho = -1, \\ q - 2 & \text{if } \rho = 1. \end{cases}$$

The proof of this theorem needs the study of the local error (using rooted trees, elementary differentials, ...) and of the error propagation (see Eq. 15). Details are given in Hairer, Lubich and Roche<sup>7</sup>.

Optimal convergence results for the index 3 problem have been obtained only very recently<sup>10</sup> under the assumption (S). In view of the application to Hamiltonian systems (section 4) we also include new results for the case  $|\rho| = 1$ .

**Theorem (index 3).** Consider the index 3 problem (Eq. 10, Eq. 11) and assume that the initial values are consistent. If the Runge-Kutta method satisfies  $B(p)$ ,  $C(q)$ ,  $D(r)$ , and  $|\rho| \leq 1$ , then the global error satisfies for  $hn \leq \text{Const}$

$$y(t_n) - y_n = \mathcal{O}(h^\eta), \quad z(t_n) - z_n = \mathcal{O}(h^\zeta), \quad u(t_n) - u_n = \mathcal{O}(h^\nu),$$

where

$$\eta = \begin{cases} \min(p, 2q - 2, q + r) & \text{if (S) holds,} \\ \min(p, q + 1) & \text{if } -1 \leq \rho < 1, \tau \geq 1, q \geq 3, \\ q & \text{else,} \end{cases}$$

$$\zeta = \begin{cases} q & \text{if } |\rho| < 1, \\ q - 1 & \text{if } \rho = -1, \\ q - 2 & \text{if } \rho = 1, \end{cases}$$

$$\nu = \begin{cases} q - 1 & \text{if } |\rho| < 1, \\ q - 3 & \text{if } \rho = -1, \\ q - 4 & \text{if } \rho = 1. \end{cases}$$

In the above two theorems the stage order  $q$  has to be sufficiently large such that the numerical solution remains close to the exact solution. Otherwise the nonlinear system may not have a solution.

4. Singular Perturbation Problems

As mentioned in the introduction, the direct approach for solving DAEs provides much insight into the numerical solution of singular perturbation problems. Let us illustrate this at two examples.

Consider first the problem

$$\begin{aligned} y' &= f(y, z), \\ \epsilon z' &= g(y, z), \end{aligned} \quad 0 < \epsilon \ll 1 \tag{17}$$

where  $f$  and  $g$  are sufficiently differentiable. Under suitable assumptions on  $g$  (e.g.,  $g_z(y, z)v \leq -\|v\|^2$ ) and on the initial values, the solution of Eq. 17 possesses an asymptotic expansion of the form

$$y(t) = y_0(t) + \epsilon y_1(t) + \epsilon^2 y_2(t) + \dots, \quad z(t) = z_0(t) + \epsilon z_1(t) + \epsilon^2 z_2(t) + \dots$$

Inserting these expansions into Eq. 17 and collecting equal powers of  $\epsilon$  we obtain

$$\begin{aligned} y_0' &= f(y_0, z_0) \\ 0 &= g(y_0, z_0) \end{aligned} \tag{18a}$$

$$\begin{aligned} y_1' &= f_y(y_0, z_0)y_1 + f_z(y_0, z_0)z_1 \\ z_0' &= g_y(y_0, z_0)y_1 + g_z(y_0, z_0)z_1. \end{aligned} \tag{18b}$$

We see that Eq. 18a constitutes an index 1 problem (Eq. 5) for  $y_0(t), z_0(t)$ . The Eqs. 18a and 18b together are an index 2 problem (Eq. 7) with  $(y_0, z_0, y_1)$  in the role of  $y$  and  $z_1$  in the role of  $z$ .

The main point is now that the numerical solution of a Runge-Kutta method applied to Eq. 17 also has an expansion of the form<sup>6,9</sup>

$$y_n = y_n^0 + \epsilon y_n^1 + \epsilon^2 y_n^2 + \dots, \quad z_n = z_n^0 + \epsilon z_n^1 + \epsilon^2 z_n^2 + \dots,$$

and the coefficients  $y_n^0, z_n^0, \dots$  are exactly the numerical solution of the Runge-Kutta method applied to Eq. 18. Consequently, the convergence results for the index 1 problem yield an estimate for  $y_0(t_n) - y_n^0, z_0(t_n) - z_n^0$  those for the index 2 problem yield an estimate for  $y_1(t_n) - y_n^1, z_1(t_n) - z_n^1$ , etc. Since the higher order terms in the asymptotic expansions can be neglected<sup>6,9</sup>, this leads to sharp convergence statements for the singular perturbation problem of Eq. 17.

As a second example we consider the problem<sup>11</sup>

$$\begin{aligned} q_1' &= p_1, & p_1' &= -\frac{1}{\epsilon^2} \frac{q_1}{\sqrt{q_1^2 + q_2^2}} (\sqrt{q_1^2 + q_2^2} - 1), \\ q_2' &= p_2, & p_2' &= -\frac{1}{\epsilon^2} \frac{q_2}{\sqrt{q_1^2 + q_2^2}} (\sqrt{q_1^2 + q_2^2} - 1), \end{aligned} \tag{19}$$

which describes a stiff spring pendulum (mass point suspended at a massless spring with Hooke's constant  $1/\epsilon^2, 0 < \epsilon \ll 1$ ). If we introduce the new variable  $\lambda$  by

$$\epsilon^2 \lambda = \frac{\sqrt{q_1^2 + q_2^2} - 1}{\sqrt{q_1^2 + q_2^2}},$$

Eq. 19 becomes

$$\begin{aligned} q_1' &= p_1, & p_1' &= -q_1 \lambda, \\ q_2' &= p_2, & p_2' &= -q_2 \lambda - 1. \end{aligned}$$

The so obtained system is a DAE of index 1 if  $\epsilon > 0$ , and it is of index 3 for the limit case  $\epsilon = 0$  ( $(q_1, q_2)$  corresponds to  $y, (p_1, p_2)$  to  $z$ , and  $\lambda$  to  $u$  in Eq. 10). Since the Runge-Kutta method (direct approach) is invariant under the above transformation, we expect that for small  $\epsilon$  the numerical method behaves similarly as for an index 3 problem. A rigorous analysis of this fact in a more general context has been presented by Lubich<sup>11</sup>.

The following numerical experiment illustrates this behaviour for the 3-stage methods GAUSS and RADAU IIA, whose coefficients satisfy the following conditions:

$$\begin{aligned} \text{GAUSS} & \quad B(2s), C(s), D(s), \rho = (-1)^s, \\ \text{RADAU IIA} & \quad B(2s-1), C(s), D(s-1), (S), \rho = 0. \end{aligned}$$

We put  $\epsilon = 0.001$ , take the initial values<sup>7</sup>

$$\begin{aligned} q_1(0) &= 1 - 3\epsilon^4 + \mathcal{O}(\epsilon^8), & q_2(0) &= 0, \\ p_1(0) &= \mathcal{O}(\epsilon^8), & p_2(0) &= 0, \end{aligned} \tag{20}$$

and integrate Eq. 19 on the interval  $[0, 0.5]$  with several different constant stepsize. The initial values are chosen such that the solution of Eq. 19 does not contain highly oscillatory terms (smooth motion<sup>11</sup>). This allows the use of stepsize which are significantly larger than  $\epsilon$ .

Fig. 1 shows the global errors as a function of  $h$ . Since we have used a double logarithmic scale, a function  $\mathcal{O}h^r$  appears as a straight line with slope  $r$ . We see that for the RADAU IIA method the errors behave like  $\mathcal{O}(h^2), \mathcal{O}(h^3), \mathcal{O}(h^5)$  for the components  $\lambda, p_1, q_2$ , respectively. For the GAUSS method they behave like  $\mathcal{O}(h^0), \mathcal{O}(h^2)$  for the components  $\lambda, p_2$  (at least for sufficiently large  $h$ ). The error of the position coordinate  $q_2$  oscillates around a line of slope 4, indicating a  $\mathcal{O}(h^4)$  behaviour. The errors for  $q_1, p_1$  behave similarly to those for  $q_2, p_2$  and are not plotted. This experiment confirms the results predicted by the theorem (index 3) of the previous section. As a consequence, methods satisfying condition (S) (like RADAU IIA) are preferred. However, for a long time integration the situation may be different.

Eq. 19 represents a Hamiltonian system

$$\begin{aligned} q' &= \frac{\partial H}{\partial p}(p, q), & p' &= -\frac{\partial H}{\partial q}(p, q) \end{aligned} \tag{21}$$

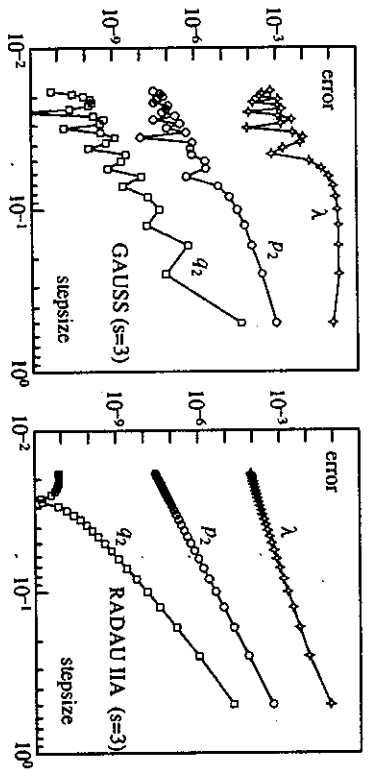


Fig. 1. Global error as a function of the stepsize for Eq. 19,  $h = 0.5/n$ ,  $n = 1, 2, 3, \dots$

with Hamiltonian function

$$H(p, q) = \frac{1}{2}(p_1^2 + p_2^2) + \frac{1}{2s^2}(\sqrt{q_1^2 + q_2^2} - 1)^2 + q_2.$$

Recently, much research has been devoted to the numerical integration of such systems (see, for example, the survey article by Sanz-Serna<sup>13</sup>). In order to retain the qualitative properties of the flow of Eq. 21, it is important for the numerical scheme to be symplectic. For implicit Runge-Kutta methods this means that the coefficients have to satisfy<sup>13</sup>

$$b_i a_j + b_j a_i - b_i b_j = 0 \quad \text{for all } i, j. \quad (22)$$

It is known that the GAUSS methods satisfy Eq. 22, whereas the RADAU IIA methods do not. One can prove (multiply Eq. 22 by  $\omega_i \omega_j$  and sum over all indices) that Eq. 22 implies  $|\rho| = 1$ , a rather undesirable property for the integration of index 3 problems.

In our second experiment we integrate Eq. 19 ( $\epsilon = 0.001$ ) with the initial values of Eq. 20 and with constant stepsize  $h = 0.05$  over a long interval [0, 26] (about 3 periods). The Hamiltonian function for the numerical solution, obtained by the GAUSS and RADAU IIA methods with  $s = 3$ , is plotted in Fig. 2 (for the exact solution the Hamiltonian is constant and equals  $4.5\epsilon^6 + O(\epsilon^8)$ ). We observe that it remains between tolerable bounds for the GAUSS method, but drifts away from the exact value for the RADAU IIA method. This experiment demonstrates the different behaviour of symplectic and non symplectic integrators.

It should be mentioned that for "stiff" Hamiltonian systems, such as Eq. 19, the use of explicit integration methods is not recommended, because they require small stepsizes (usually not larger than  $\epsilon$ ). If one uses implicit symplectic integrators, such as GAUSS, the stage order  $q$  has to be sufficiently large, so that the numerical solution is precise enough (see index 3 theorem of section 3).

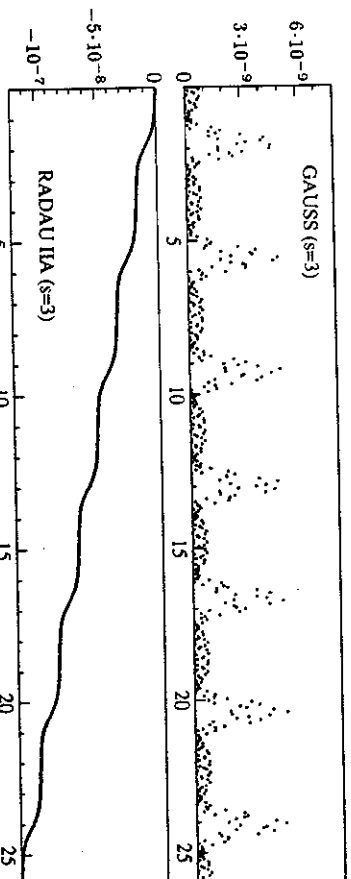


Fig. 2. Numerical Hamiltonian for Eq. 19

## 5. References

1. M. Arnold, Stability of numerical methods for differential-algebraic equations of higher index, *APNUM*, to appear.
2. K.E. Brenan, S.L. Campbell and L.R. Petzold, *Numerical Solution of Initial Value Problems in Differential-Algebraic Equations*, North-Holland, New York, 1989.
3. K.E. Brenan and L.R. Petzold, The numerical solution of higher index differential-algebraic equations by implicit Runge-Kutta methods, *SIAM J. Numer. Anal.* **26** (1989), 976-996.
4. J.C. Butcher, Coefficients for the study of Runge-Kutta integration processes, *J. Austral. Math. Soc.* **3** (1963), 185-201.
5. E. Griepentrog and R. März, *Differential-Algebraic Equations and Their Numerical Treatment*, Teubner-Texte zur Mathematik, Band 88, Leipzig, 1986.
6. 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.
7. E. Hairer, Ch. Lubich and M. Roche, *The Numerical Solution of Differential-Algebraic Systems by Runge-Kutta Methods*, Springer Lecture Notes in Mathematics 1409, Berlin, 1989.
8. E. Hairer, S.P. Nørsett and G. Wanner, *Solving Ordinary Differential Equations I. Nonsstiff problems*, Computational Mathematics 8, Springer-Verlag, Berlin, 1987.
9. E. Hairer and G. Wanner, *Solving Ordinary Differential Equations II. Stiff and Differential-Algebraic problems*, Computational Mathematics 14, Springer-Verlag, Berlin, 1991.

10. I. Jay, Convergence of Runge-Kutta methods for differential-algebraic systems of index 3, *Report, Sect. de mathématiques, Univ. de Genève, 1992*.
11. Ch. Lubich, Integration of stiff mechanical systems by Runge-Kutta methods, *ZAMP*, to appear.
12. L.R. Petzold, Order results for implicit Runge-Kutta methods applied to differential/algebraic systems, *SIAM J. Numer. Anal.* **23** (1986), 837-852.
13. J.M. Sanz-Serna, Symplectic integrators for Hamiltonian problems: an overview, *Acta Numerica* **1** (1992), 249-286.