

# Time-Step-Size-Independent Conditioning And Sensitivity To Perturbations in the Numerical Solution of Index-Three Differential Algebraic Equations\*

Carlo L. Bottasso

Dipartimento di Ingegneria Aerospaziale,  
Politecnico di Milano, Milano, Italy

Olivier A. Bauchau

D. Guggenheim School of Aerospace Engineering,  
Georgia Institute of Technology, Atlanta, GA, USA

Alberto Cardona

Cimec-Intec, Universidad Nacional del Litoral – Conicet,  
Güemes 3450, 3000 Santa Fe, Argentina

## Abstract

We propose a simple preconditioning for the equations of motion of constrained mechanical systems in index three form. The scaling transformation is applied to the displacement-velocity-multiplier and to the reduced displacement-multiplier forms. The analysis of the transformed system shows that conditioning and sensitivity to perturbations become independent of the time step size, as in the case of well behaved ordinary differential equations. The new scaling transformation is simple to implement and does not require the re-writing of the system equations as other approaches. The theoretical analysis is confirmed by numerical examples.

## 1 Introduction

It is well known that the amplification of small errors and perturbations in the solution of high index differential algebraic equations (DAEs) causes severe numerical difficulties. For example, Petzold and Lötstedt have shown in [1] that index three DAEs for constrained mechanical (multibody) systems are severely ill conditioned for small time step sizes when discretized using BDF-type formulas. Their analysis indicates that, unless corrective actions are taken, the condition number of the iteration matrix is  $\mathcal{O}(h^{-3})$ . Furthermore, errors propagate in the displacement, velocity and multiplier fields at rates which are shown to be  $\mathcal{O}(h^{-1})$ ,  $\mathcal{O}(h^{-2})$  and  $\mathcal{O}(h^{-3})$ , respectively. A related perturbation analysis that deals with the propagation of the error in the satisfaction of the constraint equations is described in Arnold [2], which also considers systems with friction. These results indicate that errors will grow very rapidly as the time step size is reduced, preventing in practice the use of time refinement procedures, and imposing a very tight tolerance on the solution of the non-linear discrete equations. Furthermore, the ill conditioning of the linear system does not favor the use of iterative solvers.

---

\**SIAM Journal on Scientific Computing*, **29**(1), pp 397–414, 2007

Several methods have been proposed in the literature to address these problems. For example, the system can be recast in index two form, by replacing the displacement level constraints with the velocity level ones. While the index two form is more robust to perturbations [2], this approach suffers from the well known drift effect, which calls for further ad-hoc corrective actions.

An alternative approach is to re-write the governing equations so as to include, together with the position level constraints, also their derivatives, as in the GGL [3] and Embedded Projection (EP) [4] methods. These methods however require additional multipliers which imply a somewhat increased problem size and, hence, a possibly higher computational cost.

Petzold and Lötstedt [1] discuss a simple scaling transformation of the index three governing equations which yields a condition number  $\mathcal{O}(h^{-2})$  and an improvement of one order in the errors for all solution fields. Although the sensitivity to perturbations is reduced with respect to the unscaled problem, difficulties can still be expected in practice.

In this paper we propose a solution to the conditioning and error propagation problems that is based on a left and right preconditioning. The left preconditioning amounts to a scaling of the dynamic equilibrium and constraint equations, while the right preconditioning amounts to a scaling of the unknowns. The proposed procedure is trivially implemented in an existing code and does not require the re-writing of the equations of motion as in the GGL and EP methods.

The new scaling transformation is applied first to the equations in the standard displacement-velocity-multiplier form, using a scaling of the sole Lagrange multipliers. A similar scaling was previously considered in Cardona and Geradin [5] in the context of the Hilbert, Hughes and Taylor (HHT) scheme. The analysis of the present paper shows that errors propagate in the displacement, velocity and multiplier fields at rates equal to  $\mathcal{O}(h^0)$ ,  $\mathcal{O}(h^{-1})$  and  $\mathcal{O}(h^0)$ , respectively, while the condition number is  $\mathcal{O}(h^{-1})$ . Next, it is shown that if also the velocities are scaled together with the Lagrange multipliers, one can achieve perfect time step size independence ( $\mathcal{O}(h^0)$ ) for error propagation and conditioning.

Finally, we analyze the displacement-multiplier form, obtained by static condensation of the velocities at each time step. This form of the equations is more computationally interesting than the full form for applications denoted by a large number of degrees of freedom, since it leads to the solution of smaller linear problems. Even in this case, the preconditioning of the reduced form leads to perfect time step size independence for both the conditioning and the error propagation in the displacement-multiplier solution variables.

These results indicate that index three DAEs are as easy to integrate as well behaved ordinary differential equations (ODEs), once they are recast in one of the  $\mathcal{O}(h^0)$  forms.

The paper is organized as follows. In Section 2 we formulate the equations of motion of multi-body systems in index three three-field form for BDF-type schemes, and we derive their sensitivity to perturbations and conditioning. Preconditioning by scaling of the discretized equations is presented in Section 3; the left preconditioning strategy is detailed in §3.1, while the left and right preconditioning is described in §3.2. The classical pendulum problem in §3.3 numerically confirms the results of the analysis and concludes this section. The two field form of the problem is discussed in Section 4; the left preconditioning is described in §4.2, the left and right preconditioning is given in §4.2, and the pendulum problem verifies the analysis in §4.3. Finally, a more complex contact-impact problem for a flexible multibody system is presented in Section 5 to show the importance of scaling in complex problems involving time step size refinement. The paper is closed by Section 6, which reports the main conclusions and findings of this work.

## 2 Sensitivity to perturbations and conditioning of multi-body system equations

The governing DAEs for multibody systems are written as follows

$$\mathbf{u}' = \mathbf{v}, \quad (1a)$$

$$\mathbf{v}' = \mathbf{f} + \mathbf{G}\boldsymbol{\lambda}, \quad (1b)$$

$$\boldsymbol{\phi} = \mathbf{0}, \quad (1c)$$

where equations (1a) are the kinematic equations, equations (1b) represent the equations of dynamic equilibrium,  $\mathbf{u}$  are generalized coordinates,  $\mathbf{v}$  generalized velocities,  $\boldsymbol{\lambda}$  are Lagrange multipliers that enforce the constraints (1c), and finally  $\mathbf{G}^T = \boldsymbol{\phi}_{,\mathbf{u}}$  is the constraint Jacobian. Following [1], we have considered for simplicity an identity mass matrix in equation (1b) without loss of generality.

The solution of equation (1) is sought by means of a  $k$ -step BDF method [6], where temporal derivatives in equation (1a) and (1b) are replaced with the following difference approximation

$$(\bullet)'_n = \frac{1}{h} \sum_{i=0}^k \alpha_i (\bullet)_{n-i}. \quad (2)$$

Although the analysis is here restricted to BDF methods, similar results and conclusions can be reached for other integration schemes for DAEs. For example, the case of the Newmark family of integrators is analyzed in Bottasso, Dopico, and Trainelli [7].

Using a Newton-type method, at each time step one solves a linear system of equations in the form

$$\mathbf{A}\mathbf{z} = \mathbf{b}, \quad (3)$$

where the iteration matrix is

$$\mathbf{A} = h\mathbf{J}_n = \begin{bmatrix} \alpha_0 \mathbf{I} & -h\mathbf{I} & \mathbf{0} \\ h\mathbf{X} & \alpha_0 \mathbf{I} + h\mathbf{Y} & -h\mathbf{G} \\ h\mathbf{G}^T & \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad (4)$$

with

$$\mathbf{X} = -\mathbf{f}_{,\mathbf{u}} - \mathbf{G}_{,\mathbf{u}}\boldsymbol{\lambda}, \quad (5a)$$

$$\mathbf{Y} = -\mathbf{f}_{,\mathbf{v}}. \quad (5b)$$

The three block rows of this matrix correspond to the kinematic, equilibrium and constraint equations of system (1), while the three block columns correspond to the  $\mathbf{u}$ ,  $\mathbf{v}$  and  $\boldsymbol{\lambda}$  variables, respectively.

Petzold and Lötstedt [1] show that, by using Gaussian elimination with partial pivoting, the accuracy in the  $i$ -th component of the solution can be estimated as

$$|\Delta z_i| \leq r\varepsilon \sum_j |(\mathbf{A}^{-1})_{ij}| \|\mathbf{A}\|_\infty \|\mathbf{z} + \Delta\mathbf{z}\|_\infty, \quad (6)$$

where  $r$  is an unknown coefficient,  $\varepsilon$  is the machine accuracy, and

$$(\mathbf{A} + \Delta\mathbf{A})(\mathbf{z} + \Delta\mathbf{z}) = \mathbf{b}, \quad (7)$$

with

$$\|\Delta\mathbf{A}\|_\infty \leq r\varepsilon \|\mathbf{A}\|_\infty. \quad (8)$$

The leading terms in the inverse of the iteration matrix (4) were derived in [1] as

$$\mathbf{A}^{-1} = (h\mathbf{J}_n)^{-1} = \frac{1}{\alpha_0} \begin{bmatrix} \mathbf{I} - \mathbf{T} & \gamma(\mathbf{I} - \mathbf{T})\mathbf{R}^{-1} & \gamma^{-1}\mathbf{R}^{-1}\mathbf{G}\mathbf{S} \\ -\gamma^{-1}\mathbf{T} & (\mathbf{I} - \mathbf{T})\mathbf{R}^{-1} & \gamma^{-2}\mathbf{R}^{-1}\mathbf{G}\mathbf{S} \\ -\gamma^{-2}\mathbf{S}\mathbf{G}^T & -\gamma^{-1}\mathbf{S}\mathbf{G}^T\mathbf{R}^{-1} & \gamma^{-3}\mathbf{S} \end{bmatrix}, \quad (9)$$

where  $\gamma = h/\alpha_0$  and

$$\mathbf{R} = \mathbf{I} + \gamma\mathbf{Y} + \gamma^2\mathbf{X}, \quad (10a)$$

$$\mathbf{S} = (\mathbf{G}^T\mathbf{R}^{-1}\mathbf{G})^{-1}, \quad (10b)$$

$$\mathbf{T} = \mathbf{R}^{-1}\mathbf{G}\mathbf{S}\mathbf{G}^T. \quad (10c)$$

Therefore, using (6) and (9), we conclude that the roundoff errors in the solution are

$$\Delta u_i = \mathcal{O}(h^{-1}), \quad (11a)$$

$$\Delta v_i = \mathcal{O}(h^{-2}), \quad (11b)$$

$$\Delta \lambda_i = \mathcal{O}(h^{-3}). \quad (11c)$$

This shows that, as  $h \rightarrow 0$ , the accuracy in the Lagrange multipliers deteriorates quickly. Furthermore, it is readily verified that

$$\|\mathbf{A}\|_\infty = \mathcal{O}(h^0), \quad (12a)$$

$$\|\mathbf{A}^{-1}\|_\infty = \mathcal{O}(h^{-3}). \quad (12b)$$

Hence, the condition number  $C = \|\mathbf{A}\|_\infty \|\mathbf{A}^{-1}\|_\infty$  becomes

$$C = \mathcal{O}(h^{-3}), \quad (13)$$

and the iteration matrix becomes severely ill conditioned for small  $h$ . The  $\mathcal{O}(h^{-3})$  dependence of roundoff errors and conditioning on the time step size poses limitations to the practical use of variable step size solvers and time refinement procedures.

### 3 Preconditioning by scaling

#### 3.1 Left preconditioning

As suggested in [1], this situation can be improved by considering a left diagonal scaling transformation of the equations. The scaled system can be written as

$$\tilde{\mathbf{A}}\mathbf{z} = \tilde{\mathbf{b}}, \quad (14)$$

where

$$\tilde{\mathbf{A}} = \mathbf{D}\mathbf{A}, \quad (15a)$$

$$\tilde{\mathbf{b}} = \mathbf{D}\mathbf{b}, \quad (15b)$$

and  $\mathbf{D}$  is a diagonal scaling matrix defined as

$$\mathbf{D} = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & h^{-1}\mathbf{I} \end{bmatrix}. \quad (16)$$

The inverse of the iteration matrix is in this case

$$\tilde{\mathbf{A}}^{-1} = \frac{1}{\alpha_0} \begin{bmatrix} \mathbf{I} - \mathbf{T} & \gamma(\mathbf{I} - \mathbf{T})\mathbf{R}^{-1} & \alpha_0\mathbf{R}^{-1}\mathbf{G}\mathbf{S} \\ -\gamma^{-1}\mathbf{T} & (\mathbf{I} - \mathbf{T})\mathbf{R}^{-1} & \alpha_0\gamma^{-1}\mathbf{R}^{-1}\mathbf{G}\mathbf{S} \\ -\gamma^{-2}\mathbf{S}\mathbf{G}^T & -\gamma^{-1}\mathbf{S}\mathbf{G}^T\mathbf{R}^{-1} & \alpha_0\gamma^{-2}\mathbf{S} \end{bmatrix}, \quad (17)$$

and, using again (6), we find

$$\Delta u_i = \mathcal{O}(h^0), \quad (18a)$$

$$\Delta v_i = \mathcal{O}(h^{-1}), \quad (18b)$$

$$\Delta \lambda_i = \mathcal{O}(h^{-2}). \quad (18c)$$

The situation is improved with respect to the original (unscaled) system, but the error in the Lagrange multipliers still grows rapidly as  $h \rightarrow 0$  even for the left preconditioned problem. Similarly, we find  $\|\tilde{\mathbf{A}}\|_\infty = \mathcal{O}(h^0)$  and  $\|\tilde{\mathbf{A}}^{-1}\|_\infty = \mathcal{O}(h^{-2})$ , which leads to the conclusion that the condition number is

$$C = \mathcal{O}(h^{-2}), \quad (19)$$

again not a favorable behavior.

### 3.2 Left and right preconditioning

A further improvement on the situation expressed by equation (18) can be obtained by considering a left and right preconditioning, i.e a scaling of the equations together with a scaling of the unknowns. In particular, since the Lagrange multipliers are affected by the largest roundoff errors according to (18), the governing equations are here re-written in the following form

$$\mathbf{u}' = \mathbf{v}, \quad (20a)$$

$$\mathbf{v}' = \mathbf{f} + s\mathbf{G}\hat{\boldsymbol{\lambda}}, \quad (20b)$$

$$\boldsymbol{\phi} = \mathbf{0}, \quad (20c)$$

where  $\hat{\boldsymbol{\lambda}} = \boldsymbol{\lambda}/s$  are scaled Lagrange multipliers, and the scaling factor is  $s = \mathcal{O}(h^{-2})$ . The idea is now to solve the problem in terms of the scaled multipliers  $\hat{\boldsymbol{\lambda}}$  and generalized coordinates  $\mathbf{u}$  and velocities  $\mathbf{v}$ , and to recover the values of  $\boldsymbol{\lambda}$  once at convergence.

Consider first for simplicity  $s = 1/h^2$ . The linear iteration problem for the left and right preconditioned system can be written as

$$\hat{\mathbf{A}}\hat{\mathbf{z}} = \hat{\mathbf{b}}, \quad (21)$$

where

$$\hat{\mathbf{A}} = \mathbf{D}_L\mathbf{A}\mathbf{D}_R, \quad (22a)$$

$$\hat{\mathbf{z}} = \mathbf{D}_R^{-1}\mathbf{z}, \quad (22b)$$

$$\hat{\mathbf{b}} = \mathbf{D}_L\mathbf{b}. \quad (22c)$$

The left diagonal scaling matrix  $\mathbf{D}_L$  is defined as

$$\mathbf{D}_L = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & h\mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & h^{-1}\mathbf{I} \end{bmatrix}. \quad (23)$$

The scaling of the third block rows is the same as in [1], while the scaling of the second block rows is needed to ensure  $\|\hat{\mathbf{A}}\|_\infty = \mathcal{O}(h^0)$ , which is beneficial for the condition number. The right diagonal transformation  $\mathbf{D}_R$  which scales the Lagrange multipliers is defined as

$$\mathbf{D}_R = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & h^{-2}\mathbf{I} \end{bmatrix}. \quad (24)$$

This gives the following form for the inverse of the iteration matrix

$$\hat{\mathbf{A}}^{-1} = \frac{1}{\alpha_0} \begin{bmatrix} \mathbf{I} - \mathbf{T} & \alpha_0^{-1}(\mathbf{I} - \mathbf{T})\mathbf{R}^{-1} & \alpha_0\mathbf{R}^{-1}\mathbf{G}\mathbf{S} \\ -\gamma^{-1}\mathbf{T} & \alpha_0^{-1}\gamma^{-1}(\mathbf{I} - \mathbf{T})\mathbf{R}^{-1} & \alpha_0\gamma^{-1}\mathbf{R}^{-1}\mathbf{G}\mathbf{S} \\ -\alpha_0^2\mathbf{S}\mathbf{G}^T & -\alpha_0\mathbf{S}\mathbf{G}^T\mathbf{R}^{-1} & \alpha_0^3\mathbf{S} \end{bmatrix}. \quad (25)$$

Using (6), the roundoff errors for the three solution fields are in this case

$$\Delta u_i = \mathcal{O}(h^0), \quad (26a)$$

$$\Delta v_i = \mathcal{O}(h^{-1}), \quad (26b)$$

$$\Delta \hat{\lambda}_i = \mathcal{O}(h^0). \quad (26c)$$

Notice that the roundoff errors in the displacement field and in the scaled Lagrange multipliers are now insensitive to the time step size, while the condition number is

$$C = \mathcal{O}(h^{-1}). \quad (27)$$

The Lagrange multipliers are recovered at convergence of the Newton process from the scaled multipliers using  $\boldsymbol{\lambda} = s\hat{\boldsymbol{\lambda}}$ . The scaled multipliers will be affected by an error due to machine accuracy and to the Newton termination tolerance, error which is  $\mathcal{O}(h^0)$  for the results above. Since  $s$  is  $\mathcal{O}(h^{-2})$ , then  $\boldsymbol{\lambda}$  will be affected by an error  $\mathcal{O}(h^{-2})$  due to the recovery process. However, the reduced accuracy of the multipliers can not affect the iterative Newton process, since this is formulated in terms of the scaled variables. Hence, the Newton iterations can be carried out until as tight a tolerance as necessary has been achieved.

A more general form for the multiplier scaling factor can be derived based on a dimensional argument. Recall that in general the dynamic equilibrium equations (1b) express the balance of inertial forces, elastic forces, constraint reactions and externally applied forces. To reflect the different nature of these forces,  $s$  is selected as

$$s = k_{\text{ave}} + \frac{m_{\text{ave}}}{h^2}, \quad (28)$$

where  $k_{\text{ave}}$  represents an average stiffness term for the system and  $m_{\text{ave}}$  an average mass term. To understand the dimensional effect of this scaling, let  $u_{\text{ave}}$  denote a representative displacement of the system,  $F_{\text{ave}}^{\text{elas}}$  a representative value of the elastic forces and finally  $F_{\text{ave}}^{\text{iner}}$  a representative value of the inertial forces. For large time step sizes,  $s \approx k_{\text{ave}}$ , and  $\hat{\boldsymbol{\lambda}} \approx \boldsymbol{\lambda}/k_{\text{ave}}$ . Then the ratio of the constraint reactions and of the elastic forces is approximately  $s\mathbf{G}\hat{\boldsymbol{\lambda}}/F_{\text{ave}}^{\text{elas}} \approx k_{\text{ave}}\mathbf{G}\hat{\boldsymbol{\lambda}}/(k_{\text{ave}}u_{\text{ave}}) \approx \mathbf{G}\hat{\boldsymbol{\lambda}}/u_{\text{ave}}$ ; this means that the scaled Lagrange multipliers will be roughly of the same order of magnitude as the displacements, balancing the magnitude of all the unknowns of the problem. For very small time step sizes,  $s \approx m_{\text{ave}}/h^2$ , and  $\hat{\boldsymbol{\lambda}} \approx \boldsymbol{\lambda}/(m_{\text{ave}}/h^2)$ . It then follows that  $s\mathbf{G}\hat{\boldsymbol{\lambda}}/F_{\text{ave}}^{\text{iner}} \approx (m_{\text{ave}}\mathbf{G}\hat{\boldsymbol{\lambda}}/h^2)/(m_{\text{ave}}u_{\text{ave}}/h^2) \approx \mathbf{G}\hat{\boldsymbol{\lambda}}/u_{\text{ave}}$ ; this means that even in this case the scaled Lagrange multipliers will be roughly of the same order of magnitude as the displacements, balancing again the magnitude of all the unknowns of the problem.

The roundoff errors and the conditioning of the left and right scaled system are much improved with respect to the unscaled or the left scaled ones. Yet, the velocity field still shows a mild

dependence on the time step size, which in turn induces the same behavior in the conditioning. This effect can be eliminated by considering a scaling of the unknowns that includes, together with the scaling of the Lagrange multipliers, also a scaling of the velocities, as suggested by Arnold [8]. This amounts to the following redefinition of the right diagonal transformation  $\mathbf{D}_R$ :

$$\mathbf{D}_R = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & h^{-1}\mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & h^{-2}\mathbf{I} \end{bmatrix}. \quad (29)$$

The inverse of the iteration matrix is readily found as

$$\hat{\mathbf{A}}^{-1} = \frac{1}{\alpha_0} \begin{bmatrix} \mathbf{I} - \mathbf{T} & \alpha_0^{-1}(\mathbf{I} - \mathbf{T})\mathbf{R}^{-1} & \alpha_0\mathbf{R}^{-1}\mathbf{G}\mathbf{S} \\ -\alpha_0\mathbf{T} & (\mathbf{I} - \mathbf{T})\mathbf{R}^{-1} & \alpha_0^2\mathbf{R}^{-1}\mathbf{G}\mathbf{S} \\ -\alpha_0^2\mathbf{S}\mathbf{G}^T & -\alpha_0\mathbf{S}\mathbf{G}^T\mathbf{R}^{-1} & \alpha_0^3\mathbf{S} \end{bmatrix}, \quad (30)$$

which does not depend on  $h$ . Hence, the roundoff errors in all solution variables and the conditioning are  $\mathcal{O}(h^0)$ .

### 3.3 Numerical example

We consider the problem of the pendulum, with equations

$$u'_x = v_x, \quad (31a)$$

$$u'_y = v_y, \quad (31b)$$

$$v'_x = \lambda u_x, \quad (31c)$$

$$v'_y = \lambda u_y - 1, \quad (31d)$$

$$0 = \frac{1}{2}(u_x^2 + u_y^2 - 1), \quad (31e)$$

where  $u_x, u_y$  are the Cartesian coordinates of the point mass, and  $v_x, v_y$  the mass velocity components. Bar length, point mass and acceleration of gravity are all equal to 1. The point mass is initially at rest with  $u_x = 1, u_y = 0$ , and falls under the action of gravity. The integration of equation (31) is performed using the 2-step BDF formula until time  $t = 1 \cdot 10^{-3}$ , using the time step sizes  $h = \{1 \cdot 10^{-4}, 1 \cdot 10^{-5}, 1 \cdot 10^{-6}, 1 \cdot 10^{-7}, 1 \cdot 10^{-8}\}$ .

At each time step, the non-linear discrete equations are solved using Newton method. At the  $j$ -th Newton iteration of time step  $n$ , we solve the linear system of equations

$$\mathbf{A}_n^j \mathbf{z}_n^j = \mathbf{b}_n^j, \quad (32)$$

which yields the corrections  $\mathbf{z}_n^j$ . Typically, the norm of the corrections will decrease at each Newton step, until the accumulation of roundoff errors leads to a saturation value. If further iterations are carried out, one typically observes an oscillation of the correction norm around its saturation value. An example of this behavior is given in Figure 1, which shows  $\|\mathbf{z}_n^j\|$  vs. the Newton iteration  $j$  for the scheme without preconditioning and  $h = 1 \cdot 10^{-7}$ . In this case, the saturation is reached at the fourth iteration, after which the correction norm has very small oscillations around a saturation value of approximately  $4 \cdot 10^{-3}$ . To quantify this effect of roundoff errors, in this work iterations are arrested when the Newton corrections stop decreasing, i.e. when we detect the condition

$$\|\mathbf{z}_n^{j+1}\| \geq \|\mathbf{z}_n^j\|. \quad (33)$$

Hence, the magnitude of the last decreasing Newton correction gives an indication of the tightest achievable convergence of the Newton iterations, which can not be further improved no matter how many iterations one carries out.

At first, we consider the scaling of the sole Lagrange multipliers based on the simple choice  $s = 1/h^2$ , and we use the definition of  $\mathbf{D}_R$  given in (24). Figure 2 shows the maximum throughout each simulation of the 2-norm of the last decreasing Newton corrections versus the time step size, i.e.  $\max_n \| \mathbf{z}_n^j \|$  vs.  $h$ . The results obtained with the left preconditioned problem are shown with the  $\square$  symbol, while the ones of the left and right preconditioning with the  $\circ$  symbol. Notice that with left preconditioning alone it is in fact impossible to solve the non-linear iteration with a tight tolerance for small values of the time step.

Figure 3 shows the maximum throughout each simulation of the absolute value of the last decreasing Newton corrections by field type, i.e.  $\max_n \max(|z_{u_x,n}^j|, |z_{u_y,n}^j|)$ ,  $\max_n \max(|z_{v_x,n}^j|, |z_{v_y,n}^j|)$ ,  $\max_n |z_{\lambda,n}^j|$ . These results are in accordance with the analysis of the previous sections.

Finally, Figure 4 shows the condition number  $C$  vs.  $h$ . The condition number is computed by means of the singular value decomposition of the iteration matrix at convergence, i.e. at the last decreasing Newton iteration. Here again, a substantial improvement is observed for the left and right preconditioning with respect to the left preconditioning alone.

Next, we consider the previous scaling of the multipliers together with the scaling of the velocities, with  $\mathbf{D}_R$  as given in (29). Figure 5 shows the maximum throughout each simulation of the 2-norm of the last decreasing Newton corrections versus the time step size, while Figure 6 shows the condition number  $C$  vs.  $h$ . In both cases, the solution of the left and right preconditioned problem does not depend on  $h$ , as predicted by the analysis.

## 4 Displacement based form

For large scale applications denoted by many degrees of freedom, the velocity-displacement-multiplier form discussed in the previous section leads to the solution of large linear systems. For example, this is the case when the governing DAEs (1) are obtained by spatial discretization using the finite element method of a flexible mechanical system. A substantial performance increase can be obtained by implementing the integration scheme in displacement-multiplier form, which is simply achieved by statically eliminating the velocity unknowns at each time step. In fact, considering again the  $k$ -step BDF method, velocities at time  $t_n$  are computed from the kinematic equations (1a) using the difference approximation (2) as

$$\mathbf{u}'_n = \mathbf{v}_n = \frac{1}{h} \sum_{i=0}^k \alpha_i \mathbf{u}_{n-i}. \quad (34)$$

Inserting this expression into the accelerations of the dynamic equilibrium equations (1b), we get

$$\mathbf{v}'_n = \frac{1}{h} \sum_{i=0}^k \alpha_i \mathbf{v}_{n-i} = \frac{1}{h} \left( \frac{\alpha_0}{h} \sum_{i=0}^k \alpha_i \mathbf{u}_{n-i} + \sum_{i=1}^k \alpha_i \mathbf{v}_{n-i} \right). \quad (35)$$

Using Newton method, one solves the reduced system of linear equations

$$\mathbf{B}\mathbf{y} = \mathbf{c}, \quad (36)$$

where the iteration matrix is in this case

$$\mathbf{B} = h\mathbf{J}_n = \begin{bmatrix} \alpha_0 \mathbf{U} & -h\mathbf{G} \\ h\mathbf{G}^T & \mathbf{0} \end{bmatrix}, \quad (37)$$

with

$$\mathbf{U} = \gamma^{-1} \mathbf{I} + \mathbf{Y} + \gamma \mathbf{X}. \quad (38)$$



The first block rows of this matrix correspond to the equilibrium equations in displacement-based form, while the second block rows correspond to the constraint conditions. The two block columns correspond to the  $\mathbf{u}$  and  $\boldsymbol{\lambda}$  variables, respectively. Once at convergence, velocities are recovered using the vector relationship (34).

The inverse of the iteration matrix is

$$\mathbf{B}^{-1} = \frac{1}{\alpha_0} \begin{bmatrix} (\mathbf{I} - \mathbf{W})\mathbf{U}^{-1} & \gamma^{-1}\mathbf{U}^{-1}\mathbf{G}\mathbf{V} \\ -\gamma^{-1}\mathbf{V}\mathbf{G}^T\mathbf{U}^{-1} & \gamma^{-2}\mathbf{V} \end{bmatrix}, \quad (39)$$

where

$$\mathbf{V} = (\mathbf{G}^T\mathbf{U}^{-1}\mathbf{G})^{-1}, \quad (40a)$$

$$\mathbf{W} = \mathbf{U}^{-1}\mathbf{G}\mathbf{V}\mathbf{G}^T. \quad (40b)$$

Using (6), and considering that

$$\lim_{h \rightarrow 0} \mathbf{U}^{-1} = \mathcal{O}(h), \quad (41a)$$

$$\lim_{h \rightarrow 0} \mathbf{V} = \mathcal{O}(h^{-1}), \quad (41b)$$

the roundoff errors in the primary solution variables are found to be

$$\Delta u_i = \mathcal{O}(h^{-1}), \quad (42a)$$

$$\Delta \lambda_i = \mathcal{O}(h^{-3}), \quad (42b)$$

while the condition number  $C$  becomes

$$C = \mathcal{O}(h^{-4}), \quad (43)$$

again very unfavorable behaviors with respect to the time step size.

## 4.1 Left preconditioning

Even in this case, the simple left preconditioning of Petzold and Lötstedt [1] gives a little benefit with respect to the unpreconditioned problem. In fact, defining the left scaling transformation as

$$\mathbf{D} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & h^{-1}\mathbf{I} \end{bmatrix}, \quad (44)$$

we find

$$\tilde{\mathbf{B}}^{-1} = \frac{1}{\alpha_0} \begin{bmatrix} (\mathbf{I} - \mathbf{W})\mathbf{U}^{-1} & \alpha_0\mathbf{U}^{-1}\mathbf{G}\mathbf{V} \\ -\gamma^{-1}\mathbf{V}\mathbf{G}^T\mathbf{U}^{-1} & \alpha_0\gamma^{-1}\mathbf{V} \end{bmatrix}, \quad (45)$$

which gives for the propagation of errors

$$\Delta u_i = \mathcal{O}(h^0), \quad (46a)$$

$$\Delta \lambda_i = \mathcal{O}(h^{-2}), \quad (46b)$$

while the conditioning is

$$C = \mathcal{O}(h^{-3}). \quad (47)$$

## 4.2 Left and right preconditioning

The left and right scaling matrices for the reduced problem are simply obtained from (21–24), and are defined as

$$\mathbf{D}_L = \begin{bmatrix} h\mathbf{I} & \mathbf{0} \\ \mathbf{0} & h^{-1}\mathbf{I} \end{bmatrix}, \quad (48)$$

and

$$\mathbf{D}_R = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & h^{-2}\mathbf{I} \end{bmatrix}, \quad (49)$$

where we have considered again the simple choice  $s = 1/h^2$ . This gives

$$\hat{\mathbf{B}}^{-1} = \frac{1}{\alpha_0} \begin{bmatrix} \alpha_0^{-1}\gamma^{-1}(\mathbf{I} - \mathbf{W})\mathbf{U}^{-1} & \alpha_0\mathbf{U}^{-1}\mathbf{G}\mathbf{V} \\ -\alpha_0\mathbf{V}\mathbf{G}^T\mathbf{U}^{-1} & \alpha_0^3\gamma\mathbf{V} \end{bmatrix}, \quad (50)$$

which yields time-step-size-independent sensitivity to perturbations of the primary variables

$$\Delta u_i = \mathcal{O}(h^0), \quad (51a)$$

$$\Delta \hat{\lambda}_i = \mathcal{O}(h^0), \quad (51b)$$

and a time-step-size-independent conditioning, i.e.

$$C = \mathcal{O}(h^0), \quad (52)$$

the best possible behavior that can be expected of a numerical method.

## 4.3 Numerical example

We consider again the model problem (31), using the same data. At the  $j$ -th Newton iteration of time step  $n$ , we solve the linear system of equations

$$\mathbf{B}_n^j \mathbf{y}_n^j = \mathbf{c}_n^j, \quad (53)$$

which yields the corrections  $\mathbf{y}_n^j$ . As in the previous case, iterations are arrested only when the Newton corrections stop decreasing.

The results of the analysis are confirmed by Figure 7 and Figure 8. In particular, the former shows the maximum throughout each simulation of the 2-norm of the last decreasing Newton corrections versus the time step size,  $\max_n \|\mathbf{z}_n^j\|$  vs.  $h$ , while the latter reports the condition number  $C$  vs.  $h$ . For the left and right preconditioning, exact time step size independence is achieved.

## 5 Numerical example: the need for scaling in contact-impact analysis of flexible multibody systems

Fig. 9a depicts a cam of length 0.4 m connected to the ground at point **A** by means of a revolute joint, which relative rotation is prescribed to be  $\phi(t) = \phi_0(1 - \cos\Omega t)$ , where  $\phi_0 = \pi/3$  rad, and  $\Omega = 2\pi/3$  rad/sec. At the tip of the cam, a rigid impactor of total mass  $M_I = 4$  kg is attached. A flexible beam of length 2.4 m is clamped to the ground at point **R**, and a point mass  $M_T = 40$  kg is attached at its free end, point **T**. At the middle point **P** of the beam, a rigid body of total mass  $M_P = 4$  kg is attached. As the impactor moves, it contacts the outer surface of this rigid body, as shown in fig. 9b. The physical properties of the flexible beam are: axial stiffness  $EA = 44000$  kN, bending stiffnesses  $I_{22} = 300$ ,  $I_{33} = 23$  kN·m<sup>2</sup>, torsional stiffness  $GJ = 28$  kN·m<sup>2</sup>, and mass per

unit span  $m = 1.6 \text{ kg/m}$ . The bending stiffnesses of the cam are  $I_{22} = 23$ ,  $I_{33} = 300 \text{ kN}\cdot\text{m}^2$ , and the axial stiffness, torsional stiffness and mass per unit span are the same as those of the beam.

When the two bodies come in contact, a contact model and friction model are used to define the normal and tangential contact forces, respectively. A linear spring with a stiffness constant  $k = 2 \text{ GN/m}$  in parallel with a viscous damper of constant  $\mu = 10^{-3}$  were used to evaluate the normal contact force. The friction force was computed with the help of the LuGre model, using the following parameters:  $\sigma_0 = 10^5 \text{ m}^{-1}$ ,  $\sigma_1 = \sigma_2 = 0 \text{ sec/m}$ ,  $v_s = 10^{-3} \text{ m/sec}$ ,  $\mu_k = 0.30$ ,  $\mu_s = 0.30$ , and  $\gamma = 2$ . More details about this modeling technique and the definition of these parameters can be found in Bauchau and Ju [9].

The simulation was conducted for a 1 sec duration. From  $t = 0$  to 0.4 sec a constant time step size,  $h = 10^{-2} \text{ sec}$ , was used and at that point, an automated time step size adaptivity algorithm was activated [10]; the desired local error was set to  $\varepsilon_{\text{loc}} = 10^{-7}$ . In light of equation (28), the scaling factor was selected as  $s = k_{\text{ave}} + m_{\text{ave}}/h^2$ , where  $k_{\text{ave}} = 10^3$  and  $m_{\text{ave}} = 0$  or 1. When  $m_{\text{ave}} = 0$ , the scaling is time step size independent, whereas it becomes time step dependent when  $m_{\text{ave}} = 1$ .

For the problem at hand, large time step sizes were used before contact occurs, i.e. when  $t < 0.494 \text{ sec}$ . This portion of the simulation is not challenging and both scaling factors performed equally well. At time  $t \approx 0.494$  the time step size was drastically reduced by the adaptivity algorithm to decrease the large contact force gradients that would result from using a constant time step. Figure 10 depicts the time step size as a function of time during the simulation, for both  $m_{\text{ave}} = 0$  and 1. As should be expected from the analysis developed in this paper, when  $m_{\text{ave}} = 0$ , the simulation fails to converge as soon as the time step size becomes small, whereas for  $m_{\text{ave}} = 1$ , convergence is achieved at each time step, despite the complex dynamic behavior of the system after impact. Figure 11 shows the corresponding history of the scaling factor. Note that scaling factors of the order of  $10^{13}$  were needed at the instant of contact; after impact, the highly vibratory response of the system required small time steps and the scaling factor was in the range of  $s \in [10^8, 10^{11}]$ .

## 6 Conclusions

We have proposed a simple scaling transformation for the index three DAEs describing constrained multibody dynamical systems. The approach amounts to a left and right preconditioning of the iteration matrix, which has the goal of making the solution less sensitive to the propagation of perturbations and of improving the condition number.

Using the proposed left preconditioning, the constraint equations are divided by the time step size as already proposed by [1], while the dynamic equilibrium equations are multiplied by the same quantity. The right preconditioning amounts to a definition of a set of modified unknowns: modified multipliers are obtained by scaling with the square of the time step size of the original Lagrange multipliers, while modified velocities are obtained by scaling with the time step size.

The new scaling was applied to the velocity-displacement-multiplier and to the reduced displacement-multiplier forms. In both cases, we obtain the remarkable results that both error propagations and conditioning are  $\mathcal{O}(h^0)$ , i.e. the numerical solution of index three DAEs behaves as in the case of regular ODEs.

The simple problem of the pendulum was used to numerically confirm the results of the analysis, while a more challenging problem dealing with the contact-impact analysis of a flexible multibody system was used to illustrate the effect of scaling when using time refinement procedures.

Although other approaches are available to deal with the problem of ill conditioning and error propagation for high index DAEs, the present approach has the possible advantages of being trivial to implement, of not requiring a re-writing of the equations of motion, and furthermore it does not introduce additional unknowns.

## References

- [1] L. PETZOLD AND P. LÖTSTEDT, *Numerical solution of nonlinear differential equations with algebraic constraints II: practical implications*, SIAM Journal on Scientific and Statistical Computing, 7 (1986), pp. 721–733.
- [2] M. ARNOLD, *A perturbation analysis for the dynamical simulation of mechanical multibody systems*, Applied Numerical Mathematics, 18 (1995), pp. 37–56.
- [3] C. GEAR, B. LEIMKUHNER, AND G. GUPTA, *Automatic integration of Euler-Lagrange equations with constraints*, Journal of Computational and Applied Mathematics, 12–13 (1985), pp. 77–90.
- [4] M. BORRI, L. TRAINELLI, AND A. CROCE, *The embedded projection method: a general index reduction procedure for constrained system dynamics*, Computer Methods in Applied Mechanics and Engineering, (2005) accepted, to appear.
- [5] A. CARDONA, AND M. GÉRARDIN, *Numerical integration of second order differential-algebraic systems in flexible mechanism dynamics*, in Computer-Aided Analysis of Rigid and Flexible Mechanical Systems, M.F.O. Seabra Pereira, and J.A.C. Ambrósio, Eds., Kluwer Academic Publishers, Dordrecht, Boston, 1994.
- [6] C. GEAR, *Simultaneous numerical solution of differential/algebraic equations*, IEEE Transactions on Circuits and Systems, 18 (1971), pp. 89–95.
- [7] C.L. BOTTASSO, D. DOPICO, AND L. TRAINELLI, *On the optimal scaling of index three DAEs in multibody dynamics*, in Proceedings of the III European Conference on Computational Mechanics, Solids, Structures and Coupled Problems in Engineering, C.A. Mota Soares et al., Eds., Lisbon, Portugal, 5–8 June 2006.
- [8] M. ARNOLD, *Private communication*, (2005).
- [9] O.A. BAUCHAU AND C.K. JU, *Modeling Friction Phenomena in Flexible Multibody Dynamics*. Computer Methods in Applied Mechanics and Engineering, to appear, 2006.
- [10] O.A. BAUCHAU, *Computational Schemes for Flexible, Nonlinear Multi-Body Systems*, Multi-body System Dynamics, 2(2) (1998), pp. 169–225.

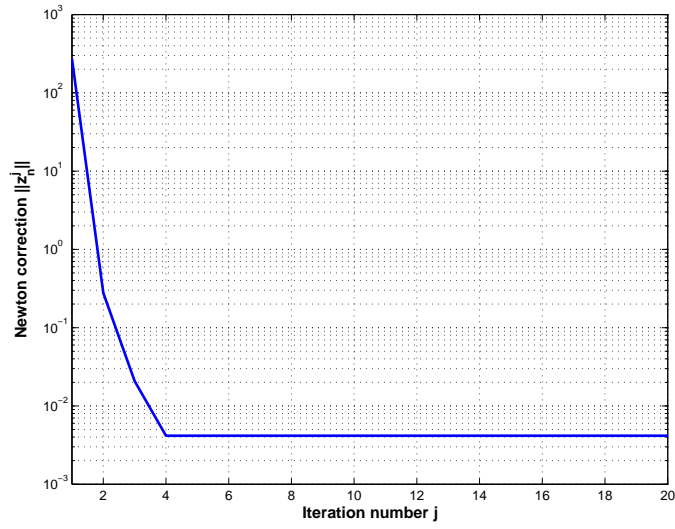


Figure 1: Convergence and saturation for the norm of the Newton corrections  $z_n^j$  plotted vs. the Newton iteration  $j$  for the scheme without preconditioning and  $h = 1 \cdot 10^{-7}$ .

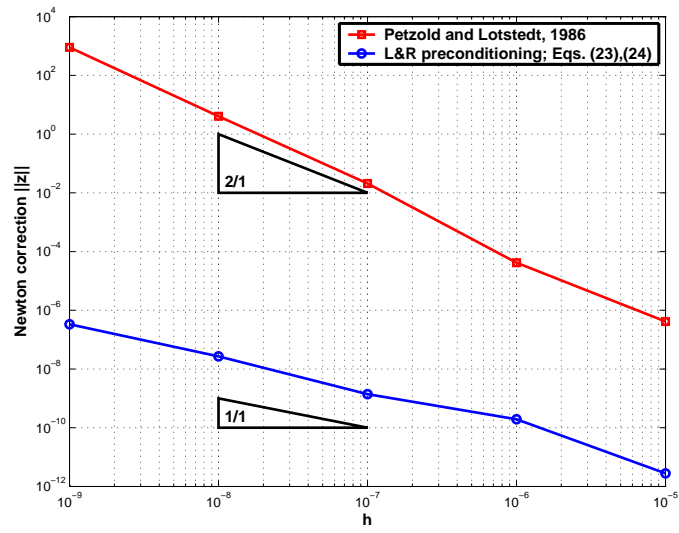


Figure 2: Last decreasing Newton corrections vs. time step size.

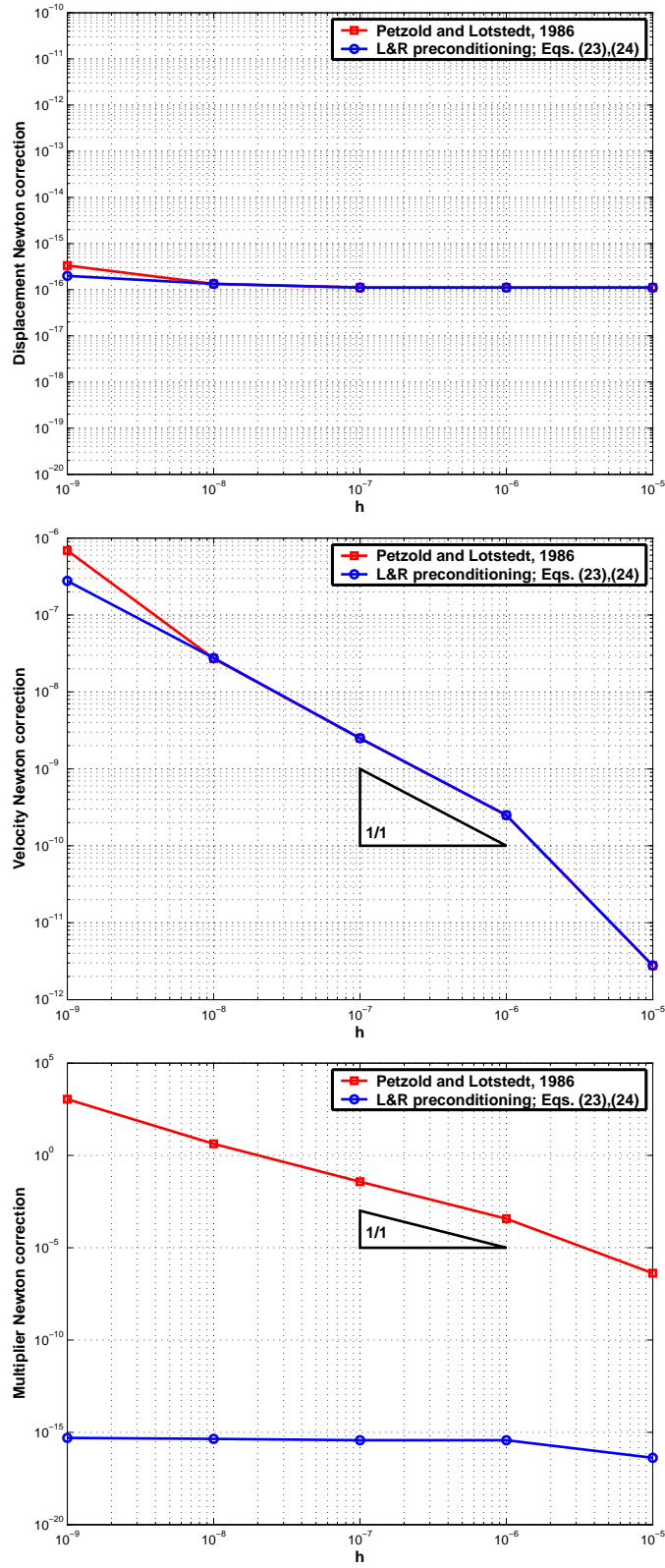


Figure 3: Last decreasing Newton corrections by field type vs. time step size. Top: displacements; center: velocities; bottom: Lagrange multipliers.

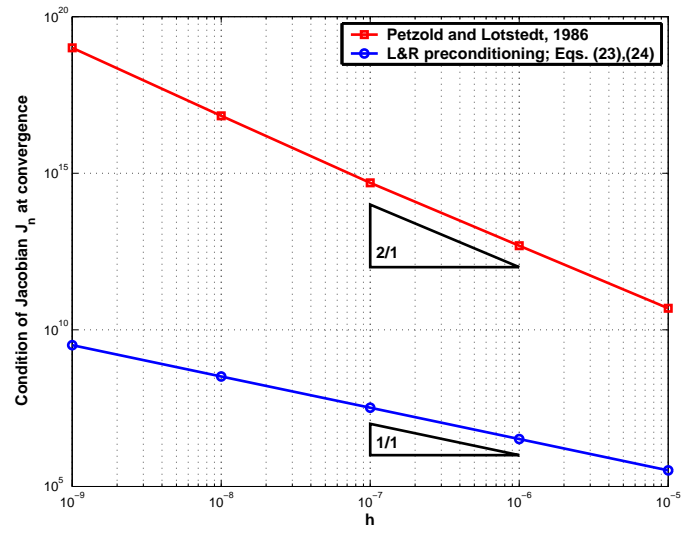


Figure 4: Condition number of Jacobian at convergence vs. time step size.



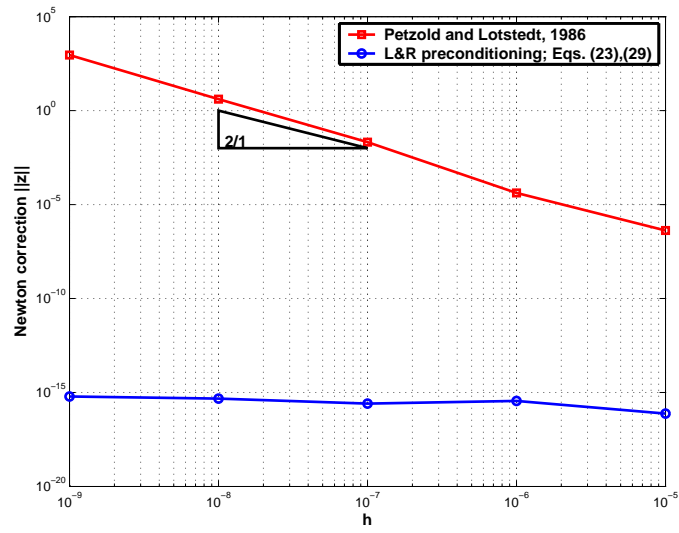


Figure 5: Last decreasing Newton corrections vs. time step size. Left and right preconditioning based on (23) and (29).

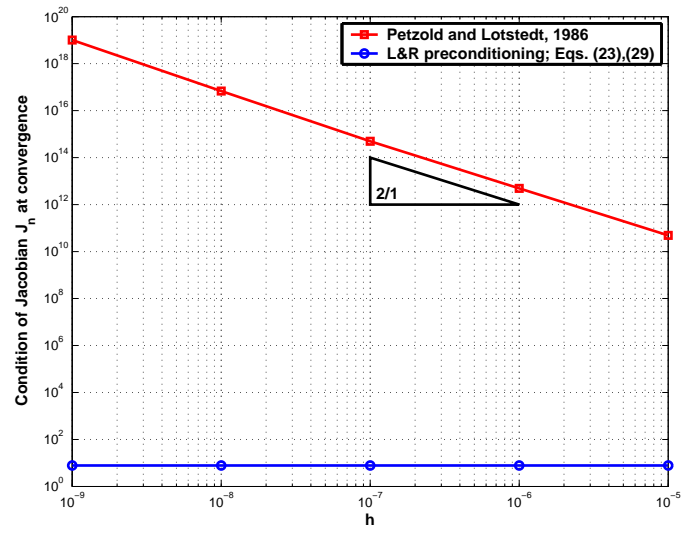


Figure 6: Condition number of Jacobian at convergence vs. time step size. Left and right preconditioning based on (23) and (29).

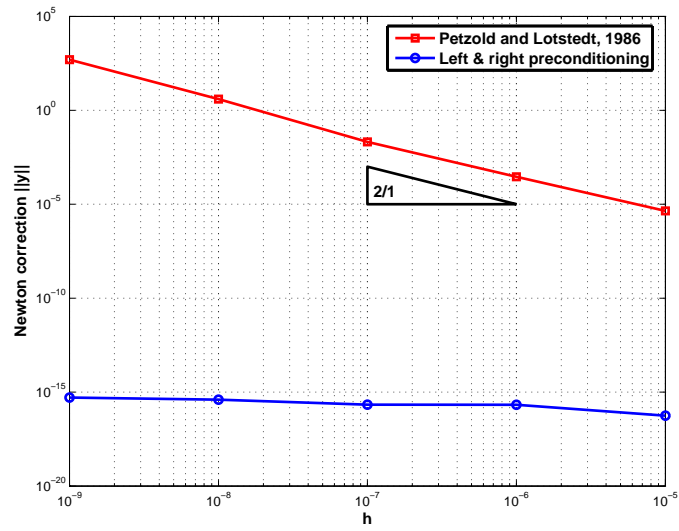


Figure 7: Last decreasing Newton corrections vs. time step size for the displacement-multiplier form.

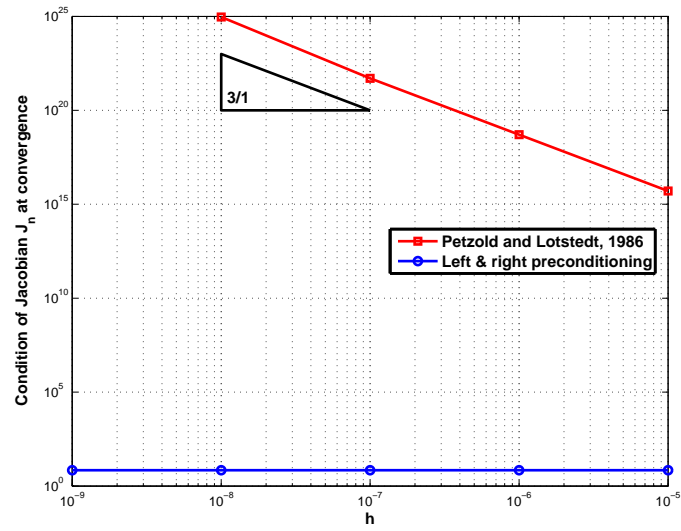


Figure 8: Condition number of Jacobian at convergence vs. time step size for the displacement-multiplier form.

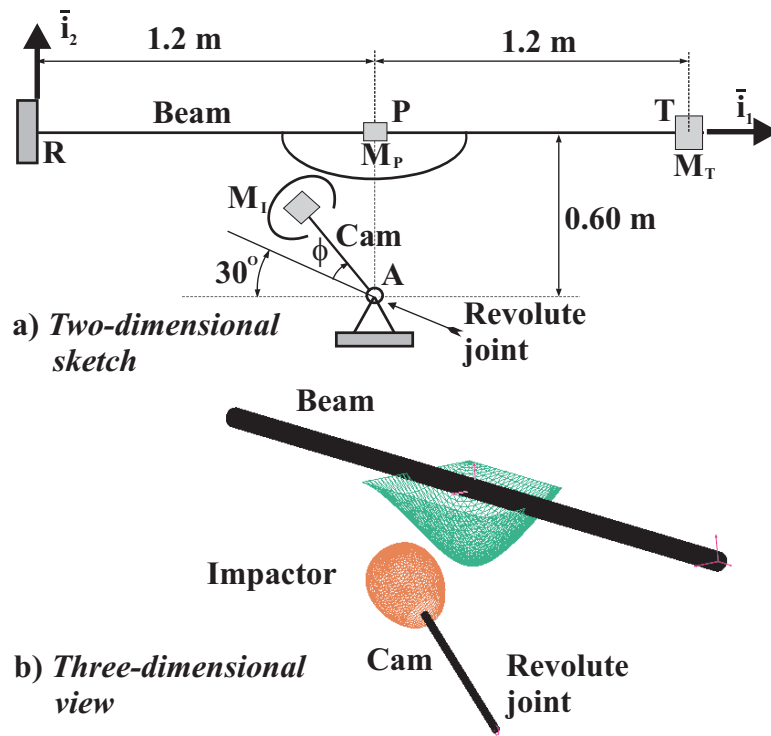


Figure 9: Configuration of the contact problem.

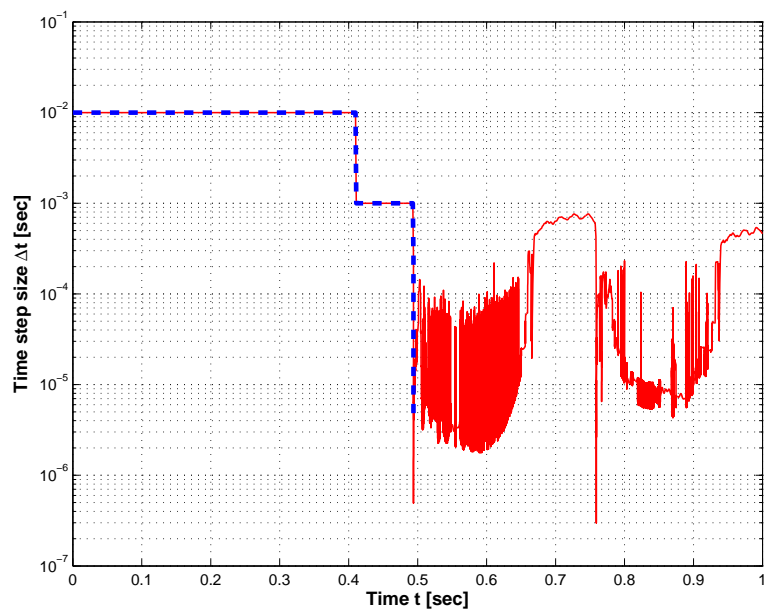


Figure 10: Time step size as a function of time.  $m_{\text{ave}} = 0$ : thick dashed line;  $m_{\text{ave}} = 1$ : thin solid line.

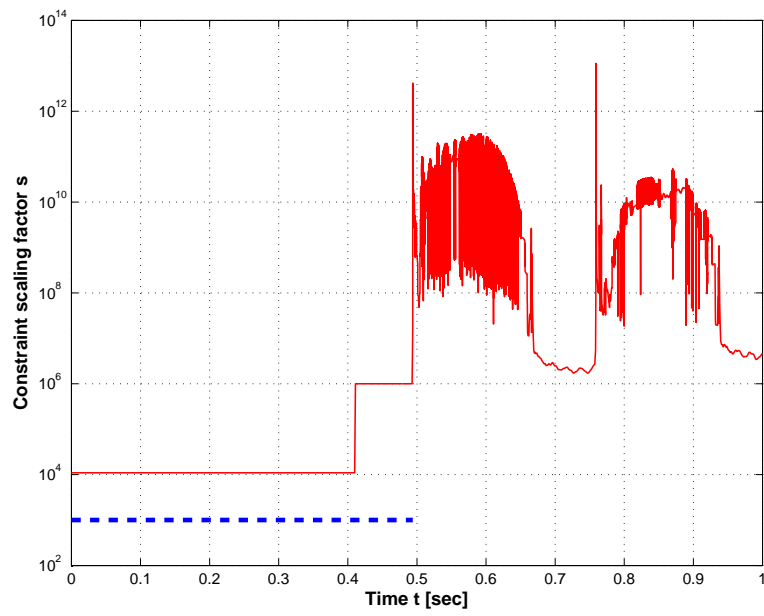


Figure 11: Constraint scaling factor  $s$ .  $m_{\text{ave}} = 0$ : thick dashed line;  $m = 1$ : thin solid line.