

# A matrix-free cone complementarity approach for solving large-scale, nonsmooth, rigid body dynamics

A.Tasora<sup>a</sup> M.Anitescu<sup>b</sup>

<sup>a</sup>Università degli Studi di Parma, Dipartimento di Ingegneria Industriale, 43100 Parma, Italy, [tasora@ied.unipr.it](mailto:tasora@ied.unipr.it)

<sup>b</sup>Mathematics and Computer Science Division, Argonne National Laboratory, 9700 South Cass Avenue, Argonne, IL 60439, USA, [anitescu@mcs.anl.gov](mailto:anitescu@mcs.anl.gov)

---

## Abstract

This paper proposes an iterative method that can simulate mechanical systems featuring a large number of contacts and joints between rigid bodies. The numerical method behaves as a contractive mapping that converges to the solution of a cone complementarity problem by means of iterated fixed-point steps with separable projections onto convex manifolds. Since computational speed and robustness are important issues when dealing with a large number of frictional contacts, we have performed special algorithmic optimizations in order to translate the numerical scheme into a matrix-free algorithm with  $O(n)$  space complexity and easy implementation. A modified version, that can run on parallel computers is discussed. A multithreaded version of the method has been used to simulate systems with more than a million contacts with friction.

*Key words:* Large-scale multibody, contact, friction, cone complementarity problem

---

## 1. Introduction

Many engineering problems involve unilateral contacts between rigid bodies, for instance, in simulations of robotic cells and part feeders, in cam followers, in masonry stability analysis, and in packaging devices such as those depicted in Fig. 1.

The dynamical simulation of such systems is complicated by the nonsmooth nature of the frictional constraints. When the number of contacts between bodies increases to thousands or millions, as in the case of granular flows in silos or in rock-soil dynamics, the computational efficiency of traditional methods can become an issue even on supercomputers.

A straightforward approach to solve this class of problems may consist of regularization schemes, that transform the discontinuities into a stiff force field. This is, for example, the approach often adopted by discrete element schemes (DEMs) because it does not require major modifications to traditional solvers based on smooth ordinary differential equations (ODEs) [15,35,36,30]. Nevertheless, although successfully used to simulate granular flows with many contacts, the regularization approach requires small time steps to achieve numerical stability. Moreover, it forces the user to introduce artificial stiffness or heuristic parameters: actually, if the deformation of the parts is negligible, a method that can use large time steps and unconditionally rigid bodies would be more welcome.

These considerations encouraged our research on a fast, robust, and unified numerical scheme that can handle complex mechanical systems made of rigid bodies with an arbitrary number of contacts and joints. Such a scheme aims at simulating mechanical systems ranging from the simplest (articulated linkages with few bilateral kinematic pairs and motors) up to the most complex (for example, a bulldozer interacting with millions of particles of sand with the tracks and the blade).

In this context, the biggest challenge comes from the discontinuous nature of the adhesion constraints and non-interpenetration constraints; in fact, the simulation of rigid contacts embeds the solution of nonsmooth equations. To this end, the straightforward application of numerical methods for ODEs or differential algebraic equations (DAEs) is inefficient. In fact, a naive approach based on piecewise integrals is virtually impossible because it would require stopping and restarting the integrator at each discontinuity to change the active set of constraints. This could work only if there were a limited number of unilateral constraints [20,21]. Otherwise, the risk of combinatorial explosion could severely affect the computational efficiency to the point where the simulation would come to a halt [45].

The nature of nonsmooth dynamics requires the adoption of a deeper mathematical framework, where concepts like set-valued functions, inclusions, and complementarity conditions are used [33]. In particular, recent time-stepping approaches construct weak solu-

tions of the differential variational inequality (DVI) that describes the continuous time motion of rigid bodies with collision, contact, and friction. Earlier numerical methods based on differential variational inequalities can be found in [26,25,24], whereas the DVI formulation has been discussed in full generality and classified by differential index only recently, in [37,29].

Two main families of solvers spawn from the DVI formulation: those that lead to an acceleration-force complementarity problem [9,31,45] and that generate velocity-impulse, complementarity-based time-stepping methods [39,6,7]. The latter case results in schemes convergent to a vector measure differential inclusion, so named because it operates on vector measures or distributions [40]. It has the advantage that it can solve a class of problems with Coulomb friction that would be unsolvable in an acceleration-force context, as the Painlevé paradox [38].

In both cases, the introduction of inequalities in time-stepping schemes for DVI, together with a polyhedral approximation of the friction cone as a faceted pyramid, leads to linear complementarity problems (LCPs) [40], which are systems of complementary inequalities to be satisfied simultaneously [14]. Such LCPs, which are hard to solve because of their inherent nonlinear nature, must be solved at each time step in order to advance the integrator [24,40].

Most literature about this topic shows how, for a large number of contacts and rigid bodies, usual LCP solution schemes have significant limitations. In fact, classical approaches to the solution of LCP problems are based on *simplex methods*, also known as *direct* or *pivoting* methods, originating from the algorithms of Lemke and Dantzig [13]. These methods may exhibit an exponential worst-case complexity [10]. Our experience shows that, in spite of algorithmic optimizations [43], simplex methods still cannot practically handle multibody systems with more than one hundred colliding bodies.

Moreover, in the three-dimensional case, typical LCP solvers can be used only at the cost of approximating the Coulomb friction cone with faceted pyramids [40,45,6]. Not only does this expedient introduce artificial anisotropy in the friction phenomenon; it also impacts negatively the performance of LCP solvers, which is already critical in general, because the finite approximation of cones results in a much larger problem.

A precise description of the friction cone constraint in three-dimensional space would imply a nonlinear complementarity problem. This is a broader class of problems in mathematical programming, for which no off-the-shelf solvers are available. A custom method must be developed.

The above-mentioned limitations of the existing LCP approaches led us to develop a novel solution method based on a fixed-point iteration with projection on a convex set and presented in [8]. That method extended the seminal work on iterative LCP solvers by Mangasarian [28] to the LCP case with conical constraints, that is, a cone complementarity problem. In the same work we presented the convergence theory

for the iteration; the scheme converges under certain conditions that do not include a small friction assumption. Applied to granular flow problems, our method demonstrated high performance and was able to solve benchmarks with up to a million dual variables.

The time-stepping scheme was proven to converge in a measure differential inclusion sense to the solution of the original continuous-time DVI [2].

In the present paper we extend our original formulation [8] in several ways.

- (i) We enhance our approach to the case of both frictional contacts and bilateral constraints, either scleronomic or rheonomic (motors, imposed trajectories, etc.). This extension cannot be obtained with full theoretical guarantees of convergence for the algorithm in [8] by simply replacing the bilateral constraint with two unilateral constraints. The latter case allows for unbounded internal forces, for which the approach in [8, Corollary 1] does not apply (the resulting constraint cone is not *pointed*).
- (ii) We present practical algorithmic details and optimizations that can be adopted to implement the method in a matrix-free, memory-efficient, reliable, fast, and robust way.
- (iii) We demonstrate the performance of our approach for configurations that include both joint and contact-with-friction constraints.

A significant side-effect of the proposed method is that it proceeds monotonically toward the solution. If implemented in real-time applications such as virtual-reality and man-in-the-loop vehicle simulations, where the requirements on precision are less severe than those on the computational times, it can be stopped prematurely before the tolerance threshold is reached.

We believe that the CAD community will welcome the availability of this solver because of its ability of simulating generic mechanisms regardless of the number of parts, joints, and frictional contacts. To this end we developed a physics library based on this time-stepping method, written in C++ and called Chrono::Engine [41], that can be used by third parties to develop simulation software. Indeed, we used it to implement a 3D graphical interface for the interactive modeling and visualization of multibody systems. Moreover, we have already simulated many types of mechanical problems, ranging from robots to granular flows with hundreds of thousands of rigid bodies with friction.

Among the most complex tests, we simulated the granular flow of a fourth-generation pebble bed nuclear reactor. Thanks to the numerical scheme, the refueling motion of 170,000 uranium spheres was simulated on a single computer, whereas DEM methods required a supercomputer and much more CPU time [42].

## 2. The model

This section presents a formulation for the nonsmooth dynamics of multibody systems in the most general case of both bilateral constraints and frictional



Fig. 1. Simulation of a palletizing device: a multibody problem with many frictional contacts.

contacts. We remark that frictional contacts embed also the case of unilateral constraints, since those represent a special case of contacts without friction.

## 2.1. System state

The position of the system at time  $t$  is represented by  $m_q$  generalized coordinates  $\mathbf{q}(t) \in \mathbb{R}^{m_q}$ . In the case of rigid bodies in three-dimensional space, these coordinates include the positions  $\mathbf{x} \in \mathbb{R}^3$  of the centers of mass of all bodies, as well as the rotations of all body frames respect to the absolute frame.

We represent rotations by means of unimodular quaternions  $\boldsymbol{\rho}(t) \in \mathbb{S}^3 \subset \mathbb{H}$ . Since quaternions are four-dimensional numbers, each rigid body in our formulation requires  $3+4=7$  scalar coordinates plus one constraint  $\|\boldsymbol{\rho}\| = 1$  enforcing the unit length of the quaternion<sup>1</sup>. With this notation, we define the position vector to be  $\mathbf{q} = \{\mathbf{x}^{1^T}, \boldsymbol{\rho}^{1^T}, \mathbf{x}^{2^T}, \boldsymbol{\rho}^{2^T}, \dots\}$ .

Generalized velocities are represented by the vector  $\mathbf{v}(t) \in \mathbb{R}^{m_v}$ , where for each body we consider the speed of the center of mass  $\dot{\mathbf{x}} \in \mathbb{R}^3$  and the angular velocity  $\boldsymbol{\omega}_i$ , expressed in local body coordinates. Therefore a system with  $n$  bodies in three dimensions is represented by  $m_v = 6n$  speed coordinates. With this notation, we define the speed vector to be  $\mathbf{v} = \{\dot{\mathbf{x}}^{1^T}, \boldsymbol{\omega}_1^{1^T}, \dot{\mathbf{x}}^{2^T}, \boldsymbol{\omega}_2^{2^T}, \dots\}^T$ .

Given the angular velocity  $\boldsymbol{\omega}_i$ , one can obtain the time derivative  $\dot{\boldsymbol{\rho}}$  of the quaternion, if needed, by building the purely imaginary quaternion  $\{0, \boldsymbol{\omega}_i\}$  and computing the quaternion product  $\dot{\boldsymbol{\rho}} = \frac{1}{2}\boldsymbol{\rho}\{0, \boldsymbol{\omega}_i\}$ . On this basis, we introduce  $\dot{\mathbf{q}} = \Gamma(\mathbf{q}, \mathbf{v})$  as the linear map that gives the time derivative of the position:

$$\begin{aligned} \Gamma(\mathbf{q}, \mathbf{v}) = \dot{\mathbf{q}} &= \{\dot{\mathbf{x}}^{1^T}, \dot{\boldsymbol{\rho}}^{1^T}, \dot{\mathbf{x}}^{2^T}, \dot{\boldsymbol{\rho}}^{2^T}, \dots\} = \\ &= \{\dot{\mathbf{x}}^{1^T}, \frac{1}{2}\boldsymbol{\rho}^1\{0, \boldsymbol{\omega}_1^1\}^T, \dot{\mathbf{x}}^{2^T}, \frac{1}{2}\boldsymbol{\rho}^2\{0, \boldsymbol{\omega}_1^2\}^T, \dots\}^T. \end{aligned} \quad (1)$$

For the time integration of the system position, different options exist for the function  $\mathbf{q}^{(t+\Delta t)} = \Lambda(\mathbf{q}^{(t)}, \mathbf{v}, \Delta t)$ , the simplest one being the first-order explicit Euler formula  $\mathbf{q}^{(t+\Delta t)} = \mathbf{q}^{(t)} + \Delta t \dot{\mathbf{q}}^{(t)}$ . Indeed, for body positions a straightforward first-order differential approximation is used:  $\mathbf{x}^{(t+\Delta t)} = \mathbf{x}^{(t)} + \Delta t \dot{\mathbf{x}}^{(t)}$ . However, if the same approach is used for quaternions, as in  $\boldsymbol{\rho}^{(t+\Delta t)} = \boldsymbol{\rho}^{(t)} + \Delta t \dot{\boldsymbol{\rho}}^{(t)}$ , an annoying situation can happen: quaternions may slowly lose the unimodularity and drift away from the  $\mathbb{S}^3$  manifold, unless some stabilization is used to enforce  $\|\boldsymbol{\rho}\| = 1$ . Therefore, we prefer to integrate the rotations using the following exponential map that preserves the unimodularity of the quaternions:

$$\boldsymbol{\rho}^{(t+\Delta t)} = \boldsymbol{\rho}^{(t)} e^{\{0, \frac{1}{2}\boldsymbol{\omega}_i \Delta t\}}, \quad (2)$$

so we get

$$\boldsymbol{\rho}^{(t+\Delta t)} = \Lambda(\mathbf{q}^{(t)}, \mathbf{v}, \Delta t) = \left\{ \begin{array}{l} \mathbf{x}^{1,(t)} + \Delta t \dot{\mathbf{x}}^{1,(t)} \\ \boldsymbol{\rho}^{1,(t)} e^{\{0, \frac{1}{2}\boldsymbol{\omega}_1^1 \Delta t\}} \\ \mathbf{x}^{2,(t)} + \Delta t \dot{\mathbf{x}}^{2,(t)} \\ \boldsymbol{\rho}^{2,(t)} e^{\{0, \frac{1}{2}\boldsymbol{\omega}_1^2 \Delta t\}} \\ \dots \end{array} \right\}. \quad (3)$$

We can explicitly compute the second factor of the quaternion product thanks to the property  $e^{\{0, \mathbf{u}\alpha\}} = \{\cos \alpha, \mathbf{u} \sin \alpha\}$  of quaternion exponentials; we obtain

$$\begin{aligned} e^{\{0, \frac{1}{2}\boldsymbol{\omega}_i \Delta t\}} &= e^{\{0, (\boldsymbol{\omega}_i / |\boldsymbol{\omega}_i|) \frac{1}{2} |\boldsymbol{\omega}_i| \Delta t\}} = \\ &= \left\{ \cos \frac{1}{2} |\boldsymbol{\omega}_i| \Delta t, \frac{\boldsymbol{\omega}_i}{|\boldsymbol{\omega}_i|} \sin \frac{1}{2} |\boldsymbol{\omega}_i| \Delta t \right\}. \end{aligned} \quad (4)$$

Since (2) preserves the norm of the quaternions, large  $\Delta t$  time steps can be used (although, once in a while, it is safer to normalize all quaternions because numerical roundoff can accumulate small errors). Yet we can demonstrate that, for  $\Delta t \rightarrow 0$ , the formula (2) still corresponds to  $\boldsymbol{\rho}^{(t+\Delta t)} = \boldsymbol{\rho}^{(t)} + \Delta t \dot{\boldsymbol{\rho}}$ . In fact,

$$\dot{\boldsymbol{\rho}} = \lim_{\Delta t \rightarrow 0} \frac{\boldsymbol{\rho}^{(t+\Delta t)} - \boldsymbol{\rho}^{(t)}}{\Delta t}. \quad (5)$$

Hence, substituting (2) in (5), we can write

$$\dot{\boldsymbol{\rho}} = \lim_{\Delta t \rightarrow 0} \boldsymbol{\rho}^{(t)} \frac{\{\cos \frac{1}{2} |\boldsymbol{\omega}_i| \Delta t, \frac{\boldsymbol{\omega}_i}{|\boldsymbol{\omega}_i|} \sin \frac{1}{2} |\boldsymbol{\omega}_i| \Delta t\} - \{1, \mathbf{0}\}}{\Delta t}. \quad (6)$$

Applying the Hôpital theorem and simplifying, we obtain  $\dot{\boldsymbol{\rho}} = \frac{1}{2}\boldsymbol{\rho}\{0, \boldsymbol{\omega}_i\}$ , as expected.

## 2.2. Bilateral constraints

Most kinematic pairs, such as revolute joints, prismatic joints, and glyphs, can be expressed by means of holonomic constraints over the relative position of two bodies. In general, we introduce a set  $\mathcal{G}_B$  of scalar equations

$$\Psi^i(\mathbf{q}, t) = 0, \quad i \in \mathcal{G}_B. \quad (7)$$

<sup>1</sup> Different methods can be used to represent the rotation. If one stores the  $3 \times 3$  rotation matrix  $A \in \text{SO}(3, \mathbb{R})$ , each body will require  $3+9=12$  scalar coordinates, which are highly redundant. On the other hand, storing only three angles (such as the three Euler angles or the three Cardano angles) could give problems of singularities. Not being affected by these limitations, quaternions are better suited for computational application.

The size of the set  $\mathcal{G}_B$  is the number of basic scalar bilateral constraints, not necessarily corresponding to complex three-dimensional mechanical joints<sup>2</sup>.

We assume that  $\Psi^i(\mathbf{q}, t)$  is smooth, so that it can be differentiated to obtain the Jacobian  $\nabla_{\mathbf{q}}\Psi^i = [\partial\Psi^i/\partial\mathbf{q}]^T$ .

Constraints must be respected also at the velocity level: the full time-derivative of the  $i$ th constraint equation is

$$\begin{aligned}\frac{d\Psi^i(\mathbf{q}, t)}{dt} &= \frac{\partial\Psi^i}{\partial\mathbf{q}}\dot{\mathbf{q}} + \frac{\partial\Psi^i}{\partial t} = \nabla_{\mathbf{q}}\Psi^{iT}\dot{\mathbf{q}} + \frac{\partial\Psi^i}{\partial t} = 0 \\ \frac{d\Psi^i(\mathbf{q}, t)}{dt} &= \nabla_{\mathbf{q}}\Psi^{iT}\Gamma(\mathbf{q}, \mathbf{v}) + \frac{\partial\Psi^i}{\partial t} = 0.\end{aligned}$$

For simplicity, from now on we will define  $\nabla\Psi^{iT} = \nabla_{\mathbf{q}}\Psi^{iT}\Gamma(\mathbf{q}, \mathbf{v})$ .

Note that the term  $\frac{\partial\Psi^i}{\partial t}$  is nonzero only for rheonomic (time-dependent) constraints such as motors and imposed trajectories.

For each bilateral constraint, there exists a Lagrange multiplier  $\gamma_B^i$  such that the force acting on the system by the  $i$ th bilateral constraint is  $\gamma_B^i\nabla\Psi^i$  [20].

### 2.3. Unilateral contact constraints

Since rigid bodies cannot overlap each other, given the set of body shapes  $\Omega = \{\Omega^1, \Omega^2, \dots, \Omega^n\}$ , we assume that there exists a set of  $\mathcal{G}_P$  distance functions  $\Phi(\mathbf{q}, \Omega)$  that must satisfy the unilateral constraint conditions:

$$\Phi^i(\mathbf{q}, \Omega) \geq 0, \quad i \in \mathcal{G}_P. \quad (8)$$

An example of such a mapping is the signed distance function [23].

For convenience, we pose the problem in terms of contact points, since in most cases we can compute a minimal set of contact normals belonging to one of two neighboring bodies: distances  $\Phi(\mathbf{q}, \Omega)$  are measured along those normals.

For example, in the case of spherical bodies with positions  $\mathbf{x}_s^j$  and radius  $r_s$ , for all body pairs  $j, k$  the signed distance function is  $\Phi = \|\mathbf{x}_s^j - \mathbf{x}_s^k\| - 2r_s$ . Note, however, that, for bodies with generic shapes, finding a proper set of contact points (and defining their  $\Phi^i$  distance functions) is not always trivial [3,4]. In fact, there could be multiple contact points, or it could even happen that defining a differentiable signed distance function is not possible, as in the case of concave shapes [3].

Nevertheless, since we are interested in enforcing non-penetration, what truly matters is that a signed distance function be defined up to some value of the penetration [4]. We thus assume that  $\Phi(\mathbf{q}, \Omega)$  can be differentially defined at least on a neighborhood of the set  $\Phi(\mathbf{q}, \Omega) \geq 0$ . Such an assumption does hold for smooth and strictly convex bodies, such as spheres [3].

<sup>2</sup> In the context of this work, bilateral constraints are always considered scalar, because complex mechanical joints can be modeled by using multiple basic scalar constraints. For example, kinematic pairs such as a ball joint require three scalar equations, a prismatic guide requires five scalar equations, and so on.

In addition, piecewise smooth bodies can be accommodated in a fixed-time-step framework, by decomposing the distance function in components attached to each pair of *features*, such as point and piecewise smooth surface or curved edge and curved edge, and using all normals attached to such pairs [19].

Also, we consider only a subset  $\mathcal{G}_A(\mathbf{q}, \Omega, \epsilon) \subset \mathcal{G}_P$  of all potential contacts, that is, only those contacts whose surfaces are under a distance threshold  $\epsilon$ :

$$\mathcal{G}_A(\mathbf{q}, \Omega, \epsilon) = \{i \mid i \in \mathcal{G}_P, \Phi^i(\mathbf{q}, \Omega) \leq \epsilon\}. \quad (9)$$

Special attention must be paid in implementing an efficient and robust collision algorithm for the generation of the  $\mathcal{G}_A(\mathbf{q}, \Omega, \epsilon)$  set of contact points.

A preliminary algorithm, called *broad-phase* collision detection, discards pairs of shapes that are farther than  $\epsilon$ , in order to avoid a combinatorial waste of time if checking for collision points with all pairs of bodies. We use the *sweep and prune* (SAP) algorithm to this end [17].

The following *narrow-phase* step operates on the pairs of bodies that passed the broad-phase check: it finds the contact points and their normals. We adopt the GJK algorithm for this purpose because it features high efficiency and robustness even in the case of non-smooth surfaces [16].

Concave shapes, if any, undergo an off-line conversion into sets of convex shapes using a convex decomposition algorithm; a *middle-phase* AABB binary-tree traversal is used to check collisions between these compounds of shapes without running into superlinear time complexity.

A thin envelope is added around all shapes using the Minkowski sums in order to allow a small amount of interpenetration. If larger overlapping occurs, the Expanding Polytope (EPA) algorithm is used [11].

### 2.4. Frictional constraints

In the following section we introduce friction by means of conic constraints, which are an extension of complementarity models discussed in [6,40].

#### 2.4.1. The Coulomb friction model

The original Coulomb model introduces static  $\mu_s$  and kinetic  $\mu_k$  friction coefficients as the only parameters to characterize the frictional phenomena at the surface. Although simple, this model was proven to be realistic and practical in many situations. Usually the kinetic coefficient is slightly lower than the static coefficient, but in this work we consider both to have the same value  $\mu$ .

If a position  $\mathbf{q}$  is feasible and the contact is active, that is,  $\Phi(\mathbf{q}, \Omega) = 0$ , then at the contact we have a normal force and a tangential force.

Let  $\mathbf{n}$  be the normal at the contact pointing from the second body to the first body, and let  $\mathbf{t}_1$  and  $\mathbf{t}_2$  be the tangents at the contact. Here  $\mathbf{n}, \mathbf{t}_1, \mathbf{t}_2$  are mutually orthogonal vectors of length one in three dimensions. The vectors  $\mathbf{n}, \mathbf{t}_1$ , and  $\mathbf{t}_2$  are a function of the position  $\mathbf{q}$ , but we ignore this fact until the end of this section.

The reaction force is impressed on the system by means of multipliers  $\hat{\gamma}_n \geq 0$ ,  $\hat{\gamma}_u$ , and  $\hat{\gamma}_v$ . The normal component of the force is  $\mathbf{F}_N = \hat{\gamma}_n \mathbf{n}$ , and the tangential component of the force is  $\mathbf{F}_T = \hat{\gamma}_u \mathbf{t}_1 + \hat{\gamma}_v \mathbf{t}_2$ .

The Coulomb model consists of the following constraints:

$$\begin{aligned} \hat{\gamma}_n &\geq 0, & \Phi(\mathbf{q}) &\geq 0, & \Phi(\mathbf{q})\hat{\gamma}_n &= 0, \\ \mu\hat{\gamma}_n &\geq \sqrt{\hat{\gamma}_u^2 + \hat{\gamma}_v^2}, & \|\mathbf{v}_T\| &\left( \mu\hat{\gamma}_n - \sqrt{\hat{\gamma}_u^2 + \hat{\gamma}_v^2} \right) &= 0, \\ \langle \mathbf{F}_T, \mathbf{v}_T \rangle &= -\|\mathbf{F}_T\| \|\mathbf{v}_T\|, \end{aligned} \quad (10)$$

where  $\mathbf{v}_T$  is the relative tangential velocity at contact. The effect of the friction over the dynamical system is defined by the friction coefficient  $\mu \in \mathbb{R}^+$ , which typically has a value between 0 and 1 for most materials.

The first part of the constraint can be restated as

$$\mathbf{F} = \mathbf{F}_N + \mathbf{F}_T = \hat{\gamma}_n \mathbf{n} + \hat{\gamma}_u \mathbf{t}_1 + \hat{\gamma}_v \mathbf{t}_2 \in \mathcal{K},$$

where  $\mathcal{K}$  is a cone in three dimensions, whose slope is  $\arctan(\mu)$ .

The constraint  $\langle \mathbf{F}_T, \mathbf{v}_T \rangle = -\|\mathbf{F}_T\| \|\mathbf{v}_T\|$  requires that the tangential force be opposite to the tangential velocity. As a result, the reaction force is dissipative. In fact, an equivalent convenient way of expressing this constraint is by using the maximum dissipation principle [40,38,39],

$$(\hat{\gamma}_u, \hat{\gamma}_v) = \operatorname{argmin}_{\sqrt{\hat{\gamma}_u^2 + \hat{\gamma}_v^2} \leq \mu \hat{\gamma}_n} (\hat{\gamma}_u \mathbf{t}_1 + \hat{\gamma}_v \mathbf{t}_2)^T \mathbf{v}_T.$$

These constraints are represented by mapping the vectors  $\mathbf{n}, \mathbf{t}_1, \mathbf{t}_2$  from contact coordinates to generalized coordinates [3].

We denote the generalized vector version of  $\mathbf{n}, \mathbf{t}_1, \mathbf{t}_2$  by  $\mathbf{D}_n, \mathbf{D}_u, \mathbf{D}_v$ .

In generalized coordinates, the Coulomb model becomes [8]

$$\mathbf{F}_N = \hat{\gamma}_n \mathbf{D}_n, \quad \mathbf{F}_T = \hat{\gamma}_u \mathbf{D}_u + \hat{\gamma}_v \mathbf{D}_v, \quad (11)$$

$$\hat{\gamma}_n \geq 0, \quad \Phi(\mathbf{q}) \geq 0, \quad \hat{\gamma}_n \Phi(\mathbf{q}) = 0, \quad (12)$$

where the tangential multipliers  $\hat{\gamma}_u, \hat{\gamma}_v$  are determined from the maximum dissipation principle

$$(\hat{\gamma}_u, \hat{\gamma}_v) = \operatorname{argmin}_{\sqrt{\hat{\gamma}_u^2 + \hat{\gamma}_v^2} \leq \mu \hat{\gamma}_n} (\hat{\gamma}_u \mathbf{D}_u + \hat{\gamma}_v \mathbf{D}_v)^T \mathbf{v}.$$

The last relation is obtained from the identities  $\mathbf{D}_u^T \mathbf{v} = \mathbf{t}_1^T \mathbf{v}_T$  and  $\mathbf{D}_v^T \mathbf{v} = \mathbf{t}_2^T \mathbf{v}_T$ .

## 2.5. The Overall dynamical model

The other dynamical data needed for the model are the mass matrix  $M(\mathbf{q})$ , which is symmetric positive definite; the external force  $\mathbf{f}_e(t, \mathbf{q}, \mathbf{v})$ ; and the inertial force  $\mathbf{f}_c(\mathbf{q}, \mathbf{v})$ , containing the centrifugal and Coriolis forces.

We can define the total force

$$\mathbf{f}_t(t, \mathbf{q}, \mathbf{v}) = \mathbf{f}_e(t, \mathbf{q}, \mathbf{v}) + \mathbf{f}_c(\mathbf{q}, \mathbf{v}). \quad (13)$$

Assume now that we have multiple contact constraints  $\Phi^i(\mathbf{q}, \Omega) \geq 0$ ,  $i \in \mathcal{G}_A$  and multiple bilateral constraints  $\Psi^i(\mathbf{q}, t) = 0$ ,  $i \in \mathcal{G}_B$ . Note that the unilateral condition  $\hat{\gamma}_n \geq 0, \Phi \geq 0, \Phi \hat{\gamma}_n = 0$  can be written as a complementarity constraint  $\hat{\gamma}_n \geq 0 \perp \Phi \geq 0$ .

The continuous model is a differential variational inequality [37]:

$$\begin{aligned} M(\mathbf{q}^l) \frac{d\mathbf{v}}{dt} &= \sum_{i \in \mathcal{G}_A} (\hat{\gamma}_n^i \mathbf{D}_n^i + \hat{\gamma}_u^i \mathbf{D}_u^i + \hat{\gamma}_v^i \mathbf{D}_v^i) + \\ &+ \sum_{i \in \mathcal{G}_B} \hat{\gamma}_B^i \nabla \Psi^i + \mathbf{f}_t(t, \mathbf{q}, \mathbf{v}) \end{aligned}$$

$$\dot{\mathbf{q}} = \Gamma(\mathbf{q}, \mathbf{v})$$

$$\Psi^i(\mathbf{q}, t) = 0 \quad i \in \mathcal{G}_B$$

$$\hat{\gamma}_n^i \geq 0 \quad \perp \quad \Phi^i(\mathbf{q}, \Omega) \geq 0, \quad i \in \mathcal{G}_A$$

$$\begin{aligned} (\hat{\gamma}_u^i, \hat{\gamma}_v^i) &= \operatorname{argmin}_{\mu \hat{\gamma}_n^i \geq \sqrt{(\hat{\gamma}_u^i)^2 + (\hat{\gamma}_v^i)^2}} \\ &\mathbf{v}^T (\hat{\gamma}_u \mathbf{D}_u^i + \hat{\gamma}_v \mathbf{D}_v^i). \end{aligned} \quad (14)$$

Unfortunately, the introduction of the Coulomb friction model may lead to an inconsistent model. It is known [9] that paradoxical configurations exist for which such a model does not have a solution in terms of unknown accelerations and reaction forces. Such configurations are called Painlevé paradoxes [39]. Nevertheless, a weaker formulation of the problem can be solved in terms of vector measures, using a nonsmooth time-stepping scheme where reaction impulses are the unknowns at each time step [39].

To this end we define the following stepping scheme, with time step  $h$ , known positions  $\mathbf{q}^{(l)}$ , and velocity  $\mathbf{v}^{(l)}$ ; the scheme is an equation problem with equilibrium constraints, where the unknowns are  $\mathbf{q}^{(l+1)}$ ,  $\mathbf{v}^{(l+1)}$ , and constraint impulses  $\gamma_n = h\hat{\gamma}_n$ ,  $\gamma_u = h\hat{\gamma}_u$ ,  $\gamma_v = h\hat{\gamma}_v$ ,  $\gamma_B = h\hat{\gamma}_B$ :

$$\begin{aligned} M^{(l)}(\mathbf{v}^{(l+1)} - \mathbf{v}^{(l)}) &= \sum_{i \in \mathcal{G}_A} (\gamma_n^i \mathbf{D}_n^i + \gamma_u^i \mathbf{D}_u^i + \gamma_v^i \mathbf{D}_v^i) + \\ &+ \sum_{i \in \mathcal{G}_B} (\gamma_B^i \nabla \Psi^i) + h \mathbf{f}_t(t^{(l)}, \mathbf{q}^{(l)}, \mathbf{v}^{(l)}) \end{aligned} \quad (15)$$

$$0 = \frac{1}{h} \Psi^i(\mathbf{q}^{(l)}) + \nabla \Psi^i \mathbf{v}^{(l+1)} + \frac{\partial \Psi^i}{\partial t}, \quad i \in \mathcal{G}_B$$

$$0 \leq \frac{1}{h} \Phi^i(\mathbf{q}^{(l)}) + \nabla \Phi^i \mathbf{v}^{(l+1)} \quad (17)$$

$$\perp \quad \gamma_n^i \geq 0, \quad i \in \mathcal{G}_A$$

$$(\gamma_u^i, \gamma_v^i) = \operatorname{argmin}_{\mu \gamma_n^i \geq \sqrt{(\gamma_u^i)^2 + (\gamma_v^i)^2}} \quad i \in \mathcal{G}_A$$

$$[\mathbf{v}^T (\gamma_u \mathbf{D}_u^i + \gamma_v \mathbf{D}_v^i)] \quad (18)$$

$$\mathbf{q}^{(l+1)} = \Lambda(\mathbf{q}^{(l)}, \mathbf{v}^{(l+1)}, h). \quad (19)$$

To simplify notation, we denoted  $M(\mathbf{q}^l)$  by  $M^{(l)}$ .

In previous work, we have shown that the scheme is convergent, as the time step  $h$  goes to 0, to the solution of a measure differential inclusion [2].

For the special case of zero friction, the subproblem simplifies to (15-17), that is, a linear complementarity problem. Such problems can be solved by Lemke's algorithm [14,6]. Introducing the Coulomb friction (18), however, turns the problem into a nonlinear complementarity problem that poses more difficulties. If the nonlinear constraint cone (the Coulomb cone) is approximated by a piecewise linear cone, the subproblem (15-18) becomes again an LCP solvable by Lemke's algorithm [6]. Nevertheless, in [5] we have also demonstrated that, as the number of constraints in the prob-

lem increases, the computational cost of typical LCP solvers increases far faster than linearly with the size of the problem. Moreover, the approximation of friction cones by means of faceted pyramids would introduce unwanted anisotropy.

To overcome these difficulties, we modified the time-stepping scheme by relaxing the constraint (17) as

$$0 \leq \frac{1}{h} \Phi^i(\mathbf{q}^{(l)}) + \nabla \Phi^{i^T} \mathbf{v}^{(l+1)} - \mu^i \sqrt{(\mathbf{D}_u^{i,T} \mathbf{v})^2 + (\mathbf{D}_v^{i,T} \mathbf{v})^2} \perp \gamma_n^i \geq 0, \quad i \in \mathcal{G}_A. \quad (20)$$

This results in a cone complementarity problem that can be solved with a fixed-point iteration approach, as demonstrated in our earlier work [8].

## 2.6. Cone complementarity formulation

Developing the optimality conditions for the equilibrium constraint in (18), we obtain that there exists a Lagrange multiplier  $\lambda^i$  such that, for any  $i \in \mathcal{G}_A$ ,

$$\begin{aligned} \lambda^i \gamma_u^i &= -\mathbf{D}_u^{i,T} \mathbf{v}, \quad \lambda^i \gamma_v^i = -\mathbf{D}_v^{i,T} \mathbf{v}, \\ \lambda^i &\geq 0 \perp \mu^i \gamma_n^i - \sqrt{(\gamma_u^i)^2 + (\gamma_v^i)^2} \geq 0. \end{aligned} \quad (21)$$

The first two equations imply that  $\lambda^i \sqrt{(\gamma_u^i)^2 + (\gamma_v^i)^2} = \sqrt{(\mathbf{D}_u^{i,T} \mathbf{v})^2 + (\mathbf{D}_v^{i,T} \mathbf{v})^2}$ , while the complementarity constraint implies that

$$0 = \lambda^i \sqrt{(\gamma_u^i)^2 + (\gamma_v^i)^2} \left( \mu^i \gamma_n^i - \sqrt{(\gamma_u^i)^2 + (\gamma_v^i)^2} \right)$$

and, in turn, that

$$\mu^i \gamma_n^i \sqrt{(\mathbf{D}_u^{i,T} \mathbf{v})^2 + (\mathbf{D}_v^{i,T} \mathbf{v})^2} = \lambda^i \left( (\gamma_u^i)^2 + (\gamma_v^i)^2 \right). \quad (22)$$

We now define, for all potential contacts, the vectors

$$\mathbf{u}_A^i = \left\{ \frac{1}{h} \Phi^i(\mathbf{q}^{(l)}) + \nabla \Phi^{i^T} \mathbf{v}^{(l+1)}, \mathbf{D}_u^{i,T} \mathbf{v}, \mathbf{D}_v^{i,T} \mathbf{v} \right\}^T \quad (23)$$

$$\gamma_A^i = \{ \gamma_n^i, \gamma_u^i, \gamma_v^i \}^T, \quad i \in \mathcal{G}_A. \quad (24)$$

We calculate the scalar product using (20),(21):

$$\begin{aligned} \langle \mathbf{u}_A^i, \gamma_A^i \rangle &= \gamma_n^i \left( \frac{1}{h} \Phi^i + \nabla \Phi^{i^T} \mathbf{v} \right) + \gamma_u^i \mathbf{D}_u^{i,T} \mathbf{v} \\ &\quad + \gamma_v^i \mathbf{D}_v^{i,T} \mathbf{v} \\ &= \mu^i \gamma_n^i \sqrt{(\mathbf{D}_u^{i,T} \mathbf{v})^2 + (\mathbf{D}_v^{i,T} \mathbf{v})^2} \\ &\quad - \lambda^i \left( (\gamma_u^i)^2 + (\gamma_v^i)^2 \right) \\ &= 0 \quad \Rightarrow \quad \mathbf{u}_A^i \perp \gamma_A^i. \end{aligned} \quad (25)$$

We recall that the *dual cone* of a convex cone  $\mathcal{K}$  is the set  $\mathcal{K}^* = \{ \mathbf{x} | \forall \mathbf{y} \in \mathcal{K} \langle \mathbf{y}, \mathbf{x} \rangle \geq 0 \}$  and that the *polar cone* is defined as  $\mathcal{K}^\circ = -\mathcal{K}^*$ .

We now define the friction cone  $\mathcal{FC}^i$  such that the Coulomb friction model is satisfied if  $\gamma_A^i \in \mathcal{FC}^i$ :

$$\mathcal{FC}^i = \left\{ x, y, z \in \mathbb{R}^3 | \mu^i x \geq \sqrt{y^2 + z^2} \right\}.$$

Then, from (18), (20), and (25), the frictional contact constraints can be expressed by means of the following *cone complementarity constraints*:

$$-\mathbf{u}_A^i \in \mathcal{FC}^{i^\circ} \perp \gamma_A^i \in \mathcal{FC}^i, \quad i \in \mathcal{G}_A. \quad (26)$$

To obtain a unified formalism, we can represent also the bilateral constraints (16) in terms of cone complementarity constraints. Of course, multipliers  $\gamma_B^i$ , with  $i \in \mathcal{G}_B$ , are not restrained into some special subset of  $\mathbb{R}$ , but even  $\mathbb{R}$  itself is a convex cone. Thus we can introduce the scalar

$$u_B^i = \frac{1}{h} \Psi^i(\mathbf{q}^{(l)}) + \nabla \Psi^{i^T} \mathbf{v}^{(l+1)} + \frac{\partial \Psi}{\partial t}, \quad i \in \mathcal{G}_B, \quad (27)$$

which allows us to write the bilateral constraints (16) as

$$-u_B^i \in \mathcal{BC}^{i^\circ} \perp \gamma_B^i \in \mathcal{BC}^i, \quad i \in \mathcal{G}_B, \quad (28)$$

where  $\mathcal{BC}^i = \{ \mathbb{R} \}$  and  $\mathcal{BC}^{i^\circ} = \{ 0 \}$ , and  $\langle u_B^i, \gamma_B^i \rangle = 0$  is always satisfied for  $-u_B^i \in \mathcal{BC}^{i^\circ}$ .

We now define the vector

$$\tilde{\mathbf{k}}^{(l)} = M^{(l)} \mathbf{v}^{(l)} + h \mathbf{f}_i(t^{(l)}, \mathbf{q}^{(l)}, \mathbf{v}^{(l)}). \quad (29)$$

Then, equations (29), (28), and (26), together with (15) and the definition of  $\mathbf{u}_A^i$ ,  $\gamma_A^i$ ,  $u_B^i$ , and  $\gamma_B^i$ , result in the following problem:

$$\begin{aligned} M^{(l)} \mathbf{v}^{(l+1)} &= \sum_{i \in \mathcal{G}_A} (\gamma_n^i \mathbf{D}_n^i + \gamma_u^i \mathbf{D}_u^i + \gamma_v^i \mathbf{D}_v^i) + \\ &\quad + \sum_{i \in \mathcal{G}_B} (\gamma_B^i \nabla \Psi^i) + \tilde{\mathbf{k}}^{(l)}, \\ -\mathbf{u}_A^i &\in \mathcal{FC}^{i^\circ} \perp \gamma_A^i \in \mathcal{FC}^i, \quad i \in \mathcal{G}_A \\ -u_B^i &\in \mathcal{BC}^{i^\circ} \perp \gamma_B^i \in \mathcal{BC}^i, \quad i \in \mathcal{G}_B. \end{aligned} \quad (30)$$

If we want to obtain a cone complementarity problem as expressed in the typical form  $-\mathcal{K}^\circ \ni \mathbf{f}(\mathbf{a}) \perp \mathbf{a} \in \mathcal{K}$ , a more compact formulation of the problem (30) is necessary. To this end we denote by  $n_A$  and  $n_B$  the number of elements in the sets  $\mathcal{G}_A$  and  $\mathcal{G}_B$ , respectively. Then, we define the following vectors  $\mathbf{b}_A \in \mathbb{R}^{3n_A}$ ,  $\gamma_A \in \mathbb{R}^{3n_A}$ ,  $\mathbf{b}_B \in \mathbb{R}^{n_B}$ , and  $\gamma_B \in \mathbb{R}^{n_B}$ :

$$\begin{aligned} \mathbf{b}_A &= \left\{ \frac{1}{h} \Phi^{i_1}, 0, 0, \frac{1}{h} \Phi^{i_2}, 0, 0, \dots, \frac{1}{h} \Phi^{i_{n_A}}, 0, 0 \right\}^T \\ \gamma_A &= \left\{ \gamma_n^{i_1}, \gamma_u^{i_1}, \gamma_v^{i_1}, \gamma_n^{i_2}, \gamma_u^{i_2}, \gamma_v^{i_2}, \dots, \gamma_n^{i_{n_A}}, \gamma_u^{i_{n_A}}, \gamma_v^{i_{n_A}} \right\}^T \\ \mathbf{b}_B &= \left\{ \frac{1}{h} \Psi^1 + \frac{\partial \Psi^1}{\partial t}, \frac{1}{h} \Psi^2 + \frac{\partial \Psi^2}{\partial t}, \dots, \frac{1}{h} \Psi^{n_B} + \frac{\partial \Psi^{n_B}}{\partial t} \right\}^T \\ \gamma_B &= \{ \gamma_B^1, \gamma_B^2, \dots, \gamma_B^{n_B} \}^T \\ \mathbf{u}_A &= \left\{ \mathbf{u}_A^1, \mathbf{u}_A^2, \dots, \mathbf{u}_A^{n_A} \right\}^T, \quad \mathbf{u}_B = \{ u_B^1, u_B^2, \dots, u_B^{n_B} \}^T. \end{aligned} \quad (31)$$

It is useful to merge these vectors, joining data from both frictional constraints and bilateral constraints, obtaining vectors with  $n_E = 3n_A + n_B$  scalar elements:

$$\mathbf{b}_E = \left\{ \mathbf{b}_A^T, \mathbf{b}_B^T \right\}^T, \quad \gamma_E = \left\{ \gamma_A^T, \gamma_B^T \right\}^T, \quad \mathbf{u}_E = \left\{ \mathbf{u}_A^T, \mathbf{u}_B^T \right\}^T. \quad (32)$$

For each frictional contact  $i \in \mathcal{G}_A$  we also define the following three-column matrix:

$$D^i = \left[ \mathbf{D}_n^i | \mathbf{D}_u^i | \mathbf{D}_v^i \right]. \quad (33)$$

As before, it is useful to merge all Jacobians from both frictional constraints and bilateral constraints in a single, large matrix having  $n_E$  columns and  $m_v$  rows:

$$D_E = \left[ D^{i_1} | D^{i_2} | \dots | D^{i_{n_A}} | \nabla \Psi^1 | \nabla \Psi^2 | \dots | \nabla \Psi^{n_B} \right]. \quad (34)$$

From the definitions (23),(27), (31), (32), (33), and (34) one can see that

$$\mathbf{u}_\varepsilon = D_\varepsilon^T \mathbf{v}^{(l+1)} + \mathbf{b}_\varepsilon. \quad (35)$$

Also, premultiplying by  $M^{(l)-1}$  equation (15), one gets

$$\mathbf{v}^{(l+1)} = M^{(l)-1} D_\varepsilon \gamma_\varepsilon + M^{(l)-1} \tilde{\mathbf{k}}. \quad (36)$$

Hence it is possible to substitute (36) into (35) to obtain

$$\mathbf{u}_\varepsilon = D_\varepsilon^T M^{(l)-1} D_\varepsilon \gamma_\varepsilon + D_\varepsilon^T M^{(l)-1} \tilde{\mathbf{k}} + \mathbf{b}_\varepsilon. \quad (37)$$

To make the expressions more compact, we introduce the following:

$$N = D_\varepsilon^T M^{(l)-1} D_\varepsilon \quad (38)$$

$$\mathbf{r} = D_\varepsilon^T M^{(l)-1} \tilde{\mathbf{k}} + \mathbf{b}_\varepsilon \quad (39)$$

In this way, we can write

$$\mathbf{u}_\varepsilon = N \gamma_\varepsilon + \mathbf{r}. \quad (40)$$

Consider the multidimensional cone obtained by performing the direct sum of all  $\mathcal{FC}$  and  $\mathcal{BC}$  cones and its embedding in the corresponding vector space direct sums:

$$\Upsilon = \left( \bigoplus_{i \in \mathcal{G}_A} \mathcal{FC}^i \right) \oplus \left( \bigoplus_{i \in \mathcal{G}_B} \mathcal{BC}^i \right). \quad (41)$$

From (41), (40), (31), the complementarity relationship in (30), and the property of convex cones [22],  $\Upsilon = \bigoplus_i \Upsilon^i \Rightarrow \Upsilon^\circ = \bigoplus_i \Upsilon^{i,\circ}$ , we can write the problem as a cone complementarity problem (CCP):

$$(N \gamma_\varepsilon + \mathbf{r}) \in -\Upsilon^\circ \quad \perp \quad \gamma_\varepsilon \in \Upsilon. \quad (42)$$

The separable structure of the cone will allow us to define an algorithm based on block matrices, with relatively small blocks (dimension no larger than 3 in the case of the contact problem).

### 2.7. Physical effects of the relaxation

In [2] we demonstrated that, for  $h \rightarrow 0$ , the solution of the time-stepping scheme with the relaxed constraints (20) will approach the solution of the same measure differential inclusion as the scheme that uses the unrelaxed constraints (17). In addition, iterates produced by the modified scheme approach the ones of the original scheme *even for one time step* at fixed  $h$ , provided that  $\mu^i \gamma_n^i \sqrt{(\mathbf{D}_u^{i,T} \mathbf{v})^2 + (\mathbf{D}_v^{i,T} \mathbf{v})^2} \ll 1$ , that is, with either low friction or low tangential speed [5]. We note that this regime happens frequently in granular flow applications, such as in the simulation of the refueling of pebble bed nuclear reactors [34] and dense packing of granular matter. On the other hand, the good behavior of our scheme occurs even beyond this regime, for reasons we now explain.

Note that the  $\Phi/h$  term achieves constraint stabilization. When the term is positive and the constraint is active, it biases the normal impulse to be smaller than the one at exact contact and allows for ‘‘soft landing’’ projection onto the contact manifold. When the term

is negative, it biases the normal impulse to be larger than the one at exact contact and allows restoration to the contact manifold. The square root term in (20), which can be written also as  $\mu \|\mathbf{v}_T^i\|$ , does not appear in the original model and, as discussed in [2], does result in a larger value of the normal velocity. Therefore, it could potentially produce a departure from the prediction of the original model. As we discuss below, however, the constraint stabilization term substantially alleviates this effect. Indeed, the normal speed is  $\mathbf{v}_T^i = \nabla \Phi^{i,T} \mathbf{v}$ , so it follows that

$$\Phi^i(\mathbf{q}) = \mu^i \|\mathbf{v}_T^i\| h, \quad \mathbf{v}_N^i = 0, \quad (43)$$

satisfies the complementarity constraints (20) exactly. Hence, a solution to (20) that involves the term  $\mu^i \|\mathbf{v}_T^i\|$  need not result in an increase value of the normal velocity as long as the gap  $\Phi^i(\mathbf{q})$  is about  $\mu^i \|\mathbf{v}_T^i\| h$ . All the meaningful simulations that we have carried out indicate that this is indeed what happens, and thus the solution approaches the solution of the original scheme insofar as normal velocity. Hence, the model effectively includes a separating boundary layer of size  $\mu^i \|\mathbf{v}_T^i\| h$  that plays the role of zero-order effective compliance. Smaller values of  $\Phi^i(\mathbf{q})$  result in higher than steady-state impulse, whereas larger values of  $\Phi^i(\mathbf{q})$  result in smaller than steady-state impulse, both of which push the value of the gap function to the value  $\mu^i \|\mathbf{v}_T^i\| h$ . This effect is immediately seen for one contact, which leaves from rest with zero tangential velocity or which lands in an inelastic fashion. This is harder to prove for multiple contacts, though it always held in our experiments. There is one caveat, however, as shown in [2]. That is, if one starts with  $\Phi^i(\mathbf{q}) \ll \mu^i \|\mathbf{v}_T^i\| h$ , then the normal component of the velocity is much larger than the one predicted by the original scheme and then what is needed to satisfy approximately (43) in a few steps, so the contact suffers an artificial takeoff. This corresponds to the only case in which we have seen our argument fail, the one of high initial tangential velocity at a contact. Since this case occurs rarely in practice, our relaxation approximates in most cases the solution that would have been obtained by the original, unrelaxed scheme.

### 3. The iterative method

To solve the CCP (42), we propose a fixed-point iteration with the following form:

$$\begin{aligned} \gamma_\varepsilon^{r+1} &= \lambda \Pi_\Upsilon (\gamma_\varepsilon^r - \omega B^r (N \gamma_\varepsilon^r + \mathbf{r} + K^r (\gamma_\varepsilon^{r+1} - \gamma_\varepsilon^r))) + \\ &\quad + (1 - \lambda) \gamma_\varepsilon^r, \\ r &= 0, 1, 2, \dots \end{aligned} \quad (44)$$

This iteration uses block matrices  $B^r$  and  $K^r$ . Matrix  $B$  is null except for blocks on the diagonal; in our implementation blocks that correspond to the  $i$ th frictional contact are scaled identity matrices  $B_A^i = \eta_A^i I$ ,  $B_A^i \in \mathbb{R}^3$ , while blocks that correspond to the  $i$ th bilateral constraint are scalars  $B_B^i = \eta_B^i$ ,  $B_B^i \in \mathbb{R}$ . Matrix  $K$  is an upper or lower block matrix, with null blocks on the diagonal corresponding to the blocks of  $B$ ; in



$$\Pi_{\Upsilon} \left\{ \begin{array}{ll} \forall i \in \mathcal{G}_{\mathcal{B}} & \Pi_{\mathcal{BC}^i} = \gamma_{\mathcal{B}}^i \\ \forall i \in \mathcal{G}_{\mathcal{A}} & \Pi_{\mathcal{FC}^i} = \gamma_{\mathcal{A}}^i \\ \gamma_r < \mu_i \gamma_n & \Pi_{\mathcal{FC}^i} = \{0, 0, 0\}^T \\ \gamma_r < -\frac{1}{\mu_i} \gamma_n & \Pi_{\mathcal{FC}^i} = \{0, 0, 0\}^T \\ \gamma_r > \mu_i \gamma_n, \gamma_r > -\frac{1}{\mu_i} \gamma_n & \Pi_{\mathcal{FC}_n^i} = \frac{\gamma_r \mu_i + \gamma_n}{\mu_i^2 + 1} \\ & \Pi_{\mathcal{FC}_u^i} = \gamma_u \frac{\mu_i \Pi_{\mathcal{FC}_n^i}}{\gamma_r} \\ & \Pi_{\mathcal{FC}_v^i} = \gamma_v \frac{\mu_i \Pi_{\mathcal{FC}_n^i}}{\gamma_r}. \end{array} \right. \quad (46)$$

## 4. Implementation

The CCP method proposed here can be applied to the simulation of multibody systems with a large number of parts and contacts because, where an upper limit on the number of iteration is enforced, the iteration (44) can run in  $O(n)$  space and  $O(n)$  time.

Previous sections showed that generic multibody problems with frictional contacts, expressed with the system (15)–(19), embed the cone complementarity problem (42), which can be solved by the iterative method (44).

Given (31), one can consider the final time-stepping scheme as a sequence of three main operations: a CCP problem that finds unknown reactions  $\gamma_{\mathcal{E}}$  (47a), an affine scaling (47b) that gives the new speeds  $\mathbf{v}^{(l+1)}$ , and a position update (47c):

$$(N\gamma_{\mathcal{E}} + \mathbf{r}) \in -\Upsilon^o \perp \gamma_{\mathcal{E}} \in \Upsilon \quad (47a)$$

$$\mathbf{v}^{(l+1)} = M^{-1} (\tilde{\mathbf{k}} + D_{\mathcal{E}}\gamma_{\mathcal{E}}) \quad (47b)$$

$$\mathbf{q}^{(l+1)} = \Lambda(\mathbf{q}^{(l)}, \mathbf{v}^{(l+1)}, h). \quad (47c)$$

The biggest computational overhead is caused by the first problem, that is, the CCP (47a). In fact, (47c) is immediate, and (47b) can be computed quickly because in most cases the matrix  $M$  is diagonal and its inverse  $M^{-1}$  can be precomputed easily.

The convergence theory about the iterative scheme (44) leaves some degrees of freedom in choosing  $\eta_{\mathcal{A}}^i, \eta_{\mathcal{B}}^i$  values that build the diagonal blocks of the iteration matrix  $B$ . A trivial choice could be to use the same  $\eta_{\mathcal{A}}^i = \eta_{\mathcal{B}}^i = \xi$  value for all diagonal blocks, that is,  $B = \xi I$ , and then use the overrelaxation parameter  $\omega$  to control the convergence. However, this may slow convergence in systems with large mass ratios, even with an optimal  $\omega$ . A more practical approach, which copes better with systems affected by uneven masses, is inspired by the Gauss-Jacobi idea of using the inverse of the diagonal of the system matrix  $N$ , so we use  $\eta_{\mathcal{B}}^i = \frac{1}{\nabla \Psi^{i,T} M^{-1} \nabla \Psi^i}$  and  $\eta_{\mathcal{A}}^i = \frac{1}{\bar{g}_i}$ , where  $\bar{g}_i$  is the average of the diagonal values of the  $i$ th block of the  $N$  matrix. We note that  $\bar{g}_i$  can be computed easily from the trace of the  $3 \times 3$  matrix  $D^{i,T} M^{-1} D^i$ , as

$$\bar{g}_i = \frac{\text{Trace}(D^{i,T} M^{-1} D^i)}{3}. \quad (48)$$

We recall that the matrix  $N$  is a product of large matrices;  $N = D_{\mathcal{E}}^T M^{-1} D_{\mathcal{E}}$ , and it is full even if  $D$  and  $M$  are sparse. For systems with a large number of contacts, the size of  $N$  would be prohibitive and clearly would not satisfy the goal of  $O(n)$  space complexity. To this end, direct multiplication of vectors and matrices in (44) must be avoided; otherwise the effort and the space requirement would be superlinear in the number of constraint.

For the reasons above, a scheme that does not need the explicit building of  $N$ ,  $B$ , and  $K$  has been developed, exploiting the sparsity of  $M$  and  $D$ .

The  $K$  matrix in (44) can be chosen freely, within the convergence limits posed by assumptions [A1]–[A3]. Among the most noticeable options, we have the case  $K = 0$ , which results in a scheme like a projected Gauss-Jacobi, or the case where  $K$  is built by using the lower blocks of  $N$ .

### 4.1. Optimized projected Gauss-Jacobi CCP

Considering the case  $K = 0$ , and recalling Eq. (38), we can express the  $r$ th step of the iteration (44) as an inner loop with index  $i = 1 \dots n_{\mathcal{A}}$  on all  $n_{\mathcal{A}}$  friction cones  $\mathcal{FC}^i$ :

$$\delta_{\mathcal{A}}^{i,r+1} = \gamma_{\mathcal{A}}^{i,r} - \omega \eta_{\mathcal{A}}^i \left( D^{i,T} M^{-1} \left( \sum_{z=1}^{n_{\mathcal{A}}} D^z \gamma_{\mathcal{A}}^{z,r} + \sum_{z=1}^{n_{\mathcal{B}}} \nabla \Psi^z \gamma_{\mathcal{B}}^{z,r} + \tilde{\mathbf{k}} \right) + \mathbf{b}_{\mathcal{A}}^i \right) \quad (49)$$

$$\gamma_{\mathcal{A}}^{i,r+1} = \lambda \Pi_{\mathcal{FC}^i} (\delta_{\mathcal{A}}^{i,r+1}) + (1 - \lambda) \gamma_{\mathcal{A}}^{i,r}, \quad (50)$$

followed by an inner loop with index  $i = 1 \dots n_{\mathcal{B}}$  on all  $n_{\mathcal{B}}$  bilateral constraints (which we do not report because it is like (49)–(50) except for  $\mathcal{B}$  instead of  $\mathcal{A}$  subscripts).

However, for each iteration, the previous loop  $i = 1 \dots n_{\mathcal{A}}$  would require quadratic time in terms of potential contacts  $n_{\mathcal{A}}$  because of the presence of the summations  $\sum_{z=1}^{n_{\mathcal{A}}}$ . This major source of slow performance can be eliminated if one computes the algorithm in incremental form. In fact, from (36) it follows that

$$\mathbf{v}^r = M^{-1} \left( \sum_{z=1}^{n_{\mathcal{A}}} D^z \gamma_{\mathcal{A}}^{z,r} + \sum_{z=1}^{n_{\mathcal{B}}} \nabla \Psi^z \gamma_{\mathcal{B}}^{z,r} + \tilde{\mathbf{k}} \right). \quad (51)$$

Substituting (51) into (49), we can write

$$\delta_{\mathcal{A}}^{i,r+1} = \gamma_{\mathcal{A}}^{i,r} - \omega \eta_{\mathcal{A}}^i (D^{i,T} \mathbf{v}^r + \mathbf{b}_{\mathcal{A}}^i). \quad (52)$$

Considering the optimizations above, we can express the final CCP algorithm with the pseudocode of Algorithm 1.

In the proposed algorithm, for achieving high performance, some auxiliary data can be precomputed before starting the iteration. Specifically, we introduce the  $m_v \times 3$  matrix  $E_{\mathcal{A}}^i = M^{-1} D^i$  and the vector  $\mathbf{E}_{\mathcal{B}}^i = M^{-1} \nabla \Psi^i$ .

The iterations, usually stopped when an approximation threshold has been reached, can be also prematurely aborted when  $r$  exceeds a limit  $r_{max}$  on the

**Algorithm 1:** Solve complementarity - PGJ CCP

```

(1) // Pre-compute some data for friction constraints
(2) for  $i := 1$  to  $n_A$ 
(3)    $E_A^i = M^{-1}D^i$ 
(4)    $\eta_A^i = \frac{3}{\text{Trace}(D^{i,T}E_A^i)}$ 
(5) // Pre-compute some data for bilateral constraints
(6) for  $i := 1$  to  $n_B$ 
(7)    $E_B^i = M^{-1}\nabla\Psi^i$ 
(8)    $\eta_B^i = \frac{1}{\nabla\Psi^{i,T}E_B^i}$ 
(9)
(10) // Initialize impulses
(11) if warm start with initial guess  $\gamma_\varepsilon^*$ 
(12)    $\gamma_\varepsilon^0 = \gamma_\varepsilon^*$ 
(13) else
(14)    $\gamma_\varepsilon^0 = \mathbf{0}$ 
(15)
(16) // Initialize speeds
(17)  $\mathbf{v}^0 = \sum_{i=1}^{n_A} E_A^i \gamma_A^{i,0} + \sum_{i=1}^{n_B} E_B^i \gamma_B^{i,0} + M^{-1}\tilde{\mathbf{k}}$ 
(18)
(19) // Main iteration loop
(20) for  $r := 0$  to  $r_{max}$ 
(21)   // Loop on frictional constraints
(22)   for  $i := 1$  to  $n_A$ 
(23)      $\delta_A^{i,r+1} = (\gamma_A^{i,r} - \omega\eta_A^i (D^{i,T}\mathbf{v}^r + \mathbf{b}_A^i));$ 
(24)      $\gamma_A^{i,r+1} = \lambda\Pi_\Upsilon(\delta_A^{i,r+1}) + (1-\lambda)\gamma_A^{i,r};$ 
(25)   // Loop on bilateral constraints
(26)   for  $i := 1$  to  $n_B$ 
(27)      $\delta_B^{i,r+1} = (\gamma_B^{i,r} - \omega\eta_B^i (\nabla\Psi^{i,T}\mathbf{v}^r + b_B^i));$ 
(28)      $\gamma_B^{i,r+1} = \lambda\Pi_\Upsilon(\delta_B^{i,r+1}) + (1-\lambda)\gamma_B^{i,r};$ 
(29)   // Update speeds
(30)    $\mathbf{v}^{r+1} = \sum_{i=1}^{n_A} E_A^i \gamma_A^{i,r+1} + \sum_{i=1}^{n_B} E_B^i \gamma_B^{i,r+1} + M^{-1}\tilde{\mathbf{k}}$ 
(31)
(32) return  $\gamma_\varepsilon, \mathbf{v}$ 

```

maximum number of iterations if the simulation must meet hard-real-time requirements.

With minimal modifications to the  $\Pi_\Upsilon(\cdot)$  operator, the proposed method can be easily adapted to the case of friction in 2D or the case of generic unilateral constraints.

In our simulations, we chose  $\omega = 1$  and  $\lambda = 1$ , except for the  $K = 0$  case, where we used  $\omega = 0.2$ . We cannot guarantee a priori that this will satisfy condition [A3], but it did for all our simulations. In addition, the matrix sequences  $K^r$  and  $B^r$  were constant. We can therefore claim that Theorem 1 does apply and, since the sequence did not diverge, any accumulation point is a solution of the cone complementarity problem (47a). In addition, our proofs of the theoretical results allow for similar conclusions if  $\omega$  varies from iteration to iteration. Therefore, we could ensure that at some iteration the appropriate  $\omega$  is chosen after decreasing its value a few times until assumption [A3] holds. It can be shown that if the value of  $\omega$  is halved each time [A3] does not hold and the respective iteration is rejected, then [A3] will eventually be satisfied after a finite number of steps. In our experiments, however, the values we have chosen for  $\omega$  and  $\lambda$  have worked for all iterations without need of further adjustment.

#### 4.2. Optimized projected Gauss-Seidel CCP

Another option is to take  $K$  as the lower block structure of  $N$ , which results in a scheme similar to a pro-

**Algorithm 2:** Solve complementarity - PGS CCP

```

(1) // Pre-compute some data for friction constraints
(2) for  $i := 1$  to  $n_A$ 
(3)    $E_A^i = M^{-1}D^i$ 
(4)    $\eta_A^i = \frac{3}{\text{Trace}(D^{i,T}E_A^i)}$ 
(5) // Pre-compute some data for bilateral constraints
(6) for  $i := 1$  to  $n_B$ 
(7)    $E_B^i = M^{-1}\nabla\Psi^i$ 
(8)    $\eta_B^i = \frac{1}{\nabla\Psi^{i,T}E_B^i}$ 
(9)
(10) // Initialize impulses
(11) if warm start with initial guess  $\gamma_\varepsilon^*$ 
(12)    $\gamma_\varepsilon^0 = \gamma_\varepsilon^*$ 
(13) else
(14)    $\gamma_\varepsilon^0 = \mathbf{0}$ 
(15)
(16) // Initialize speeds
(17)  $\mathbf{v} = \sum_{i=1}^{n_A} E_A^i \gamma_A^{i,0} + \sum_{i=1}^{n_B} E_B^i \gamma_B^{i,0} + M^{-1}\tilde{\mathbf{k}}$ 
(18)
(19) // Main iteration loop
(20) for  $r := 0$  to  $r_{max}$ 
(21)   // Loop on frictional constraints
(22)   for  $i := 1$  to  $n_A$ 
(23)      $\delta_A^{i,r+1} = (\gamma_A^{i,r} - \omega\eta_A^i (D^{i,T}\mathbf{v}^r + \mathbf{b}_A^i));$ 
(24)      $\gamma_A^{i,r+1} = \lambda\Pi_\Upsilon(\delta_A^{i,r+1}) + (1-\lambda)\gamma_A^{i,r};$ 
(25)      $\Delta\gamma_A^{i,r+1} = \gamma_A^{i,r+1} - \gamma_A^{i,r};$ 
(26)      $\mathbf{v} := \mathbf{v} + E_A^i \Delta\gamma_A^{i,r+1}.$ 
(27)   // Loop on bilateral constraints
(28)   for  $i := 1$  to  $n_B$ 
(29)      $\delta_B^{i,r+1} = (\gamma_B^{i,r} - \omega\eta_B^i (\nabla\Psi^{i,T}\mathbf{v}^r + b_B^i));$ 
(30)      $\gamma_B^{i,r+1} = \lambda\Pi_\Upsilon(\delta_B^{i,r+1}) + (1-\lambda)\gamma_B^{i,r};$ 
(31)      $\Delta\gamma_B^{i,r+1} = \gamma_B^{i,r+1} - \gamma_B^{i,r};$ 
(32)      $\mathbf{v} := \mathbf{v} + E_B^i \Delta\gamma_B^{i,r+1}.$ 
(33)
(34) return  $\gamma_\varepsilon, \mathbf{v}$ 

```

jected Gauss-Seidel. The  $K$  matrix does not need to be explicitly built: its effect is that, as soon as computed, a reaction impulse  $\gamma^i$  will be used also for computing the following  $\gamma^{i+1}$  impulse, and so on for all  $i$ , without needing to finish a single iteration. In practical terms this means that the  $\sum_{z=1}^{n_A} D^z \gamma_A^{z,r}$  term in Eqs. (49) and (51) is split in  $\sum_{z=1}^{i-1} D^z \gamma_A^{z,r+1} + \sum_{z=i}^{n_A} D^z \gamma_A^{z,r}$ , and the same for bilateral constraints; so the difference from Algorithm 1 is that after the update of a single multiplier we immediately update  $\mathbf{v}^{(l+1)}$ , as shown in Algorithm 2. Note that, to update the speeds, we avoid the full summation (36) and add only the contributions caused by the change of the single multiplier after the projection (50):

$$\Delta\gamma_A^{i,r+1} = \gamma_A^{i,r+1} - \gamma_A^{i,r}; \quad (53)$$

$$\mathbf{v}^{(l+1),i+1} = \mathbf{v}^{(l+1),i} + M^{-1}D^i \Delta\gamma_A^{i,r+1}. \quad (54)$$

Hence, the computational overhead is not different from Algorithm 1.

Numerical tests show that this scheme converges faster than the case of  $K = 0$ ; moreover,  $K = 0$  is more prone to slow convergence in case of redundant constraints.<sup>3</sup>

<sup>3</sup> Redundancy in constraints is frequent in simulations of degenerate contact situations, such as flat surface against flat surface, where the collision engine may create a large amount of superfluous contact points.

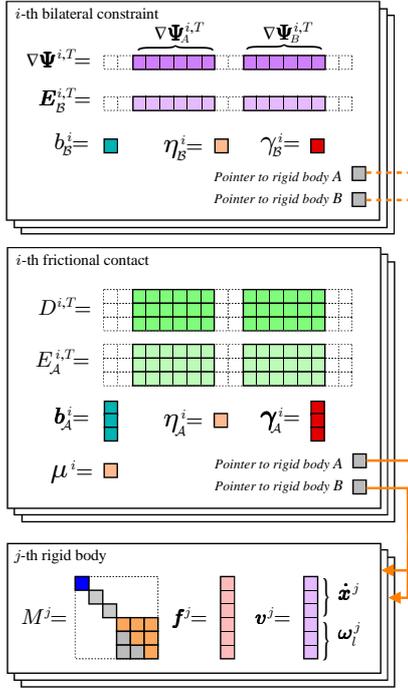


Fig. 3. Sparse data structures used by the solver.

## 5. Optimizations and improvements to the method

The proposed method can be further developed by introducing some algorithmic and theoretical improvements, obtaining different flavours of the original scheme. This section discusses the most significant optimizations.

### 5.1. Transient data structures

Figure 3 shows how the multibody model is represented by structures that are placed in memory. This transient data can be allocated on the heap during run-time, creating lists with unlimited numbers of constraints and rigid bodies.

Basically, each object that builds up the lists of bilateral constraints encapsulates the pointers to the two connected bodies, the Lagrange multiplier  $\gamma_B^i$ , the constraint residual  $b_B^i$ , the scalar value  $\eta_B^i$ , and the Jacobian  $\nabla\Psi^i$ .

The Jacobian should be a long row with  $m_v$  elements, but this would mean that the storage requirement for each constraint depends on the number of bodies, hence leading to an algorithm with superlinear space complexity. Instead, we store only the two small portions of  $\nabla\Psi^i$  that are not null, that is, the two parts  $\nabla\Psi_A^i$  and  $\nabla\Psi_B^i$  corresponding to the two connected bodies, hence requiring only a fixed number of 12 elements per Jacobian. When the  $\nabla\Psi^{i,T} \mathbf{v}^r$  multiplication must be performed, it is computed as  $\nabla\Psi_A^{i,T} \mathbf{v}_A^r + \nabla\Psi_B^{i,T} \mathbf{v}_B^r$ . Similar considerations apply to the data structure for frictional contacts, except that also  $\mu^i$  coefficients must be stored and that Jacobians are made of two larger blocks, each with three rows. Thanks to this major optimization, the memory requirement per constraint and per contact is constant,

so we obtain an optimal  $O(n_{\mathcal{E}})$  space complexity.<sup>4</sup>

Looking at Algorithms 1 and 2, one can see that the solver can operate directly on these structures. Hence, no temporary matrices are necessary. Our method is thus a matrix-free method.

### 5.2. Stabilization factor

The stabilization terms  $\frac{1}{h}\Psi^i(\mathbf{q}^{(l)})$  and  $\frac{1}{h}\Phi^i(\mathbf{q}^{(l)})$  in (16) and (17) are used to avoid constraint drifting during the time integration [4]. In fact, if these terms were missing, equations (16) and (17) would simply enforce the closure of constraints at the speed level, but errors might slowly accumulate in constraint positions after several integration steps.<sup>5</sup>

If the model were linear (that is, the matrix  $D_{\mathcal{E}}$  were constant in space, and, implicitly, in time), these terms would close constraint gaps, if any, in a single step. We experienced that, for bilateral constraints, this approach works well even if the time integration step  $h$  is small. However, this is not always the case when dealing with unilateral constraints, for the reason that follows. Because of numerical issues and nonlinearities, it is not possible to avoid some amount of penetration in contact constraints. For a contact with penetration, the term  $\frac{1}{h}\Phi^i(\mathbf{q}^{(l)})$  has a negative value. Hence, the inequality (17) requires that  $\nabla\Phi^i \mathbf{v}^{(l+1)}$  (the speed of detachment of the contact) be large enough to bring the two surfaces at zero distance in a single step  $h$ . After the time integration step, at the next step the two surfaces might keep the  $\frac{1}{h}\Phi^i(\mathbf{q}^{(l)})$  separation speed, hence causing a bouncy behavior. Apart from being not predictable, this side effect gets worse if small timesteps  $h$  are used, because even the smallest interpenetration might cause macroscopic effects that can be perceived as bouncy motions even in contacts that should have no restitution.

We investigated different strategies to improve the original stabilization method. One idea was to scale the terms by  $K^i$  factors, with  $0 < K^i < 1$ :

$$\frac{K^i}{h}\Psi^i(\mathbf{q}^{(l)}) \quad , \quad \frac{K^i}{h}\Phi^i(\mathbf{q}^{(l)}) .$$

This diminishes, but does not exclude, the risk of exceeding the interpenetration correction in a single step. Moreover, it has the drawback of creating artificial softness in stacked objects, since contact stabilization is somewhat delayed in successive frames.

A second approach, which we successfully used in many tests, involves using a clamping  $A$  as the maximum orthogonal speed for penetration recovery. This does not eliminate the risk of popping out from an intersecting contact, but at least the residual speed of separation, if any, is often negligible (comparable with

<sup>4</sup> We experienced that in many cases of granular flow simulations, the number of contacts tends to scale linearly with the number of rigid bodies  $n$ . Hence, under those circumstances the algorithm shows approximate linear space complexity  $O(n)$  also in the number of bodies.

<sup>5</sup> This would happen either because of the numerical integration, whose finite precision cannot catch all geometric nonlinearities, or because of numerical truncation and errors.

the parameter  $A$ , which can be adjusted by the user) and independent of the time step  $h$ . Thus, the stabilization terms would become

$$\frac{1}{h}\Psi^i(\mathbf{q}^{(l)}) \quad , \quad \max\left(\frac{1}{h}\Phi^i(\mathbf{q}^{(l)}), -A\right).$$

Note that in this case, if the contact surfaces are separated, we have  $\Phi^i > 0$  and the  $A$  clamping has no effect; thus it behaves as the original scheme with the  $\frac{1}{h}\Phi^i(\mathbf{q}^{(l)})$  term.

Until now we discussed how the stabilization term can correct the contact penetration errors when the surface distance is negative. However, it also has a useful side effect also when surfaces are not yet in contact. In fact, Eq. (17) can be interpreted as follows: the contact constraint is enforced only if the two surfaces are approaching fast enough to close the positive gap  $\Phi^i(\mathbf{q}^{(l)})$  in a  $h$  time step. This allows us to include in the multibody system also contact constraints that are not yet in contact but simply within a “warning envelope”; later, the cone complementarity solver will do the rest.

The third approach is more expensive in terms of CPU effort, because it avoids constraint drifting by solving an additional complementarity problem over body positions, at each time step [12]. In this case, one solves the speed CCP problem using no stabilization terms on bilateral constraints, except for contacts that are not yet in contact, where it is useful to have

$$\max\left(\frac{1}{h}\Phi^i(\mathbf{q}^{(l)}), 0\right),$$

so that contacts with clearance are still allowed to approach until contact, because of positive  $\frac{1}{h}\Phi^i$ , while surfaces already in contact are forced to have a separation speed greater than or equal to 0, regardless of the amount of penetration. Later, after the time step advancement, one performs the following poststabilization step, that is a linear complementarity problem that corrects the positions  $\mathbf{q}$ :

$$M\Delta\mathbf{q} = \sum_{i \in \mathcal{G}_A} (\beta_n^i \mathbf{D}_n^i) + \sum_{i \in \mathcal{G}_B} (\beta_b^i \nabla \Psi^i) \quad (55)$$

$$0 = \Psi^i(\mathbf{q}^{(l)}) + \nabla \Psi^{iT} \Delta\mathbf{q}, \quad i \in \mathcal{G}_B \quad (56)$$

$$0 \leq \Phi^i(\mathbf{q}^{(l)}) + \nabla \Phi^{iT} \Delta\mathbf{q} \perp \beta_n^i \geq 0, \quad i \in \mathcal{G}_A \quad (57)$$

This problem can be solved by using the same solver used for the speed CCP problem. The poststabilization idea performs better in the case of very large penetrations and ill-posed initial conditions.

In Fig. 4 a benchmark shows that the effort in reducing the maximum penetration error  $\|\epsilon_q\|_\infty$  with poststabilization iterations is almost the same as that of the basic approach using only the stabilization coefficient in the speed CCP. However some iterations on the speed CCP still must be performed. For instance, in our complex simulations involving dense granular flow, the number of iterations of the poststabilization is comparable to the number of iterations for the speed CCP, so computational efforts are almost doubled.

A family of methods can be generated by balancing the total overhead toward the speed CCP iterations

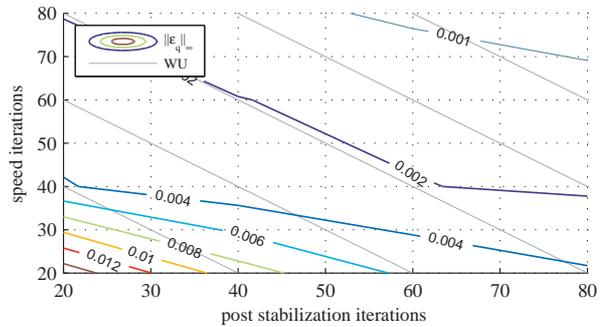


Fig. 4. Penetration error: tradeoff between speed iterations and poststabilization iterations, compared to WU overhead (GPU working units).

or the poststabilization iterations. In the extreme case where no iterations for the speed CCP are performed and iterations are performed only for the poststabilization, this method becomes similar to the *position based dynamics* algorithm, proposed in [27] as a robust method for real-time simulations.

Note that in this section we did not discuss the case of collisions with restitution and we assumed the simulation of contacts with fully plastic collisions. The reader interested in modifications to the method, for simulating also of the restitution phase, can read [1].

### 5.3. Parallelization of the algorithm

Algorithm 2 is inherently sequential, so Algorithm 1 can be easier to implement on computational architectures exploiting parallel processing because it does not feature data dependency for most of the inner loop.

In Algorithm 1, loops at rows 23 and 27 can be easily executed in parallel because they do not need to write at the same memory address at the same time. In the best scenario, one could run up to  $n_A + n_B$  independent threads, each performing an update of its multiplier  $\gamma$ .

However, the sums at row 30 are more critical if parallelized in the same way, that is, on the basis of a thread per contact, since there is the risk that multiple threads will need to update the same element of the speed vector at the same time (this situation happens, for instance, if two contacts that refer to the same body are processed simultaneously). This problem can be solved in many ways, for example, by using a vector of Boolean mutexes, one per rigid body, which can prevent one thread to modify the speed of a rigid body if some other thread is writing into it. Otherwise, if the number of physical threads is very large as in GPU architectures, it is worth performing the sums in 30 using parallel-reduction algorithms.

We implemented this parallel algorithm on both a multicore processor and on an NVIDIA Tesla C870 GPU board featuring a massively multithreaded architecture, thanks to a stream processor that is capable of processing hundreds of threads in parallel. In the latter case, we experienced a speedup of 15 times respect to the serial implementation [44].

## 6. Examples

We present two benchmarks that we used to test the performance of our algorithm and an application to the simulation of the refueling of a nuclear reactor.

### 6.1. Forklift truck simulator

As a complex mechanical system, a forklift truck represents a significant benchmark for our algorithm because it entails most aspects that we discussed in this article: rigid bodies, bilateral constraints, applied forces, and unilateral contact constraints with friction. In detail, our truck model is made of seven rigid bodies (the frame, three wheels, the steering strut, the mast, the carriage with the forks) connected by six revolute and prismatic joints. The tilting of the mast, the lifting of the carriage, and the steering are obtained by introducing rheonomic bilateral constraints. This is a simplified approach, but more detailed models of the actuators, with hydraulic components and feedback controllers, are also possible within this framework. Motors and brakes provide torques to the wheels; motors, brakes, and actuators can be controlled in real time by the user, using a joystick or a keyboard, or by using automatic procedures.

In order to also test frictional constraints, a seventh rigid body is added: a wooden pallet of EUR/ISO1 type: the truck can pick and move such a pallet with the forks. Specific collision shapes have been defined to detect the contact points between the pallets, the environment, and the truck. Collision shapes have been used also for frame of the forklift and its overhead guard, so that we can simulate the roll-over of the vehicle and other hazardous events.

On our test system, a dual-core Centrino T2600 2.17 GHz with 2 GB of RAM, the proposed algorithm is able to simulate the truck in faster than real time, using a fixed timestep  $h = 0.005$  s.

To assess the efficiency and capability of our matrixless approach, we simulated an increasing number of forklift trucks, up to 1,600 vehicles with 1,600 pallets, see Fig.5. In the largest simulation scenario, with 12,800 rigid bodies, the algorithm must handle 54,400 bilateral constraints and, on average, 19,000 frictional contacts: this means more than 110,000 primal and dual variables.

We used Algorithm 1 with a limit of 20 iterations. Table 1, averaged over 100 steps, shows that the CPU overhead grows almost linearly with the number of dual variables  $3n_A + n_B$ , that is, somehow proportional to the number of bodies.

Note that the frame rate in the case of hundreds of trucks, although not real time as for few trucks, is still fast and interactive.

Increasing the number of iterations results in an improvement of the precision: going from 20 iterations to 80 iterations, the largest errors in constraint position and in constraint velocity decrease, respectively, from 0.0310 mm and 0.024 m/s to 0.006 mm and 0.002 m/s. These results about precision are not dependent on

Table 1  
Time-step performance.

No. of Trucks	CCP Solve [ms]	Collision [ms]	Bodies	$3n_A + n_B$
1	0.22	0.1	8	70
400	135	39	3200	28000
800	268	79	6400	56000
1200	396	130	9600	84000
1600	550	200	12800	112000

the number of trucks, because the convergence of the solver is not affected by the increasing complexity of this type of benchmark, where the dynamics of each vehicle is uncoupled. On the other hand, our results indicate the efficiency of our algorithms, which is due to our customization to rigid-body dynamics structure. For example, the most efficient off-the-shelf algorithms in [32], which are either of either the interior-point or the projected gradient type, still need around 20 s per time step for about 10,000 dual variables. Those algorithms solve the same problem as here. If linear scaling would hold for those methods – a big assumption in their favor – our algorithm will still be more than 100 times faster.

How precision and convergence can be affected by systems with more stringent topology is investigated in the following example.

### 6.2. Dense granular packing benchmark

Dense stacking of granular material is one of the hardest problems involving nonsmooth rigid body dynamics. Indeed, the benchmark described in this section involves the simulation of the progressive stacking of many convex shapes in an empty box, with different settings.

The rigid bodies have a mass  $m = 10$  kg, the moments of inertia are  $I_{xx} = I_{yy} = I_{zz} = 10.24$  kgm<sup>2</sup>, and the friction coefficient is  $\mu = 0.4$ . The horizontal section of the box measures 20 m×20 m. We performed tests with different types of colliding shapes, but for the graphs presented here we used spheres with a radius  $r = 1.6$  m. At the beginning of the simulation the box is filled by 220 spheres, with randomized positions and with an initial volume fraction of 0.4, on average. After the spheres have settled, we obtained the plots.

Figures 6 and 7 show a typical convergence pattern of the PGS CCP algorithm:  $\epsilon_v$  is the violation of the constraints at the speed level. For increasing  $\omega$ , the iteration shows a faster convergence. However, large  $\omega$  values may lead to nonmonotonic behavior (that does not necessarily lead to divergence). We found that a good tradeoff between convergence speed and a monotonic nondivergent iteration, for scenarios involving equally sized spheres, is  $\omega = 1$ .

Figure 8 shows that the  $\lambda$  parameter acts like a smoothing factor over the iteration. The higher the parameter the slower is the convergence, but the lower is the risk of nonmonotonic or divergent patterns. The optimal value depends on the type of simulation; we experienced that, on average, good default values can



Fig. 5. Benchmark: simulation of thousands of forklifts.

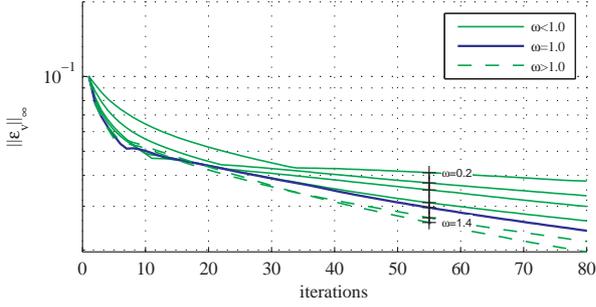


Fig. 6. Convergence of the residual for varying  $\omega$ , for a sample time step in the 220-sphere benchmark.

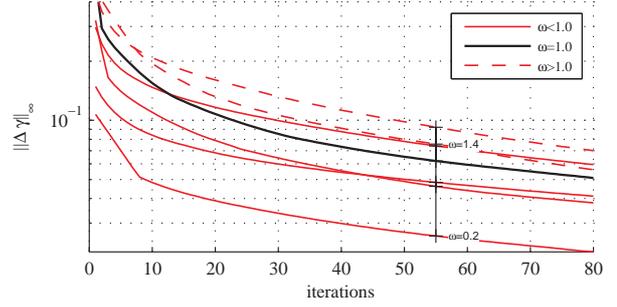


Fig. 7. Convergence of  $\Delta\gamma$  for varying  $\omega$ , for a sample time step in the 220-sphere benchmark.

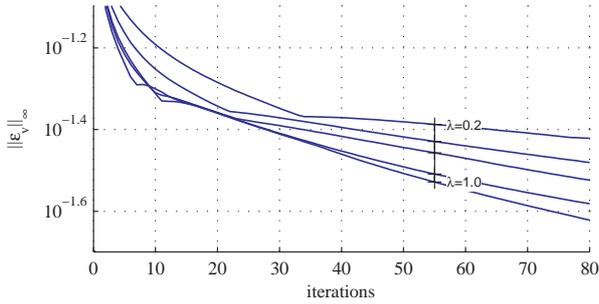


Fig. 8. Effect of the smoothing parameter  $\lambda$ , for fixed  $\omega = 1$ .

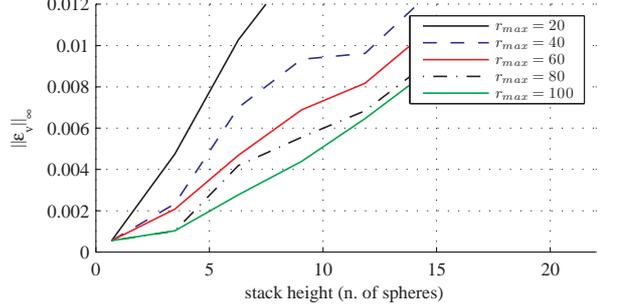


Fig. 9. Effect of the vertical size of the stack on the convergence, for varying number of iterations.

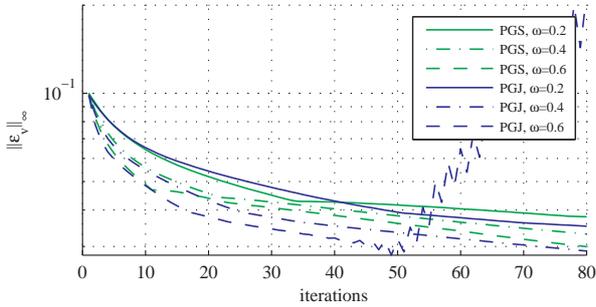


Fig. 10. Comparison of the convergence of the PGS and PGJ algorithms during the granular stacking benchmark, and example of divergence.

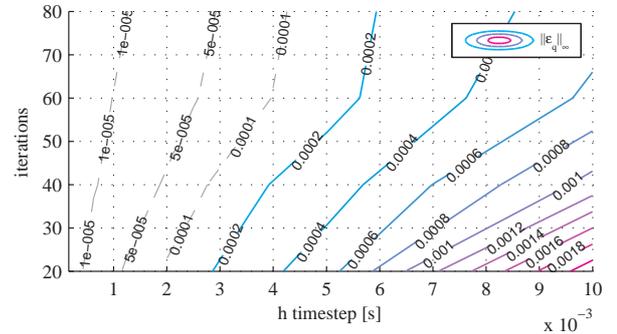


Fig. 11. Penetration error: combined effect of time step and number of iterations during the benchmark.

be chosen in the range  $\lambda = 0.8 \div 1.0$ .

The convergence can be largely affected by the topology of the mechanical system; the best-case scenario being many single objects on a flat plane, and the worst

case being the objects stacked in a vertical row. This is shown in Fig. 9, where we performed simulations with the same number of spheres but with different sizes of boxes. The precision of the iteration, for a fixed number

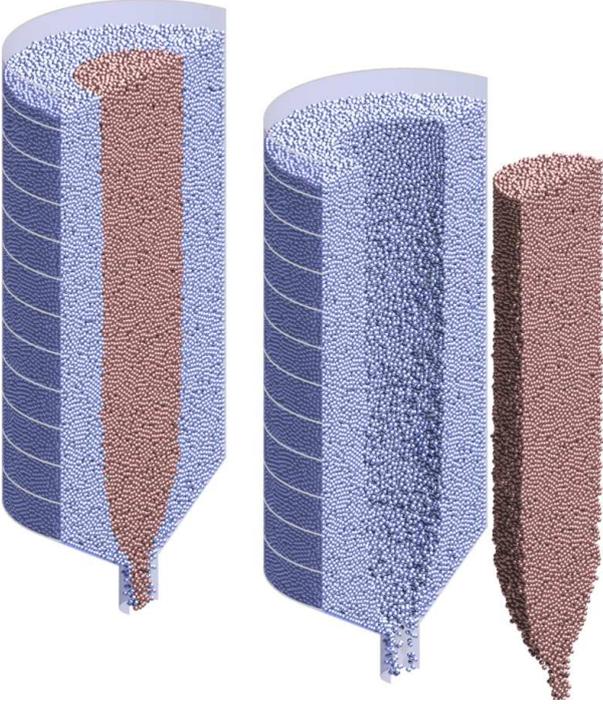


Fig. 12. Simulation of bidisperse granular flow in a PBR nuclear reactor (170,000 bodies). The inner column of graphite pebbles shows almost no dispersion in the surrounding fissile material.

of iterations, deteriorates proportionally to the height of the stack.

Figure 10 shows that, for low  $\omega$ , the convergence of the PGJ CCP algorithm is comparable to the convergence of the PGS CCP algorithm. However the PGJ CCP method is more likely to run into divergence, especially in the case of redundant constraints. Therefore, it should be used with low  $\omega$ , for example  $\omega = 0.2$  in this case.

In Fig. 11 we show the results of various simulations with different time steps and iteration numbers, while monitoring the average penetration error in constraints,  $\epsilon_p$ . The figure shows that a good precision in satisfying the constraints at the position level can be achieved both by increasing the number of iterations and by decreasing the time step  $h$ .

These plots are obtained for a sample time step, after 5 seconds from the beginning of the granular stacking simulation. Hence, previous time steps and different scenarios can lead to slightly different graphs. Although this is only a numerical benchmark, it is a worst-case scenario, and its results are indicative about the behavior of the solution method when dealing with practical engineering problems that share the same theoretical difficulties (masonry stability, soil compaction, etc.).

### 6.3. Refueling cycle in a pebble bed nuclear reactor

A significant application, which may benefit from the robustness and the speed of the method, is the simulation of the granular flow in the pebble bed nuclear reactor PBR [18].

The PBR reactor features a fourth-generation design based on a slow recirculation of fuel pebbles in a large

silo: actinides are coated and packed with graphite moderator in the spherical pebbles, each with a typical diameter of 60 mm, while the helium coolant flows between the pebbles. To increase the efficiency, a central column of spheres could contain only graphite, to flatten the neutron flux. Spheres are slowly extracted from the bottom, reprocessed, and reinserted at the top. The simulation of the downward granular flow can be useful in estimating statistical parameters such as the void fraction or the dispersion of the vertical column, hence guiding more efficient designs that can maximize the burnup of the actinides. A past attempt at simulating a PBR reactor required one week on a 64-processor supercomputer at Sandia National Laboratories, using the discrete element method (DEM) [34]. Unlike DEM methods, our approach does not introduce stiff force fields, and larger time steps are allowed, so we could perform the simulation on a single laptop computer in few hours. On average, this problem involved 170,000 rigid bodies, more than 500,000 frictional contacts, leading to more than two millions of primal and dual variables (Fig. 12).

In [42] we presented results of these simulations and validation against experimental data. Because of the high stack of spheres, this example falls in the class of problems with slow convergence already discussed in the previous benchmark: this advocates for future research efforts that could improve the performance of the solver by leveraging on multiscale and domain decomposition implementations. Nevertheless, we point out that even under this circumstance we have shown in [42] that we correctly compute macroscopic parameters of the granular flow in the pebble bed reactor, such as porosity, while needing only a few hours on a laptop.

## 7. Conclusions

We presented a formulation for multibody systems with large amounts of bilateral constraints and frictional contacts, and we developed an iterative method for this purpose. Our approach poses the problem as a convex optimization that can be solved as a cone complementarity problem. The proposed fixed-point iteration has been tailored to feature high performance even in large simulation scenarios: special care has been devoted in optimizing algorithms and formulas and avoiding matrix storage, hence obtaining an  $O(n)$  space complexity.

The method converges under standard assumptions on the configuration of the system, resulting in a robust algorithm that can simulate systems with millions of multipliers.

Improvements have been presented for the stabilization of contact constraints, resulting in matrix-free schemes that can correct large interpenetrations of rigid bodies without running into numerical problems.

We implemented a multibody system based on the method presented in this paper: the Chrono::Engine library [41]. Aiming at high performance, we expanded substantial efforts optimizing the C++ source code.

This software has been already used to simulate complex systems that were hardly tractable with other applications: granular flows in silos, interaction of wheels with sand and pebbles, size segregation devices, and other problems with a large number of bodies.

The method has been recently ported also on parallel stream-kernel GPU hardware, obtaining a remarkable computational efficiency [44].

## Acknowledgments

We thank Erwin Coumans for hints about collision detection algorithms, and Dan Negrut for comments on an earlier version of the paper. Mihai Anitescu was supported by Contract DE-AC02-06CH11357 of the U.S. Department of Energy.

## References

- [1] M. Anitescu. A fixed time-step approach for multibody dynamics with contact and friction. In *2003 IEEE/RSJ International Conference on Intelligent Robots and Systems, 2003. (IROS 2003). Proceedings*, volume 4, 2003.
- [2] M. Anitescu. Optimization-based simulation of nonsmooth rigid multibody dynamics. *Math. Program.*, 105(1):113–143, 2006.
- [3] M. Anitescu, J. F. Cremer, and F. A. Potra. Formulating 3d contact dynamics problems. *Mechanics of Structures and Machines*, 24(4):405–437, 1996.
- [4] M. Anitescu and G. D. Hart. A constraint-stabilized time-stepping approach for rigid multibody dynamics with joints, contact and friction. *International Journal for Numerical Methods in Engineering*, 60(14):2335–2371, 2004.
- [5] M. Anitescu and G. D. Hart. A fixed-point iteration approach for multibody dynamics with contact and friction. *Mathematical Programming, Series B*, 101(1)(ANL/MCS-P985-0802):3–32, 2004.
- [6] M. Anitescu and F. A. Potra. Formulating dynamic multi-rigid-body contact problems with friction as solvable linear complementarity problems. *Nonlinear Dynamics*, 14:231–247, 1997.
- [7] M. Anitescu, F. A. Potra, and D. Stewart. Time-stepping for three-dimensional rigid-body dynamics. *Computer Methods in Applied Mechanics and Engineering*, 177:183–197, 1999.
- [8] M. Anitescu and A. Tasora. An iterative approach for cone complementarity problems for nonsmooth dynamics. *Computational Optimization and Applications*, pages Accepted for printing, DOI 10.1007/s10589-008-9223-4, 2008.
- [9] D. Baraff. Issues in computing contact forces for non-penetrating rigid bodies. *Algorithmica*, 10:292–352, 1993.
- [10] D. Baraff. Fast contact force computation for nonpenetrating rigid bodies. In *Computer Graphics (Proceedings of SIGGRAPH)*, pages 23–34, 1994.
- [11] G. V. D. Bergen and G. J. A. Bergen. *Collision Detection in Interactive 3D Environments*. Morgan Kaufmann, 2004.
- [12] M. Cline and D. Pai. Post-stabilization for rigid body simulation with contact and constraints. In *IEEE International Conference on Robotics and Automation*, volume 3, pages 3744–3751, 2003.
- [13] R. Cottle and G. Dantzig. Complementary pivot theory of mathematical programming. *Linear Algebra and Its Applications*, 1:103–125, 1968.
- [14] R. W. Cottle, J.-S. Pang, and R. E. Stone. *The Linear Complementarity Problem*. Academic Press, Boston, 1992.
- [15] B. R. Donald and D. K. Pai. On the motion of compliantly connected rigid bodies in contact: a system for analyzing designs for assembly. In *Proceedings of the Conf. on Robotics and Automation*, pages 1756–1762. IEEE, 1990.
- [16] S. K. E.G. Gilbert, D.W. Johnson. A fast procedure for computing the distance between complex objects in three-dimensional space. *Robotics and Automation*, 4(2):193–203, 1988.
- [17] C. Ericson. *Real-Time Collision Detection*. Elsevier, 2005.
- [18] H. D. Gougar. *Advanced core design and fuel management for pebble-bed reactors*. Ph.D thesis, Penn State University, Department of Nuclear Engineering, 2004.
- [19] G. D. Hart. *A Constraint-stabilized Time-stepping Approach for Piecewise Smooth Multibody Dynamics*. Ph.D thesis, University of Pittsburgh, Department of Mathematics, April 2007.
- [20] E. J. Haug. *Computer Aided Kinematics and Dynamics of Mechanical Systems*. Allyn and Bacon, Boston, 1989.
- [21] E. J. Haug, S. Wu, and S. Yang. Dynamic mechanical systems with coulomb friction, stiction, impact and constraint addition-deletion. *Mechanisms and Machine Theory*, 21(5):407–416, 1986.
- [22] J.-B. Hiriart-Urruty and C. Lemarechal. *Convex Analysis and Minimization Algorithms*. Springer Verlag, Berlin, 1993.
- [23] Y. J. Kim, M. C. Lin, and D. Manocha. Deep: Dual-space expansion for estimating penetration depth between convex polytopes. In *Proceedings of the 2002 International Conference on Robotics and Automation*, volume 1, pages 921–926. Institute for Electrical and Electronics Engineering, 2002.
- [24] P. Lotstedt. Mechanical systems of rigid bodies subject to unilateral constraints. *SIAM Journal of Applied Mathematics*, 42(2):281–296, 1982.
- [25] M. D. P. Marques. *Differential Inclusions in Nonsmooth Mechanical Problems: Shocks and Dry Friction*, volume 9 of *Progress in Nonlinear Differential Equations and Their Applications*. Birkhäuser Verlag, Basel, 1993.
- [26] J. J. Moreau. Standard inelastic shocks and the dynamics of unilateral constraints. In G. D. Piero, F. Macieri, and S. Verlag, editors, *Unilateral Problems in Structural Analysis*, pages 173–221, New York, 1983. CISM Courses and Lectures no. 288.
- [27] M. Müller, B. H. M. Hennix, and J. Ratcliff. Position based dynamics. In *Proceedings of Virtual Reality Interactions and Physical Simulations*, pages 71–80, 2006.
- [28] K. Murty. *Linear Complementarity, Linear and Nonlinear Programming*. Helderman Verlag, Berlin, 1988.
- [29] J. Pang and D. Stewart. Solution dependence on initial conditions in differential variational inequalities. *Mathematical Programming*, 116(1):429–460, 2009.
- [30] J.-S. Pang, V. Kumar, and P. Song. Convergence of time-stepping method for initial and boundary-value frictional compliant contact problems. *SIAM J. Numer. Anal.*, 43(5):2200–2226, 2005.
- [31] J.-S. Pang and J. C. Trinkle. Complementarity formulations and existence of solutions of dynamic multi-rigid-body contact problems with coulomb friction. *Math. Program.*, 73(2):199–226, 1996.
- [32] C. Petra, B. Gavrea, M. Anitescu, and F. Potra. A computational study of the use of an optimization-based method for simulating large multibody systems. *Optimization Methods and Software*, 24(6):871–894, 2009.
- [33] F. Pfeiffer and C. Glocker. *Multibody Dynamics with Unilateral Contacts*. John Wiley, New York City, 1996.
- [34] C. Rycroft, G. Grest, J. Landry, and M. Bazant. Analysis of granular flow in a pebble-bed nuclear reactor. *Physical Review E*, 74, 021306, 2006.
- [35] P. Song, P. Kraus, V. Kumar, and P. Dupont. Analysis of rigid-body dynamic models for simulation of systems with frictional contacts. *Journal of Applied Mechanics*, 68(1):118–128, 2001.
- [36] P. Song, J.-S. Pang, and V. Kumar. A semi-implicit time-stepping model for frictional compliant contact problems. *International Journal of Numerical Methods in Engineering*, 60(13):267–279, 2004.

- [37] D. Stewart and J.-S. Pang. Differential variational inequalities. *Mathematical Programming*, 113(2):345–424, 2008.
- [38] D. E. Stewart. Convergence of a time-stepping scheme for rigid body dynamics and resolution of Painleve’s problems. *Archive Rational Mechanics and Analysis*, 145(3):215–260, 1998.
- [39] D. E. Stewart. Rigid-body dynamics with friction and impact. *SIAM Review*, 42(1):3–39, 2000.
- [40] D. E. Stewart and J. C. Trinkle. An implicit time-stepping scheme for rigid-body dynamics with inelastic collisions and Coulomb friction. *International Journal for Numerical Methods in Engineering*, 39:2673–2691, 1996.
- [41] A. Tasora. Chrono::Engine project, web page. [www.deltaknowledge.com/chronoengine](http://www.deltaknowledge.com/chronoengine), 2006.
- [42] A. Tasora and M. Anitescu. A convex complementarity approach for simulating large granular flow. *Journal of Computational Nonlinear Dynamics*, 2009. To appear.
- [43] A. Tasora, E. Manconi, and M. Silvestri. Un nuovo metodo del semplice per il problema di complementarit lineare mista in sistemi multibody con vincoli unilateri. In *Proceedings of AIMETA 05*, Firenze, Italy, 2005.
- [44] A. Tasora, D. Negrut, and M. Anitescu. Large-scale parallel multi-body dynamics with frictional contact on the graphical processing unit. *Proceedings of the Institution of Mechanical Engineers, Part K: Journal of Multi-body Dynamics*, 222(4):315–326, 2008.
- [45] J. Trinkle, J.-S. Pang, S. Sudarsky, and G. Lo. On dynamic multi-rigid-body contact problems with Coulomb friction. *Zeithschrift fur Angewandte Mathematik und Mechanik*, 77:267–279, 1997.

(To be removed before publication) The submitted manuscript has been created by the University of Chicago as Operator of Argonne National Laboratory (“Argonne”) under Contract No. DE-AC02-06CH11357 with the U.S. Department of Energy. The U.S. Government retains for itself, and others acting on its behalf, a paid-up, nonexclusive, irrevocable worldwide license in said article to reproduce, prepare derivative works, distribute copies to the public, and perform publicly and display publicly, by or on behalf of the Government.