



**HAL**  
open science

# The mathematical foundations of physical systems modeling languages

Albert Benveniste, Benoît Caillaud, Mathias Malandain

► **To cite this version:**

Albert Benveniste, Benoît Caillaud, Mathias Malandain. The mathematical foundations of physical systems modeling languages. Annual Reviews in Control, 2020, 50, pp.72-118. 10.1016/j.arcontrol.2020.08.001 . hal-03045498

**HAL Id: hal-03045498**

**<https://inria.hal.science/hal-03045498>**

Submitted on 8 Dec 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# The Mathematical Foundations of Physical Systems Modeling Languages

Albert Benveniste, Benoit Caillaud, Mathias Malandain

Inria Rennes, Campus de Beaulieu, 35042 Rennes cedex, France; email: surname.name@inria.fr

---

## Abstract

Modern modeling languages for general physical systems, such as Modelica, Amesim, or Simscape, rely on Differential Algebraic Equations (DAEs), i.e., constraints of the form  $f(x', x, u) = 0$ . This drastically facilitates modeling from first principles of the physics, as well as the reuse of models. In this paper, we develop the mathematical theory needed to establish the development of compilers and tools for DAE-based physical modeling languages on solid mathematical bases.

Unlike Ordinary Differential Equations (ODEs, of the form  $x' = g(x, u)$ ), DAEs exhibit subtle issues because of the notion of *differentiation index* and related *latent equations*—ODEs are DAEs of index zero, for which no latent equation needs to be considered. Prior to generating execution code and calling solvers, the compilation of such languages requires a nontrivial *structural analysis* step that reduces the differentiation index to a level acceptable by DAE solvers.

The models supported by tools of the Modelica class involve multiple modes, with mode-dependent DAE-based dynamics and state-dependent mode switching. However, multimode DAEs are much more difficult to handle than DAEs, especially because of the events of mode change. Unfortunately, the large literature devoted to the mathematical analysis of DAEs does not cover the multimode case, typically saying nothing about mode changes. This lack of foundations causes numerous difficulties to the existing modeling tools. Some models are well handled, others are not, with no clear boundary between the two classes.

In this paper, we develop a comprehensive mathematical approach supporting compilation and code generation for this class of languages. Its core is the *structural analysis of multimode DAE systems*. As a byproduct of this structural analysis, we propose sound criteria for accepting or rejecting multimode models. Our mathematical development relies on *nonstandard analysis*, which allows us to cast hybrid system dynamics to discrete-time dynamics with infinitesimal step size, thus providing a uniform framework for handling both continuous dynamics and mode change events.

*Keywords:* physical systems, modeling and simulation, differential algebraic equations, structural analysis, multimode DAE, constructive semantics

This work was supported by the FUI ModeliScale DOS0066450/00 French national grant (2018-2021) and the Inria IPL ModeliScale large scale initiative (2017-2021, <https://team.inria.fr/modeliscale/>).

---

## Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	Overall motivations . . . . .	4
1.2	Some illustration examples . . . . .	6
1.2.1	An ideal clutch . . . . .	6
1.2.2	A Cup-and-Ball game . . . . .	9
1.3	Related work . . . . .	10
1.4	Contributions . . . . .	12
<b>2</b>	<b>The ideal clutch</b>	<b>15</b>
2.1	Separate Analysis of Each Mode . . . . .	15
2.2	Mode Transitions . . . . .	17
2.3	Nonstandard Semantics . . . . .	18
2.4	Nonstandard structural analysis . . . . .	20
2.5	Standardization . . . . .	22
2.5.1	Within continuous modes . . . . .	22
2.5.2	At the instants of mode change . . . . .	23
<b>3</b>	<b>The Cup-and-Ball example</b>	<b>25</b>
3.1	Rejecting the submitted model . . . . .	26
3.2	Correcting the model . . . . .	27
3.3	Nonstandard structural analysis . . . . .	27
3.4	Standardization . . . . .	28
3.5	Handling transient modes . . . . .	31
3.6	Consequences for the modeling language . . . . .	32
<b>4</b>	<b>Issues and systematic approach</b>	<b>32</b>
<b>5</b>	<b>Background on Structural Analysis</b>	<b>34</b>
5.1	Structural analysis of algebraic equations . . . . .	34
5.1.1	Structural nonsingularity of algebraic equations . . . . .	34
5.1.2	Structural analysis . . . . .	35
5.1.3	Local versus non-local use of structural analysis . . . . .	37
5.1.4	Equations with existential quantifiers . . . . .	39
5.2	The $\Sigma$ -method for DAE systems . . . . .	41
<b>6</b>	<b>Structural analysis of multimode DAE systems</b>	<b>45</b>
6.1	Defining multimode DAE/dAE systems . . . . .	45
6.1.1	Syntax and meaning . . . . .	45
6.1.2	Long versus transient modes . . . . .	47
6.2	Structural analysis: intuition . . . . .	48
6.3	Structural Analysis of long modes . . . . .	50
6.3.1	Constructive Semantics . . . . .	51
6.3.2	Auxiliary algorithms . . . . .	54
6.3.3	Executing a nonstandard instant . . . . .	55
6.3.4	Important properties . . . . .	56

6.3.5	Generic Blocks within a long mode	58
6.4	Structural Analysis: general case	59
6.4.1	Index reduction with Difference Arrays	59
6.4.2	Executing a nonstandard instant	62
6.4.3	Continuations of transient modes	63
<b>7</b>	<b>Background on nonstandard analysis</b>	<b>66</b>
7.1	Intuitive introduction	66
7.2	Nonstandard domains	67
7.3	Nonstandard reals and integers	68
7.4	Internal functions and sets	69
7.5	Integrals	70
7.6	ODE	70
7.7	Infinitesimal calculus	71
<b>8</b>	<b>The toolkit supporting standardization</b>	<b>71</b>
8.1	Impulse Analysis	72
8.1.1	The rules of impulse analysis	73
8.1.2	Impulse analysis of systems of equations for restarts	74
8.1.3	Example: mode change $\gamma : T \rightarrow F$ of the clutch	76
8.1.4	Example: mode change $\gamma : F \rightarrow T$ of the clutch	76
8.1.5	Using impulse analysis in code generation	77
8.2	Nonstandard Structural Analysis	79
8.3	DAE of index zero	81
8.4	Long modes of mDAE	81
8.5	Mode changes of mDAE	83
<b>9</b>	<b>Main results</b>	<b>85</b>
9.1	What if we change the interpretation of the derivative?	86
9.1.1	The clutch example	86
9.1.2	A general result	87
9.2	A correctness result for the subclass of semilinear systems	95
9.3	Numerical scheme for restarts	98
9.4	Experimental results: mode changes of a non semilinear system	99
<b>10</b>	<b>Toward a tool supporting our approach</b>	<b>101</b>
10.1	The RLDC2 model	102
10.2	Handling the RLDC2 model with Modelica tools	103
10.3	Structural analysis of the RLDC2 model	104
10.4	Perspectives	106
<b>11</b>	<b>Conclusion</b>	<b>106</b>

## 1. Introduction

### 1.1. Overall motivations

*Multimode DAE systems* constitute the mathematical framework supporting the physical modeling of systems possessing different *modes*. Each mode exhibits a different dynamics, captured by *Differential Algebraic Equations* (DAEs). Multimode DAE systems are the underlying framework of ‘object-oriented’ modeling languages such as [Modelica](#), [Amesim](#), or [Simscape](#). Corresponding models can be represented as systems of guarded equations of the form

$$\text{if } \gamma_j(x_i\text{'s and derivatives)} \text{ then } f_j(x_i\text{'s and derivatives}) = 0 \quad (1)$$

where  $x_i, i = 1, \dots, n$  denote the system variables, and, for  $j = 1, \dots, m$ ,  $\gamma_j(\dots)$  is a Boolean combination of predicates guarding the differential or algebraic equation  $f_j(\dots) = 0$ . The meaning is that, if  $\gamma_j$  has the value  $\top$  (the constant true), then equation  $f_j(\dots) = 0$  has to hold; otherwise, it is discarded. In particular, when all the predicates have the value  $\top$ , one obtains a *single-mode* DAE, that is, a classical DAE defined by the set of equations

$$f_j(x_i\text{'s and derivatives}) = 0 \quad (2)$$

where  $i=1, \dots, n$  and  $j=1, \dots, m$ . When all  $f_j$ 's have the special form  $x'_j - g_j(x_1, \dots, x_n)$ , one recovers the Ordinary Differential Equations (ODE)  $x'_j = g_j(x_1, \dots, x_n)$ . DAEs are a strict generalization of ODEs, where the so-called *state variables*  $x_1, \dots, x_n$  are implicitly related to their time derivatives  $x'_1, \dots, x'_n$ . This modeling framework is fully compositional, since systems of systems of equations of the form (1) are just systems of equations of the form (1), with no restriction.

DAE systems (having a single mode) are well understood. With comparison to ODE systems, a new difficulty arises with the notion of *differentiation index*, introduced in the late 1980's [42, 21]. For simplicity, in this introduction, we discuss it for the particular case in which System (2) involves all derivatives  $x'_i$  but no higher-order derivative. The problem addressed by this notion of differentiation index relates to the Jacobian matrix  $\mathbf{J}$  associated to System (2), defined by  $\mathbf{J}_{ij} = \partial f_j / \partial x'_i$ . If this Jacobian matrix is invertible (requiring, in particular,  $m=n$ ), then all the derivatives  $x'_i$  are uniquely determined as functions of the  $x_i$ , so that System (2) is equivalent to an ODE system. If, however, matrix  $\mathbf{J}$  is singular, then System (2) as such is no longer equivalent to an ODE system. This situation can occur, even for a square DAE system possessing a unique solution for any given initial conditions. This situation makes it difficult to design DAE solvers that work for any kind of DAE.

Now,  $f_j=0$  in (2) implies  $\frac{d}{dt}f_j=0$ , revealing that such *latent equations* come implicitly with System (2). Adding latent equations to the DAE system does not change the system solutions, but it can bring additional constraints on highest-order derivatives of the original system states—on the other hand, this can introduce ‘spurious’ derivatives (of higher order than in the original system), which are to be eliminated. Adding more latent equations while eliminating spurious derivatives eventually leads to highest-order derivatives of the original system states being uniquely determined as functions of the lower order state derivatives. The *differentiation index* [42, 21] is the smallest integer  $k$

such that it is enough to differentiate with respect to time every equation  $f_j=0$  up to order at most  $k$  for the above situation to occur.

DAE solvers exist for DAE of index 1, 2, or 3 [37]. An efficient and popular approach, called the *dummy derivatives* method [37], consists in reducing the index of the DAE system to 1, producing in addition a form in which the algebraic constraints are all preserved by the solver—this no longer holds if the index is reduced to zero, leading to an ODE. *Structural analysis* methods have been proposed [40, 48] which perform index reduction by only exploiting the bipartite graph associated to the different equations and variables of the system. Structural analyses scale up much better and provide valid results outside exceptional values for the coefficients arising in the system equations [21]; structural analysis with dummy derivatives is implemented in most Modelica tools. Efforts are still ongoing to improve the efficiency and range of applicability of these methods, but one can say that DAE systems are now reasonably well understood.

In contrast, the handling of mode changes is much less mature. No structural analysis exists for mode changes. Due to discontinuities in system trajectories, the notion of differentiation index does not apply to multimode DAE systems. The notion of solutions of such systems is not even well understood, except for some subclasses of models, such as *semilinear systems* [8]. Still, modeling languages exist that support multimode DAE systems, e.g., *Modelica*, *VHDL-AMS*, and the proprietary languages *Amesim* and *Simscape*. However, the lack of mathematical understanding of mode changes can result in a spurious handling of some physically meaningful models. Indeed, with the exception of the recent work [8], the class of “safe” models (well supported by the considered tool) is never characterized, thus leading to a “try-and-see” methodology. This situation motivates our work.

Establishing the mathematical foundations for compilers<sup>1</sup> of multimode DAE systems raises non-classical difficulties. To substantiate this claim, let us compare the following three subject matters:

1. Developing formal verification for a given class of hybrid systems [2, 45];
2. Developing existence/uniqueness results and discretization schemes, for a given class of (single-mode or multimode) DAE systems [3];
3. Developing the mathematical foundations for multiphysics/multimode DAE systems modeling languages (our focus).

For subject 1, restrictions on the class are stated, under which the proposed algorithms are proved correct, and their complexity is analyzed. Similarly, for subject 2, assumptions are stated on which discretization schemes are proved correct, and convergence rates can be given. In both cases, assumptions are formulated and it is the responsibility of the user to check their validity when performing verification or applying discretization schemes.

In contrast, for subject 3, one cannot expect a compiler to check the existence/uniqueness of solutions of a given multimode DAE system. Worse, we cannot expect the user to check this as part of his/her model design activity. In fact, the compiler must handle *any* submitted model, and it has the responsibility for accepting or rejecting a model on the sole basis of syntactic or symbolic (but never numerical) analyses. This is a demanding task that we claim is not well addressed today.

---

<sup>1</sup>By *compilation*, we mean here the suite of symbolic analyses and transformations that must be performed prior to generating simulation code.

## 1.2. Some illustration examples

To illustrate the above discussion, we now review two specific examples of multimode DAE systems.

### 1.2.1. An ideal clutch

This clutch is depicted in Figure 1. It is a simple, idealized clutch involving two rotating shafts where no motor or brake are connected. More precisely, we assume that this system is closed, with no interaction other than explicitly specified.

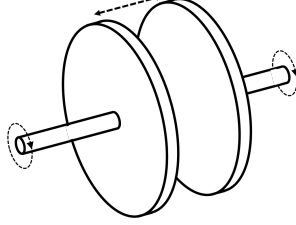


Figure 1: An ideal clutch with two shafts.

We provide hereafter in (3) a model for it, complying with the general form (1). The dynamics of each shaft  $i$  is modeled by  $\omega'_i = f_i(\omega_i, \tau_i)$  for some, yet unspecified, function  $f_i$ , where  $\omega_i$  is the angular velocity,  $\tau_i$  is the torque applied to shaft  $i$ , and  $\omega'_i$  denotes the time derivative of  $\omega_i$ . Depending on the value of the input Boolean variable  $\gamma$ , the clutch is either engaged ( $\gamma = \text{T}$ , the constant “true”) or released ( $\gamma = \text{F}$ , the constant “false”). When the clutch is released, the two shafts rotate freely: no torque is applied to them ( $\tau_i = 0$ ). When the clutch is engaged, it ensures a perfect join between the two shafts, forcing them to have the same angular velocity ( $\omega_1 - \omega_2 = 0$ ) and opposite torques ( $\tau_1 + \tau_2 = 0$ ). Here is the model:

$$\left\{ \begin{array}{ll} & \omega'_1 = f_1(\omega_1, \tau_1) \quad (e_1) \\ & \omega'_2 = f_2(\omega_2, \tau_2) \quad (e_2) \\ \text{if } \gamma \text{ then} & \omega_1 - \omega_2 = 0 \quad (e_3) \\ & \text{and } \tau_1 + \tau_2 = 0 \quad (e_4) \\ \text{if not } \gamma \text{ then} & \tau_1 = 0 \quad (e_5) \\ & \text{and } \tau_2 = 0 \quad (e_6) \end{array} \right. \quad (3)$$

When  $\gamma = \text{T}$ , equations  $(e_3, e_4)$  are active and equations  $(e_5, e_6)$  are disabled, and vice-versa when  $\gamma = \text{F}$ . If the clutch is initially released, then, at the instant of contact, the relative speed of the two rotating shafts jumps to zero; as a consequence, an impulse is expected on the torques. In addition, the model yields an ODE system when the clutch is released, and a DAE system of index 1 when the clutch is engaged, due to equation  $(e_3)$  being active in this mode.

Note that the condition  $\tau_1 + \tau_2 = 0$  is an invariant of the system in both modes, which reflects the preservation of the angular momentum due to the assumption that the system is closed. Model (3) is one possible specification. One could have instead put  $\tau_1 + \tau_2 = 0$  outside the scope of any guard, thus making explicit that this is an invariant—this is actually the form adopted for the Modelica model of Figure 2; then, mode  $\gamma = \text{F}$  only

states  $\tau_1 = 0$ . Both models are equivalent. It turns out that our approach yields the same generated code for both source codes.

```

model ClutchBasic
  parameter Real w01=1;
  parameter Real w02=1.5;
  parameter Real j1=1;
  parameter Real j2=2;
  parameter Real k1=0.01;
  parameter Real k2=0.0125;
  parameter Real t1=5;
  parameter Real t2=7;
  Real t(start=0, fixed=true);
  Boolean g(start=false);
  Real w1(start = w01, fixed=true);
  Real w2(start= w02, fixed=true);
  Real f1; Real f2;

  equation
    der(t) = 1;
    g = (t >= t1) and (t <= t2);
    j1*der(w1) = -k1*w1 + f1;
    j2*der(w2) = -k2*w2 + f2;
    0 = if g then w1-w2 else f1;
    f1 + f2 = 0;
end ClutchBasic;

```

Figure 2: Modelica code for the idealized clutch.

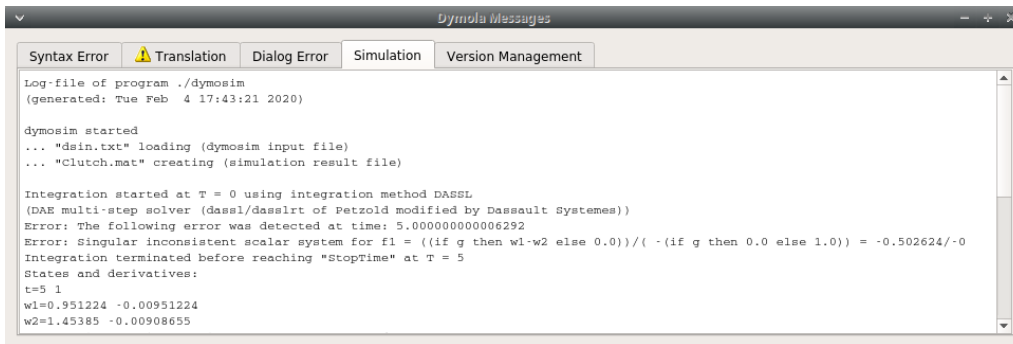


Figure 3: Division by zero exception occurring when simulating the Ideal Clutch Modelica model.

*The clutch in Modelica.* Figure 2 details the Modelica model of the Ideal Clutch system, with two shafts interconnected by an idealized clutch. It is a faithful translation in the Modelica language of the mDAE (3), with the exception that the two differential equations have been linearized. Also, the trajectory of the input guard  $\gamma$  (here called  $g$ ) has been fully defined: it takes the value T between instants  $t_1$  and  $t_2$ , and F elsewhere. The model is deemed to be structurally nonsingular by the two Modelica tools we had the opportunity to test: OpenModelica and Dymola. However, none of these tools generates correct simulation code from this model. Indeed, simulations fail precisely at the instant when the clutch switches from the uncoupled mode ( $g=false$ ) to the coupled one ( $g=true$ ). This is evidenced by a division by zero exception, as shown in Figures 3 and 4.

The root cause of this exception is that none of these tools performs a multimode structural analysis. Instead, the structure of the model is assumed to be invariant, and the structural analysis implemented in these tools is a Dummy Derivatives method, which is proved to be correct only on single-mode DAE systems; it is possibly correct on multimode



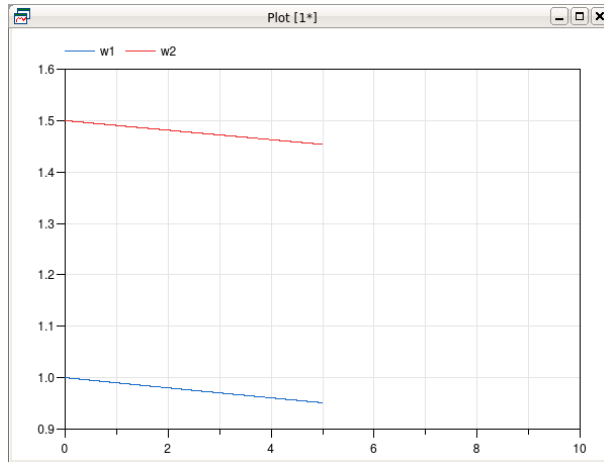


Figure 4: Trajectory of the Ideal Clutch Modelica model: it stops at  $t = 5s$ , when the exception occurs.

systems, but under the stringent assumption that the model structure (including, but not limited to, its differentiation index) is independent of the mode. The consequence is that the structural analysis methods implemented in these tools do not detect that the differentiation index jumps from 0 to 1 when the shafts are coupled, and that the structure is not invariant. The division by zero results from the pivoting of a linear system of equations that becomes singular when  $g$  becomes equal to `true`.

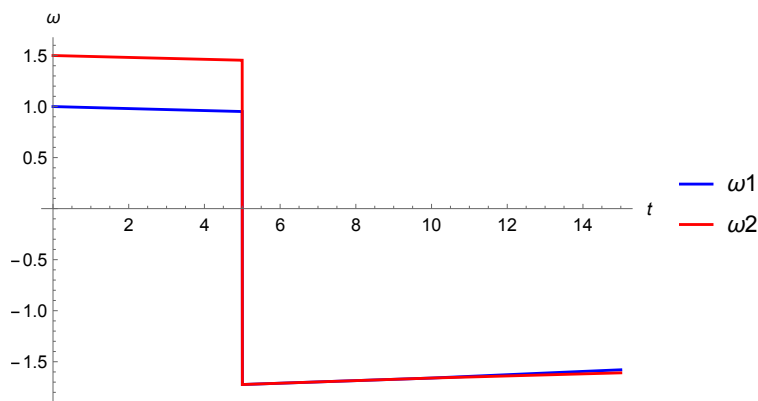


Figure 5: Solution of the Ideal Clutch model, computed by Mathematica.

*The clutch in Mathematica.* DAE systems can be solved with the Mathematica “Advanced Hybrid & DAE” toolbox.<sup>2</sup> This library is also advertised as supporting multimode DAE systems. Figure 5 shows the solution of the Ideal Clutch system, computed by

<sup>2</sup><https://www.wolfram.com/mathematica/new-in-9/advanced-hybrid-and-differential-algebraic-equations/>

Mathematica. It can be seen that the angular momentum of the mechanical system is not preserved, in contradiction with the physics. Also, a small change in the velocities before the change results in totally different values for restart, although the system is not chaotic by itself. A likely explanation is that Mathematica regards the mode switching as a consistent initialization problem, with one degree of freedom. Indeed, any choice of angular velocities satisfying the constraint  $\omega_1 = \omega_2$  is a consistent initial state for the DAE in the engaged mode.

The conclusion of these two experiments is that, clearly, some fundamental study is needed to ensure a correct handling of events of mode change by the Modelica tools. Smoothing the ‘if then else’ equation could help solving the above problem, but requires a delicate and definitely non-modular tuning, as it depends on the different time scales arising in the system. We believe that, as the tools reputedly support multimode DAE models, they should handle them correctly.

*Our objective.* For this example, *our objective is to be able to compile System(3) and produce correct simulations for it, without any need for smoothing the mode change.* This objective is addressed in Section 2.

### 1.2.2. A Cup-and-Ball game



Figure 6: The Cup-and-Ball game.

We sketch here a multimode extension of the popular example of the pendulum in Cartesian coordinates [40], namely the planar Cup-and-Ball game, whose 3D version is illustrated by Figure 6. A ball, modeled by a point mass, is attached to one end of a rope, while the other end of the rope is fixed, to the origin of the plane in the model. The ball is subject to the unilateral constraint set by the rope, but moves freely while the distance between the ball and the origin is less than its actual length. Again, we assume that the system is closed and subject to no interaction other than specified. Using Cartesian coordinates  $x$  and  $y$ , the equations of the planar Cup-and-Ball game are the following:

$$\begin{cases} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ 0 \leq L^2 - (x^2 + y^2) & (\kappa_1) \\ 0 \leq \lambda & (\kappa_2) \\ 0 = [L^2 - (x^2 + y^2)] \times \lambda & (\kappa_3) \end{cases} \quad (4)$$

where the unknowns (also called *dependent variables*) are the position  $(x, y)$  of the ball in Cartesian coordinates and the rope tension  $\lambda$ .

The subsystem  $(\kappa_1, \kappa_2, \kappa_3)$  expresses that the distance of the ball from the origin is less than or equal to  $L$ , the tension is nonnegative, and one cannot have a distance less

than  $L$  and a nonzero tension at the same time. This is known as a *complementarity condition*, written as  $0 \leq L^2 - (x^2 + y^2) \perp \lambda \geq 0$  in the *nonsmooth systems* literature [1], and is an adequate modeling of ideal valves, diodes, and contact in mechanics. Constraints  $\kappa_1$  and  $\kappa_2$  are unilateral, which does not belong to our framework of guarded equations. Therefore, using a technique communicated to us by H. Elmqvist and M. Otter, we redefine the graph of this complementarity condition as a parametric curve, represented by the following three equations:

$$\begin{aligned} s &= \text{if } \gamma \text{ then } -\lambda \text{ else } L^2 - (x^2 + y^2) \\ 0 &= \text{if } \gamma \text{ then } L^2 - (x^2 + y^2) \text{ else } \lambda \\ \gamma &= [s \leq 0] \end{aligned}$$

which allows us to rewrite the complementarity condition  $(\kappa_1, \kappa_2, \kappa_3)$  as follows:

$$\left\{ \begin{array}{ll} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ \gamma = [s \leq 0] & (k_0) \\ \text{if } \gamma \text{ then } 0 = L^2 - (x^2 + y^2) & (k_1) \\ \quad \text{and } 0 = \lambda + s & (k_2) \\ \text{if not } \gamma \text{ then } 0 = \lambda & (k_3) \\ \quad \text{and } 0 = (L^2 - (x^2 + y^2)) - s & (k_4) \end{array} \right. \quad (5)$$

Unsurprisingly (considering the reasons for the wrong handling of the clutch example), the Modelica tools also fail to handle this model correctly.

Some new issues emerge from Example (5). First, subsystem  $(\kappa_1, \kappa_2, \kappa_3)$  of (4) leaves the impact law at mode change insufficiently specified: it could be fully elastic, fully inelastic, or inbetween. Second, (5) exhibits a logico-numerical fixpoint equation in  $s$ , which we regard as problematic. While the Jacobian matrix yields a regularity criterion for systems of smooth algebraic equations, logico-numerical systems of equations are unfriendly: there is no simple criterion for the existence and/or uniqueness of solutions. Note that a sensible modification of model (5) would consist in replacing, in  $(k_0)$ ,  $s$  by its left-limit  $s^-(t) =_{\text{def}} \limsup_{u \nearrow t} s(u)$ .<sup>3</sup>

*Our objectives.* For this example, our objectives are *to address the above difficulties at compile time, and to generate correct simulation code once structural problems are fixed*. These objectives are addressed in Section 3.

### 1.3. Related work

There exists a huge litterature devoted to the multi-physics DAEs (having a single mode); see, e.g., [16] as a basic reference. Many references related to the Modelica language can also be found on the website of the Modelica organization.<sup>4</sup> We focus this review on works addressing multimode DAE systems with an emphasis on the compilation and code generation, for related models.

<sup>3</sup>The left-limit is available in Simulink under the ‘state port’ construct, whereas the `pre` keyword in Modelica serves the same purpose.

<sup>4</sup><https://www.modelica.org/publications>

A large body of work was developed in the domain of multi-body systems with contacts where impulses can occur; see for instance [44, 43] for an overview of this domain. The current research focuses on the handling of contacts between many bodies using time-stepping methods. Contact equations are usually expressed on the velocity or acceleration, and since the constraints on position level are not explicitly taken into account, a drift typically occurs. An exception is [50], where constraints on position level are enforced. It is not obvious how the specialized multimode methods for multi-body systems can be generalized to any type of multimode DAE. For electrical circuits with idealized switches, like ideal diodes, impulses can also occur. Again, a large literature is available, concentrating mostly on specialized electrical circuits, such as piecewise-linear networks with idealized switches [30]. However, it is not obvious how to generalize methods in this restricted area to general multimode DAEs.

The article by Mehrmann *et al.* [38] contains interesting results regarding numerical techniques to detect chattering between modes. It however assumes that consistent reset values are explicitly given for each mode. Such an assumption does not hold in general, particularly for models derived from first principles of the physics; the clutch and cup-and-ball examples are good illustrations of this.

In the PhD thesis of Zimmer [54], variable-structure multi-domain DAEs are defined thanks to an *ad hoc* modeling language, and a runtime interpreter is used that processes the equations at runtime, when the structure and/or the index changes. Limitations of this work are that impulsive behavior is not supported and that the user has to explicitly define the transfer of variable values from one mode to another, which is not practical for large models.

Describing variable-structure systems with *causal* state machines is discussed by Pepper *et al.* [41]. Dynamically changing the structural analysis at runtime is also performed by Hoeger [31, 32]. In [31], the author proposes a dynamic execution of John Pryce’s  $\Sigma$ -method [48].

Elmqvist *et al.* [28, 36] propose a high-level description of multimode models as an extension to the synchronous Modelica 3.3 state machines, by using continuous-time state machines having continuous-time models as ‘states’. Besides ODEs as used in hybrid automata, *acausal DAE models with physical connectors* can also be a ‘state’ of a state machine. Such a state machine is mapped into a form such that the resulting equations can be processed by standard symbolic algorithms supported by Modelica tools. The major restrictions of this approach are that mode changes with impulsive behavior are not supported and that not all types of multimode systems can be handled due to the static code generation.

Correct restart at mode changes is clearly a central issue in multi-mode DAE systems. Techniques from *nonsmooth dynamical systems* [1] address this issue by performing a warm restart at mode changes using time-stepping discretization schemes: not only the numerical scheme does not interrupt at mode changes, but these mode changes are not even detected. This is very elegant and effective. However, this relies on the handling of *complementarity conditions* (see Section 1.2.2), thus, it does not apply to general multimode DAE systems but only to a certain subclass.

We observe that, in general, the issue of correct restart at mode changes is not addressed in the above mentioned literature. It is, however, considered in the few references to follow.

Benveniste *et al.* [7, 6] considered this issue, as well as the problem of varying structure

and index, from a fundamental point of view, by relying on nonstandard analysis to capture continuous-time dynamics and mode change events in a unified framework. A first structural analysis algorithm was presented in [7], by significantly modifying the original Pantelides algorithm [40]. A proof-of-concept mockup named SunDAE was developed implementing this approach. This first attempt suffers from some issues: the proposed structural analysis does not boil down trivially to the Pantelides algorithm in the case of single-mode systems; it involves nondeterministic decisions, an unwanted feature for the mathematical foundation of compilers; and its mathematical study is incomplete.

The work of Stephan Trenn is an important contribution to the subject. In his PhD thesis [52] and his article [53], he points out the difficulty in defining piecewise smooth distributions: there is no intrinsic definition for their ‘Dirac’ part, as several definitions comply with the requirements for it. This indicates that distributions are not the ultimate answer to deal with impulsive variables in multimode DAE systems. Still, Trenn was able in [34] to define complete solutions for a class of switched DAE systems in which each mode is in *quasi-linear form*: switching conditions are time-based, not state-based. A main difference with our approach is that Trenn’s solution is complete for all variables and relies on piecewise smooth distributions, whereas our works on restart conditions consist in eliminating impulsive variables while computing restart values for state variables. (Beyond the existence of solutions, the work of Trenn *et al.* contains additional interesting results that are not relevant to our study.)

An important step forward was done in [8]. The interesting subclass of multimode DAE systems that is now called *semilinear* (see Section 9.2) was first identified. Semilinear multimode DAE systems can exhibit impulsive variables at mode changes. They generalize the ‘semilinear systems’ proposed by Trenn in the sense that switching conditions are no longer restricted to time-based ones, instead including state-based switching conditions. The analysis and discretization schemes proposed in [8] are mathematically sound. Building on this work, Martin Otter has developed the [ModiaMath](#) tool for semilinear multimode DAE systems. Hence, it makes sense to compare the schemes proposed in [8] to the ones we develop in this paper, for general mDAE systems. It turns out that our general approach coincides with the schemes proposed in [8] when applied to the subclass of semilinear systems, see Section 9.2. Our present work thus extends and significantly improves [8].

#### 1.4. Contributions

In this work, we develop a systematic approach that addresses, as particular cases, the objectives stated for the examples of Sections 1.2.1 and 1.2.2. Its main contributions can be detailed as follows.

*Structural analysis of mode changes using nonstandard analysis.* As our main contribution, we extend the notion of structural analysis to mode change events. More precisely, we develop a structural analysis for multimode DAE systems that is valid at any time, that is, for both continuous dynamics and mode changes.

Doing so raises a number of difficulties. First, mode changes occur under various kinds of dynamics; there can be isolated events or finite cascades thereof, but Zeno situations can also occur, i.e., events of mode changes can accumulate and form an infinite sequence bounded by some finite instant. As a result, sliding modes can emerge out of mode changes, and the modeling and simulation tools can only identify such situations

at runtime. Second, impulses may occur at mode changes for certain variables, which again is most of the time discovered by tools at runtime. Compilation, however, is a task performed prior to generating simulation code; hence, it encounters serious difficulties in both the nature of time at mode changes and the domain of variables, due to impulsive behaviors.

The main objective of structural analysis, akin to compilation in general, is precisely to keep away from numerical considerations, by focusing only on the structure (or the syntax) of the system of equations. This is captured by the bipartite graph associated to the considered DAE system (also referred to as ‘incidence graph’ or ‘incidence matrix’ or ‘ $\Sigma$ -matrix’ depending on the context and authors). The basic principle supporting structural analysis consists in: (1) abstracting an equation  $f(x, y, z) = 0$  as the bipartite graph having  $f, x, y, z$  as vertices and  $(f, x), (f, y), (f, z)$  as edges, and a system of equations as the union of the above bipartite graphs; and (2) exploiting the resulting graph to generate efficient simulation code. This means focusing on syntax and symbols, while abstracting away numerical aspects. We thus wish to *symbolically* manipulate both normal and impulsive values for variables, as well as both continuous time and events (in finite or infinite cascades), and possibly stranger structures for time.

We see no vehicle for supporting all of this, other than *nonstandard analysis*. Nonstandard analysis was proposed by Abraham Robinson in the 1960s to allow the explicit manipulation of ‘infinitesimals’ and ‘infinities’ in analysis [49, 23]. Impulsive values become normal citizens in nonstandard analysis. Continuous time and dynamics can be discretized using infinitesimal time steps, e.g., by setting  $x'(t) \approx \frac{x(t+\partial) - x(t)}{\partial}$  with  $\partial$  infinitesimal, thus providing perfect fidelity up to infinitesimal errors. This discretized time can then support the various kinds of time needed when modeling mode changes. Since techniques from structural analysis translate almost verbatim to discrete-time dynamics, having dynamical systems indexed by discrete (nonstandard) time in a uniform setting paves the way toward structural analysis for multimode DAE systems in their generality. Using this approach, we are able to extend the structural analysis, from single-mode to multimode DAE systems, by encompassing both modes and mode changes.

*Generating code for the restart at mode changes.* The structural analysis developed hereafter works in the nonstandard analysis domain. As such, it cannot be executed on any real-life computer. A further stage is needed to produce executable code, by mapping the nonstandard execution schemes to ordinary world (standard) schemes; this is called *standardization*. This allows us to produce mathematically justified schemes for the restart at mode changes.

*Rejecting and accepting programs on a clear basis.* Our structural analysis is precise enough to properly identify models that are structurally over- or under-specified at mode change events, or models that contain a fixpoint equation involving a subset of variables and some guard. We reject such models, with a proper diagnosis explaining the rejection. In turn, mode-dependent index/state/dynamics are not reasons for rejection *per se*. Our compilation technique is able to handle such cases.

*The paper is organized as follows.* The clutch example is informally developed in Section 2; we show how to map this model to the nonstandard domain, how the structural analysis

works (in particular, the handling of mode changes is emphasized), and how standardization is performed to generate actual simulation code. The Cup-and-Ball example is developed in Section 3, following the same guidelines; based on the case of elastic impact, it brings forth the need for considering transient modes, i.e., modes that last for zero time. The insight gained from studying these examples leads to the approach described in Section 4.

The extension of structural analysis to multimode DAE systems is then developed. In Section 5, the basics of structural analysis for algebraic systems of equations are recalled. Then, we review J. Pryce’s  $\Sigma$ -method for the structural analysis of DAE systems. We extend the  $\Sigma$ -method to the multimode case, including the handling of mode change events, in Section 6, for the subclass of systems involving only long modes—which excludes the Cup-and-Ball example with elastic impact. Motivated by the latter case, we identify the need for handling *transient* modes, in which the system spends zero time; this is addressed in Section 6.4.

Section 7 recalls the needed background on nonstandard analysis, mostly related to differential calculus. In Section 8, we complement this background with additional results related to structural analysis and impulse analysis, and we apply this machinery to the task of standardization.

Three important results are developed in Section 9.

- In Section 9.1, we prove that modifying the nonstandard expansion of the derivative (there are infinitely many possibilities) does not change the final generated code, thus showing that our approach is intrinsic.
- We prove in Section 9.2 that our approach does compute correct solutions for *semilinear* multimode DAE systems, a nontrivial subclass possibly involving impulsive behaviors and containing, as a subclass, multi-body mechanics (see [8]). We would be happy to prove that the simulation code we generate does compute the solution of any multimode DAE system; this, however, is beyond what is doable since no notion of solution is known for multimode DAE systems in general.
- We propose in Section 9.3 a numerical scheme for approximating restart values for non-impulsive variables at mode changes. This alleviates the need for performing standardization when generating code in practice; the proof of this result, however, uses standardization as a mathematical tool. This scheme is illustrated on a (non semilinear) example in Section 9.4.

Finally, we draw in Section 10 how our approach can be implemented in a tool, and we illustrate this on the RLDC2 circuit example. This example, provided to us by Sven-Erik Mattsson, is problematic for the existing Modelica tools, as it has mode-dependent index and structure. The structural analysis of the corresponding DAE system in all its modes is successfully handled by our IsamDAE tool [19]. This tool, currently under development, will support multimode DAE structural analysis following our theory. Its current version does not support impulsive variables at mode changes—the RLDC2 circuit example, however, does not exhibit this difficulty.

## 2. The ideal clutch

In this section, we introduce our approach by discussing in detail the ideal clutch example presented in Section 1.2.1. Its model was given in (3). We first analyze separately the model for each mode of the clutch (Section 2.1). Then, we discuss the difficulties arising when handling mode changes (Section 2.2). Finally, we propose a global comprehensive analysis in Sections 2.3–2.5. For convenience, we present an intuitive and informal introduction to nonstandard analysis in Section 2.3.

### 2.1. Separate Analysis of Each Mode

In the released mode, when  $\gamma$  is false in System (3), the two shafts are independent and one obtains the following two independent ODEs for  $\omega_1$  and  $\omega_2$ :

$$\begin{aligned} \omega_1' &= f_1(\omega_1, \tau_1) & (e_1) & & \tau_1 &= 0 & (e_5) \\ \omega_2' &= f_2(\omega_2, \tau_2) & (e_2) & & \tau_2 &= 0 & (e_6) \end{aligned} \quad (6)$$

In the engaged mode, however,  $\gamma$  holds true, and both the two velocities and the two torques are algebraically related:

$$\begin{aligned} \omega_1' &= f_1(\omega_1, \tau_1) & (e_1) & & \omega_1 - \omega_2 &= 0 & (e_3) \\ \omega_2' &= f_2(\omega_2, \tau_2) & (e_2) & & \tau_1 + \tau_2 &= 0 & (e_4) \end{aligned} \quad (7)$$

Equation ( $e_4$ ), that relates the two torques, is not an issue by itself, as it can be rewritten as  $\tau_2 = -\tau_1$  and used to eliminate  $\tau_2$ . Equation ( $e_3$ ), however, relates the two velocities  $\omega_1$  and  $\omega_2$  that are otherwise subject to the ODE system ( $e_1, e_2$ ), which makes System (7) a DAE instead of an ODE.

In System (7), the torques  $\tau_1$  and  $\tau_2$  are not differentiated: we call them *algebraic variables*. The derivatives of the velocities  $\omega_1$  and  $\omega_2$  appear: they are *state variables*. The  $\omega_1', \omega_2', \tau_1, \tau_2$  are called the *leading variables* of System (7).

Now, suppose one is able to uniquely determine the leading variables  $\omega_1', \omega_2', \tau_1, \tau_2$  given *consistent* values for the state variables  $\omega_1, \omega_2$ , i.e., values satisfying ( $e_3$ ). Then, by using an ODE solver, one could perform an integration step and update the current velocities  $\omega_1, \omega_2$  using the computed values for their derivatives  $\omega_1', \omega_2'$ . In this case, we say that the considered DAE is an “extended ODE” [48].

It turns out that this condition does not hold for System (7) as is. To intuitively explain the issue, we move to discrete time by applying an explicit first-order Euler scheme with constant step size  $\delta > 0$ :

$$\begin{aligned} \omega_1^\bullet &= \omega_1 + \delta \cdot f_1(\omega_1, \tau_1) & (e_1^\delta) & & \omega_1 - \omega_2 &= 0 & (e_3) \\ \omega_2^\bullet &= \omega_2 + \delta \cdot f_2(\omega_2, \tau_2) & (e_2^\delta) & & \tau_1 + \tau_2 &= 0 & (e_4) \end{aligned} \quad (8)$$

where  $\omega^\bullet(t) =_{\text{def}} \omega(t+\delta)$  denotes the forward time shift operator by an amount of  $\delta$ . Suppose we are given consistent values  $\omega_1 = \omega_2$  and we wish to use System (8), seen as a system of algebraic equations, to determine the value of the dependent variables (i.e., the unknowns)  $\tau_1, \tau_2, \omega_1^\bullet, \omega_2^\bullet$ . This attempt fails since we have only three equations  $e_1^\delta, e_2^\delta$  and  $e_4$  to determine four unknowns  $\tau_1, \tau_2, \omega_1^\bullet$  and  $\omega_2^\bullet$ ; indeed, the omitted equation ( $e_3$ ) does not involve any dependent variable.



However, since System (8) is time invariant, and assuming that the system remains in the engaged mode for at least  $\delta$  seconds, there exists an additional *latent equation* on the set of variables, namely

$$\omega_1^\bullet - \omega_2^\bullet = 0 \quad (e_3^\bullet) \quad (9)$$

obtained by shifting  $(e_3)$  forward. Now, replacing  $(e_3)$  with  $(e_3^\bullet)$  in System (8) yields a system with four equations and four dependent variables; moreover, this system is, in fact, nonsingular in a generic sense. One can now use System (8) along with equation  $(e_3^\bullet)$  to get an execution scheme for the engaged mode of the clutch. This is shown in Exec. Sch. 1 below.

---

**Execution Scheme 1** System (8)+Eq. (9).

---

**Require:** consistent  $\omega_1$  and  $\omega_2$ , i.e., satisfying  $(e_3)$ .

- |  |  |
|--|--|
| 1: Solve $\{e_1^\delta, e_2^\delta, e_3^\bullet, e_4\}$ for $(\omega_1^\bullet, \omega_2^\bullet, \tau_1, \tau_2)$ | ▷ 4 equations, 4 unknowns                  |
| 2: $(\omega_1, \omega_2) \leftarrow (\omega_1^\bullet, \omega_2^\bullet)$  | ▷ update the states $(\omega_1, \omega_2)$ |
| 3: Tick  | ▷ move to next discrete step               |
- 

Since the new values of the state variables satisfy (9) by construction, the consistency condition is met at the next iteration step (should the system remain in the same mode).

**Comment 1 (structural nonsingularity).** The implicit assumption behind Line 1 in Exec. Sch. 1 is that solving  $\{e_1^\delta, e_2^\delta, e_3^\bullet, e_4\}$  always returns a unique set of values. In our example, this is true in a ‘generic’ or ‘structural’ sense, meaning that it holds for all but exceptional values for the parameters in the considered equations (through the unspecified functions  $f_1$  and  $f_2$ ). We will repeatedly use *structural nonsingularity* (also called *structural regularity*) in the sequel. A formalization of structural reasoning and analysis is developed in Section 5. □

Observe that the same analysis could be applied to the original continuous-time dynamics from System (7) by augmenting it with the following *latent equation*:

$$\omega_1' - \omega_2' = 0 \quad (e_3') \quad (10)$$

obtained by differentiating  $(e_3)$ ; as a matter of fact, since  $(e_3)$  holds at any instant,  $(e_3')$  follows as long as the solution is smooth enough for the derivatives  $\omega_1'$  and  $\omega_2'$  to be defined. The resulting execution scheme is given in Exec. Sch. 2, that parallels Exec. Sch. 1.

---

**Execution Scheme 2** System (7)+Eq. (10).

---

**Require:** consistent  $\omega_1$  and  $\omega_2$ , i.e., satisfying  $(e_3)$ .

- |   |  |
|---|--|
| 1: Solve $\{e_1, e_2, e_3', e_4\}$ for $(\omega_1', \omega_2', \tau_1, \tau_2)$ | ▷ 4 equations, 4 unknowns                  |
| 2: ODESolve $(\omega_1, \omega_2)$  | ▷ update the states $(\omega_1, \omega_2)$ |
| 3: Tick   | ▷ move to next discretization step         |
- 

Line 1 is identical for the two schemes and is assumed to give a unique solution, generically; it fails if one omits the latent equation  $(e_3')$ . Then, although getting the next

values for the  $\omega_1$  and  $\omega_2$  is easy in Exec.Sch. 1, the situation changes in Exec.Sch. 2: the derivatives  $(\omega'_1, \omega'_2)$  are first evaluated, and then an ODE solver (here denoted by `ODESolve`) is used to update the state  $(\omega_1, \omega_2)$  for the next discretization step. Note that, when considering an exact mathematical solution, if  $\omega_1 - \omega_2 = 0$  initially holds and  $\omega'_1 - \omega'_2 = 0$ , then the linear constraint  $(e_3)$  will be satisfied at any positive time.

The replacement of  $(e_3)$  by  $(e'_3)$ , which resulted in Exec.Sch. 2, is known in the literature as *index reduction* [37]. It requires adding the (smallest set of) latent equations needed for the system in Line 1 of the execution scheme to become solvable and deterministic. The number of successive time differentiations needed for obtaining all the latent equations is called the *differentiation index* [21], which we simply refer to as the *index* in the sequel.<sup>5</sup> Finally, observe that both execution schemes rely on an algebraic equation system solver to evaluate the leading variables as a function of state variables.

We conclude this part by briefly discussing the initialization problem. Unlike for ODE systems, initialization is far from trivial for DAE systems. For the clutch example, if one considers System (7) as a standalone DAE, the initialization is performed as indicated in Exec.Sch. 3.

---

**Execution Scheme 3** Initialization of System (7)+Eq. (9).

---

1:  $(\omega_1, \omega_2) \leftarrow \text{Solve}\{e_3\}$  for  $(\omega_1, \omega_2)$  ▷ 2 unknowns for only 1 equation

---

The system for solving possesses two unknowns for only one equation, which leaves us with one degree of freedom; in mathematical terms, the set of all initial values for the 2-tuple of variables is a manifold of dimension 1, generically or structurally. For example, one may freely fix the initial common rotation speed so that  $(e_3)$  is satisfied.

## 2.2. Mode Transitions

In an attempt to reduce the full clutch model to the analysis of the DAE of each mode, one hopes that the handling of a mode change boils down to applying the initialization given in Exec.Sch. 3. If one was to handle restarts at mode changes like initializations, it would make the clutch system nondeterministic, due to the *extra* degree of freedom of Exec.Sch. 3. In contrast, the physics tells us that the state of the system is entirely determined when the clutch gets engaged. This comforts the intuition that resets at mode changes strongly differ from mere initializations.

However, inferring by hand the reset values for rotation velocities when the clutch gets engaged is definitely non-trivial. Furthermore, these values depend on the whole system model, so that the task of determining them becomes complex if external components are added. *It is therefore highly desirable, for this example, to let the compiler infer these reset values from model (3).*

This problem involves a mix of continuous-time and discrete-time dynamics, the latter consisting in one or several computation steps to restart the dynamics of the next mode. As a result, it will be convenient to cast everything in discrete time thanks to nonstandard analysis. We develop this next.

---

<sup>5</sup>There are indeed a great number of possible definitions for “the index” of a DAE, serving various purposes, see [21] for instance. This, however, is not relevant to the present work.

### 2.3. Nonstandard Semantics

Nonstandard analysis [49, 35, 5] extends the set  $\mathbb{R}$  of real numbers into a superset  ${}^*\mathbb{R}$  of *hyperreals* (also called *nonstandard reals*) that includes infinite sets of infinitely large numbers and infinitely small numbers. The key properties of hyperreals, needed for the informal discussion of the clutch example, are the following [35]:<sup>6</sup>

- There exist *infinitesimals*, defined as hyperreals that are smaller in absolute value than any real number: an infinitesimal  $\partial \in {}^*\mathbb{R}$  is such that  $|\partial| < a$  for any positive  $a \in \mathbb{R}$ . For  $x, y$  two hyperreals, write  $x \approx y$  if  $x - y$  is an infinitesimal.
- All relations, operators, and propositional formulas that are valid over  $\mathbb{R}$  are also valid over  ${}^*\mathbb{R}$ ; for example,  ${}^*\mathbb{R}$  is a totally ordered set. The arithmetic operations  $+$ ,  $\times$ , etc. can be lifted to  ${}^*\mathbb{R}$ . We say that a hyperreal  $x$  is *finite* if there exists some standard finite positive real number  $a$  such that  $|x| < a$ .
- For every finite hyperreal  $x \in {}^*\mathbb{R}$ , there exists a unique standard real number  $\text{st}(x) \in \mathbb{R}$  such that  $\text{st}(x) \approx x$ , and  $\text{st}(x)$  is called the *standard part* (or *standardization*) of  $x$ . As we shall see, standardizing more complex objects, such as functions or systems of equations, requires some care.
- Every real function lifts in a systematic way to a hyperreal function (regardless of its continuity properties).<sup>7</sup> This allows us to write  $f(x)$  where  $f$  is a real function and  $x$  is a nonstandard number.
- Let  $t \mapsto x(t), t \in \mathbb{R}$  be an  $\mathbb{R}$ -valued (standard) signal. Then:

$x$  is continuous at instant  $t \in \mathbb{R}$  if and only if, for any infinitesimal  $\partial \in {}^*\mathbb{R}$ , one has  $x(t + \partial) \approx x(t)$ ; (11)

$x$  is differentiable at instant  $t \in \mathbb{R}$  if and only if there exists  $a \in \mathbb{R}$  such that, for any infinitesimal  $\partial \in {}^*\mathbb{R}$ ,  $\frac{x(t+\partial)-x(t)}{\partial} \approx a$ . In this case,  $a = x'(t)$ . (12)

We can then consider, for a given positive infinitesimal  $\partial$ , the following time index set  $\mathbb{T} \subseteq {}^*\mathbb{R}$ :

$$\mathbb{T} = 0, \partial, 2\partial, 3\partial, \dots = \{n\partial \mid n \in {}^*\mathbb{N}\} \quad (13)$$

where  ${}^*\mathbb{N}$  denotes the set of *hyperintegers*, consisting of all natural integers augmented with additional infinite numbers called *nonstandard*. For convenience, the elements of  $\mathbb{T}$  will be generically denoted by the symbol  $\mathfrak{t}$ , the explicit expression  $n\partial$ , or even simply by  $t$  when the possible confusion with (usual) real time does not matter. The important features of  $\mathbb{T}$  are:<sup>8</sup>

- Any finite positive real time  $t \in \mathbb{R}, t \geq 0$  is infinitesimally close to some element of  $\mathbb{T}$  (informally,  $\mathbb{T}$  covers the set  $\mathbb{R}_+$  of nonnegative real numbers and can be used to index continuous-time dynamics); and

<sup>6</sup>The skeptical reader is gently referred to Section 7 for mathematical details.

<sup>7</sup>The result of this lifting may not be intuitive in general, see Section 7 for mathematical details.

<sup>8</sup>In French, one uses to qualify this as: *avoir le beurre et l'argent du beurre*.

- $\mathbb{T}$  is ‘discrete’: every instant  $n\partial$  has a predecessor  $(n-1)\partial$  and a successor  $(n+1)\partial$ , except 0 that has no predecessor.

**Definition 2 (Forward and Backward Shifts).** Let  $x$  be a nonstandard signal indexed by  $\mathbb{T}$ . We define the forward shifted signal  $x^\bullet$  through  $x^\bullet(n\partial) =_{\text{def}} x((n+1)\partial)$  and, for  $f(X)$  a function of the tuple  $X$  of signals, we set  $(f(X))^\bullet =_{\text{def}} f(X^\bullet)$  where the forward shift  $X \mapsto X^\bullet$  applies pointwise to all the components of the tuple. It will also be convenient to consider the backward shift  $\bullet x$  defined by  $\bullet x((n+1)\partial) =_{\text{def}} x(n\partial)$ , implying that an initial value for  $\bullet x(0)$  must be provided.

For example,  $f^\bullet(x, y)(t) =_{\text{def}} f(x^\bullet, y^\bullet)(t) = f(x(t+\partial), y(t+\partial))$ . Observe that the definition of the forward shift depends on the infinitesimal  $\partial$ . Thanks to (12), using forward shifts allows us to represent, up to an infinitesimal, the derivative  $x'$  of a signal by its first-order explicit Euler approximation  $\frac{1}{\partial}(x^\bullet - x)$ .

Solutions of multimode DAE systems may, however, be non-differentiable and even non-continuous at events of mode change. To give a meaning to  $x'$  at any instant, we decide to **define** it everywhere as the nonstandard first-order Euler increment

$$x' =_{\text{def}} \frac{1}{\partial}(x^\bullet - x) . \quad (14)$$

**Definition 3.** Substituting, in a multi-mode DAE system, every occurrence of  $x'(t)$  by its expansion (14) yields a difference algebraic equation (dAE) system<sup>9</sup> that we call its nonstandard semantics.

For instance, the nonstandard semantics of System (3) is the following:

$$\left\{ \begin{array}{ll} & \frac{\omega_1^\bullet - \omega_1}{\partial} = f_1(\omega_1, \tau_1) \quad (e_1^\partial) \\ & \frac{\omega_2^\bullet - \omega_2}{\partial} = f_2(\omega_2, \tau_2) \quad (e_2^\partial) \\ \text{if } \gamma \text{ then} & \omega_1 - \omega_2 = 0 \quad (e_3) \\ & \text{and } \tau_1 + \tau_2 = 0 \quad (e_4) \\ \text{if not } \gamma \text{ then} & \tau_1 = 0 \quad (e_5) \\ & \text{and } \tau_2 = 0 \quad (e_6) \end{array} \right. \quad (15)$$

The state variables are  $\omega_1, \omega_2$  whereas the leading variables are now  $\tau_1, \tau_2, \omega_1^\bullet, \omega_2^\bullet$ , in both modes  $\gamma = \text{F}$  and  $\gamma = \text{T}$ . Reproducing the reasoning of Section 2.1, one obtains, for each mode, a discrete system very much like the explicit Euler scheme of Section 2.1, except that the step size is now infinitesimal and that the variables are all nonstandard. The added value of System (15) with respect to the explicit Euler scheme of Section 2.1 is twofold: first, it is exact up to infinitesimals within each mode; second, it will allow us to carefully analyze what happens at events of modes change.

---

<sup>9</sup>Throughout this paper, we consistently use, in acronyms, the letters “D” and “d” to refer to “Differential” and “difference”, respectively.

#### 2.4. Nonstandard structural analysis

We now develop a structural analysis for the nonstandard multimode dAE system (15). As a preliminary step, we perform the structural analysis in each mode and combine the two dynamics with the due guards (the added latent equation is shown in red):

$$\left\{ \begin{array}{ll} & \frac{\omega_1^\bullet - \omega_1}{\partial} = f_1(\omega_1, \tau_1) \quad (e_1^\partial) \\ & \frac{\omega_2^\bullet - \omega_2}{\partial} = f_2(\omega_2, \tau_2) \quad (e_2^\partial) \\ \text{if } \gamma & \text{then } \omega_1 - \omega_2 = 0 \quad (e_3) \\ & \text{and } \omega_1^\bullet - \omega_2^\bullet = 0 \quad (e_3^\bullet) \\ & \text{and } \tau_1 + \tau_2 = 0 \quad (e_4) \\ \text{if not } \gamma & \text{then } \tau_1 = 0 \quad (e_5) \\ & \text{and } \tau_2 = 0 \quad (e_6) \end{array} \right. \quad (16)$$

System (16) coincides with the known structural analysis in the interior of each of the two modes  $\gamma = \text{T}$  and  $\gamma = \text{F}$ . What about the instants of mode change?

*Mode change*  $\gamma : \text{T} \rightarrow \text{F}$ . At the considered instant, we have  $\bullet\gamma = \text{T}$  and  $\gamma = \text{F}$ , where  $\bullet\gamma$  denotes the backward shift of  $\gamma$  following Definition 2. Unfolding System (16) at the two successive previous (shown in blue) and current (shown in black) instants yields, by taking the actual values for the guard at those instants into account:

$$\begin{array}{l} \text{previous} \\ \text{current} \end{array} \left\{ \begin{array}{ll} \frac{\omega_1^\bullet - \omega_1}{\partial} = f_1(\bullet\omega_1, \bullet\tau_1) \\ \frac{\omega_2^\bullet - \omega_2}{\partial} = f_2(\bullet\omega_2, \bullet\tau_2) \\ \omega_1 - \omega_2 = 0 & (\bullet e_3) \\ \omega_1 - \omega_2 = 0 \\ \bullet\tau_1 + \bullet\tau_2 = 0 \\ \frac{\omega_1^\bullet - \omega_1}{\partial} = f_1(\omega_1, \tau_1) \\ \frac{\omega_2^\bullet - \omega_2}{\partial} = f_2(\omega_2, \tau_2) \\ \tau_1 = 0 \\ \tau_2 = 0 \end{array} \right. \quad (17)$$

We assume that values are given for the state variables at the previous instant, namely  $\bullet\omega_1$  and  $\bullet\omega_2$ , and we regard System (17) as an algebraic system of equations with dependent variables  $\bullet\tau_i, \omega_i; \tau_i, \omega_i^\bullet$  for  $i = 1, 2$ , i.e., the leading variables of System (16) at the previous and current instants. Note that equation  $(\bullet e_3)$  does not involve any dependent variable: it is actually a *fact* that is inherited from the instant  $t - 2\partial$ , where  $t$  is the current instant; we can thus ignore this equation in System (17). Seen this way, System (17) is structurally nonsingular and we can solve it as two successive blocks, corresponding to the previous instant followed by the current instant: the latter is our desired code for the mode change. This yields the following code for the mode change  $\gamma : \text{T} \rightarrow \text{F}$ :

$$\left\{ \begin{array}{l} \omega_1, \omega_2 \text{ set by previous instant} \\ \frac{\omega_1^\bullet - \omega_1}{\partial} = f_1(\omega_1, \tau_1) \\ \frac{\omega_2^\bullet - \omega_2}{\partial} = f_2(\omega_2, \tau_2) \\ \tau_1 = 0 \\ \tau_2 = 0 \end{array} \right. \quad (18)$$

The above reasoning may seem an overshoot for this mode change  $\gamma : T \rightarrow F$ , since one could simply say: let the previous part of System (17) fix  $\omega_1, \omega_2$  and, then, let the current part of System (17) determine  $\tau_i, \omega_i^\bullet$  for  $i = 1, 2$ . However, the same line of reasoning will allow us to address the other, more difficult, mode change as well.

*Mode change*  $\gamma : F \rightarrow T$ . At the considered instant, we have  $\bullet\gamma = F$  and  $\gamma = T$ . We proceed as for the other case. Unfolding System (16) at the two successive previous (shown in blue) and current (shown in black) instants yields, by taking the actual values for the guard at those instants into account:

$$\begin{array}{l} \text{previous} \\ \text{current} \end{array} \left\{ \begin{array}{l} \frac{\omega_1 - \bullet\omega_1}{\partial} = f_1(\bullet\omega_1, \bullet\tau_1) \quad (\bullet e_1^\partial) \\ \frac{\omega_2 - \bullet\omega_2}{\partial} = f_2(\bullet\omega_2, \bullet\tau_2) \quad (\bullet e_2^\partial) \\ \bullet\tau_1 = 0 \\ \bullet\tau_2 = 0 \\ \frac{\omega_1^\bullet - \omega_1}{\partial} = f_1(\omega_1, \tau_1) \\ \frac{\omega_2^\bullet - \omega_2}{\partial} = f_2(\omega_2, \tau_2) \\ \omega_1 - \omega_2 = 0 \quad (e_3) \\ \omega_1^\bullet - \omega_2^\bullet = 0 \\ \tau_1 + \tau_2 = 0 \end{array} \right. \quad (19)$$

Once more, we regard System (19) as an algebraic system of equations with dependent variables  $\bullet\tau_i, \omega_i^\bullet; \tau_i, \omega_i$  for  $i = 1, 2$ , i.e., the leading variables of System (16) at the previous and current instants. In contrast to the previous mode change, System (19) is no longer structurally nonsingular since the following subsystem possesses five equations and only four dependent variables  $\omega_1, \omega_2, \bullet\tau_1, \bullet\tau_2$ :

$$\left\{ \begin{array}{l} \frac{\omega_1 - \bullet\omega_1}{\partial} = f_1(\bullet\omega_1, \bullet\tau_1) \quad (\bullet e_1^\partial) \\ \frac{\omega_2 - \bullet\omega_2}{\partial} = f_2(\bullet\omega_2, \bullet\tau_2) \quad (\bullet e_2^\partial) \\ \bullet\tau_1 = 0 \\ \bullet\tau_2 = 0 \\ \omega_1 - \omega_2 = 0 \quad (e_3) \end{array} \right. \quad (20)$$

This subsystem exhibits a conflict. We propose to resolve it by applying the following causality principle:

**Principle 1 (causality).** *What was done at the previous instant cannot be undone at the current instant.*

Applying Principle 1 leads to removing, from subsystem (20), the conflicting equation ( $e_3$ ). This yields the following code for the mode change  $\gamma : F \rightarrow T$ :

$$\left\{ \begin{array}{l} \omega_1, \omega_2, \bullet\tau_1, \bullet\tau_2 \text{ set by previous instant} \\ \frac{\omega_1^\bullet - \omega_1}{\partial} = f_1(\omega_1, \tau_1) \\ \frac{\omega_2^\bullet - \omega_2}{\partial} = f_2(\omega_2, \tau_2) \\ \omega_1^\bullet - \omega_2^\bullet = 0 \\ \tau_1 + \tau_2 = 0 \end{array} \right. \quad (21)$$

Note that the consistency equation  $(e_3) : \omega_1 - \omega_2 = 0$  has been removed from System (21), thus modifying the original model. However, this removal occurs only at mode change events  $\gamma : F \rightarrow T$ . What we have done amounts to *delaying by one nonstandard instant the satisfaction of some of the constraints in force in the new mode  $\gamma = T$* . Since our time step  $\partial$  is infinitesimal, this takes zero standard time.

The corresponding nonstandard execution scheme is summarized in Exec. Sch. 4. We use variable  $\Delta$  to encode the *context*, that is, the set of equations known to be satisfied by the state variables as a result of having executed the previous instant. For instance,  $e_3 \notin \Delta$  corresponds to the mode change  $\gamma : F \rightarrow T$ , expressing that  $(e_3)$  is not provably true if the guard was false at the previous instant. At each tick, that is, every time the scheme moves to a new instant, the context gets updated to account for the equations satisfied by the new state. The procedure ‘Solve’ solves the (algebraic) system to determine the values of the leading variables.

---

**Execution Scheme 4** for Nonstandard System (15). Comments are written in blue.

---

**Require:**  $\omega_1$  and  $\omega_2$ .

```

1: if  $\gamma$  (engaged mode) then
2:   if  $e_3 \notin \Delta$  (mode change) then
3:      $(\omega_1^\bullet, \omega_2^\bullet) \leftarrow \text{Solve} \{e_1^\partial, e_2^\partial, e_3^\bullet, e_4\}$ 
4:     Tick:  $\Delta \leftarrow \Delta \cup \{e_3\}$ 
5:   else
6:      $(\tau_1, \tau_2, \omega_1^\bullet, \omega_2^\bullet) \leftarrow \text{Solve} \{e_1^\partial, e_2^\partial, e_3^\bullet, e_4\}$ 
7:     Tick:  $\Delta$  unchanged
8:   else (released mode)
9:      $(\tau_1, \tau_2, \omega_1^\bullet, \omega_2^\bullet) \leftarrow \text{Solve} \{e_1^\partial, e_2^\partial, e_5, e_6\}$ 
10:    Tick:  $\Delta \leftarrow \Delta \setminus \{e_3\}$ 

```

---

Observe that Exec. Sch. 4 would work without changes if the guard  $\gamma$  was a predicate on the state variables  $\omega_1$  and  $\omega_2$ .

## 2.5. Standardization

Exec. Sch. 4 cannot be run in its present form, since it involves nonstandard reals. To recover executable code over the real numbers, a supplementary *standardization* step is needed. Recall that any finite nonstandard real  $x$  has a unique standard part  $\text{st}(x) \in \mathbb{R}$  such that  $x \approx \text{st}(x)$ . The standardization procedure aims at recovering the standard parts of the leading variables from their nonstandard version. We distinguish two cases: continuous evolutions within each mode, assuming that the sojourn time in each mode has (standard) positive duration, and discrete evolutions at events of mode change.

### 2.5.1. Within continuous modes

Let  $x : t \mapsto x(t)$ ,  $t \in [s, p)$ , denote the real continuous and differentiable solution in a given mode (assuming it exists). In the sequel, we use the notation  $o(\xi)$  to denote any (unspecified) nonstandard real number  $\zeta$  such that  $\zeta/\xi$  is infinitesimal. Using (12) and this notation, equations  $(e_1^\partial)$  and  $(e_2^\partial)$  rewrite as  $\omega'_i = f_i(\omega_i, \tau_i) + o(1)$ , which can be shown to standardize as the differential equations  $(e_1)$  and  $(e_2)$ , respectively. This suffices to recover the dynamics (6) in the released mode  $\gamma = F$ .

In the engaged mode  $\gamma = \text{T}$ , however, we need to handle  $(e_3^\bullet) : \omega_1^\bullet - \omega_2^\bullet = 0$ . The nonstandard characterization of derivatives writes  $\omega_i^\bullet = \omega_i + \partial.\omega_i' + o(\partial)$ , where  $\omega_i$  and  $\omega_i'$  are both standard. Since  $\omega_1 - \omega_2 = 0$  and  $\omega_1^\bullet - \omega_2^\bullet = 0$  both hold, we inherit  $\omega_1' - \omega_2' = o(1)$ , which implies  $\omega_1' - \omega_2' = 0$  since the  $\omega_i'$  are standard. We thus recover the dynamics (7) augmented with the latent equation (10) in the engaged mode  $\gamma = \text{T}$ .

To summarize, within continuous modes, we perform structural analysis as usual for DAE systems in continuous time: nothing new happens.

### 2.5.2. At the instants of mode change

Suppose we have an event of mode change at time  $t$ , meaning that  $\gamma(t) \neq \gamma(t - \partial)$ . Our aim is to use the next values  $\omega_i^\bullet(t) = \omega_i(t + \partial)$  to initialize the state variables for the next mode with the value  $\omega_i^+(t) = \omega_i^\bullet(t)$ . Thus, unlike for continuous time dynamics, the occurrence of the infinitesimal time step  $\partial$  in the expression  $\omega_i(t + \partial)$  is not an issue: it just points to the supplementary instant at which restart is performed.

However, the equations defining  $\omega_i^\bullet$  as functions of the  $\omega_i$ 's involve the hyperreal  $\partial$  *in space*: standardization must be performed to get rid of it.

*Transition  $\gamma : \text{T} \rightarrow \text{F}$ .* It is easy, as the new mode has an ODE dynamics whose state variables are initialized with corresponding exit values when leaving the previous mode.

*Transition  $\gamma : \text{F} \rightarrow \text{T}$ .* This one is more involved. As established in Exec. Sch. 4, in order to compute the reset values, we use the system of 4 equations  $\{e_1^\partial, e_2^\partial, e_3^\bullet, e_4\}$  to determine the leading variables  $(\tau_1, \tau_2, \omega_1^\bullet, \omega_2^\bullet)$ . In particular, from  $e_i^\partial$ , we get

$$\frac{\omega_i^\bullet - \omega_i}{\partial} = f_i(\omega_i, \tau_i), \quad i = 1, 2. \quad (22)$$

At this point, we must identify possible impulsive variables, which is the objective of the following *impulse analysis*.

*Impulse analysis.* Before engaging the clutch, we must assume  $\omega_1 - \omega_2 \neq 0$ . Since  $\omega_1^\bullet - \omega_2^\bullet = 0$  holds,  $\frac{(\omega_1^\bullet - \omega_2^\bullet) - (\omega_1 - \omega_2)}{\partial} = f_1(\omega_1, \tau_1) - f_2(\omega_2, \tau_2)$  cannot be a finite nonstandard real because, if it was, then the function  $t \mapsto \omega_1(t) - \omega_2(t)$  would be right-continuous, in contradiction with the above assumption. Hence, the nonstandard real  $f_1(\omega_1, \tau_1) - f_2(\omega_2, \tau_2)$  is necessarily infinite. However, we assumed continuous functions  $f_i$  and finite state  $(\omega_1, \omega_2)$ . Thus, one of the torques  $\tau_i$  must be infinite at  $t$ , and because of equation  $(e_4) : \tau_1 + \tau_2 = 0$ , both torques are in fact infinite, i.e., are *impulsive*.

*Eliminating impulsive variables.* It remains to compute the (standard) restart values for the state variables. To this end, we assume that the  $f_i$ 's are linear in the torques, i.e., each  $f_i$  has the following form:

$$f_i(\omega_i, \tau_i) = a_i(\omega_i) + b_i(\omega_i)\tau_i, \quad (23)$$

where  $b_1$  and  $b_2$  are the inverse moments of inertia of the rotating masses and  $a_1$  and  $a_2$  are damping factors divided by the corresponding moments of inertia. This yields the following system of equations, to be solved for  $\omega_1^\bullet, \omega_2^\bullet, \tau_1, \tau_2$  at the instant when  $\gamma$  switches from F to T:

$$\begin{cases} \omega_1^\bullet = \omega_1 + \partial.(a_1(\omega_1) + b_1(\omega_1)\tau_1) & (e_1^\partial) \\ \omega_2^\bullet = \omega_2 + \partial.(a_2(\omega_2) + b_2(\omega_2)\tau_2) & (e_2^\partial) \\ \omega_1^\bullet - \omega_2^\bullet = 0 & (e_3^\bullet) \\ \tau_1 + \tau_2 = 0 & (e_4) \end{cases} \quad (24)$$



It is tempting to standardize System (24) by simply setting  $\partial = 0$  in it. Unfortunately, doing this leaves us with a structurally singular system, since the two torques are then involved in only one equation.

**Comment 4. (standardizing functions vs. standardizing equations)** To get a better insight into this problem, let us replace for a while the infinitesimal  $\partial$  by a ‘small’ (standard) time step  $\delta$  and rewrite System (24) as  $F(\omega_1^\bullet, \omega_2^\bullet, \tau_1, \tau_2, \omega_1, \omega_2, \delta) = 0$ , where  $F=0$  stacks the four equations of (24) into one vector equation. Then, the function  $\delta \mapsto F(\omega_1^\bullet, \omega_2^\bullet, \tau_1, \tau_2, \omega_1, \omega_2, \delta)$  is continuous in  $\delta$  at  $\delta=0$ . Consequently, by (11), the standardization of  $F(\omega_1^\bullet, \omega_2^\bullet, \tau_1, \tau_2, \omega_1, \omega_2, \partial)$ , seen as a function, is obtained by setting  $\partial=0$  in it.

This, however, becomes incorrect if, instead of the function  $F$ , we are interested in the equation  $F=0$ . More precisely, in our example, let  $Y$  be the vector that collects the states  $\omega_1, \omega_2$ , and  $X$  be the vector that collects the dependent variables  $\omega_1^\bullet, \omega_2^\bullet, \tau_1, \tau_2$ ; then, the function  $F$  defined by System (24) is linear in  $\delta$ , i.e., it has the form  $Y \mapsto X$ , where  $X$  is solution of the algebraic equation

$$[A - \delta \times B(Y)]X = C(\delta, Y),$$

in which

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} ; \quad B = B(Y) = \begin{bmatrix} 0 & 0 & b_1(\omega_1) & 0 \\ 0 & 0 & 0 & b_2(\omega_2) \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

and vector  $C(\delta, Y)$  is appropriately fixed; from now on, we omit  $Y$  in the writing. For every  $\delta > 0$ , the matrix  $A - \delta \times B$  is invertible and we have  $X = (A - \delta \times B)^{-1}C$ .

If  $A$  was nonsingular, we could simply write  $X = (I - \delta \times D)^{-1}A^{-1}C$ , where  $D = A^{-1}B$ . Expanding the inverse around  $\delta = 0$  would yield  $(I - \delta \times D)^{-1} = I + \delta \times D + o(\delta)$ , implying  $X = (I + \delta \times D)A^{-1}C + o(\delta)$ . Therefore, the solution  $X(\partial)$  of equation  $(A - \partial \times B(Y))X = C(\partial, Y)$  would standardize as the solution  $X$  of the standard equation  $AX = C(0, Y)$ . All of this, however, requires that  $A$  be nonsingular, which is wrong in our case.<sup>10</sup>  $\square$

When  $A$  is singular (as in our case), the solution consists in eliminating the impulsive variables from System (24), namely, the two torques. This is performed by symbolic pivoting:

1. Using  $(e_4)$  yields  $-\tau_2 = \tau_1 =_{\text{def}} \tau$ ;
2. Premultiplying the system of equations

$$\begin{cases} \omega_1^\bullet = \omega_1 + \partial.(a_1(\omega_1) + b_1(\omega_1)\tau) & (e_1^\partial) \\ \omega_2^\bullet = \omega_2 + \partial.(a_2(\omega_2) - b_2(\omega_2)\tau) & (e_2^\partial) \end{cases}$$

by the row matrix  $[ \ b_2(\omega_2) \quad b_1(\omega_1) \ ]$  yields the equation

$$b_2(\omega_2)\omega_1^\bullet + b_1(\omega_1)\omega_2^\bullet = b_2(\omega_2)(\omega_1 + \partial.a_1(\omega_1)) + b_1(\omega_1)(\omega_2 + \partial.a_2(\omega_2)) .$$

---

<sup>10</sup>This informal reasoning will be formalized in Section 8.2, in the context of nonstandard analysis.

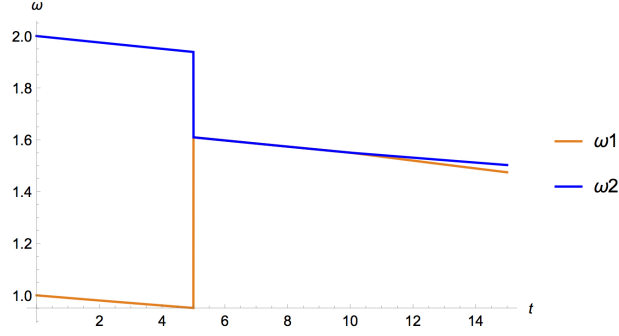


Figure 7: Simulation of the clutch model with resets, using our method. Mode change  $F \rightarrow T$  occurs at  $t = 5s$  and mode change  $T \rightarrow F$  occurs at  $t = 10s$ .

3. Using in addition  $(e_3^\bullet)$  and setting  $\omega^\bullet =_{\text{def}} \omega_1^\bullet = \omega_2^\bullet$  finally yields

$$\omega^\bullet = \frac{b_2(\omega_2)\omega_1 + b_1(\omega_1)\omega_2}{b_1(\omega_1) + b_2(\omega_2)} + \partial \frac{a_1(\omega_1)b_2(\omega_2) + a_2(\omega_2)b_1(\omega_1)}{b_1(\omega_1) + b_2(\omega_2)}. \quad (25)$$

The symbolic-only rewriting used here actually alleviates any difficulty in standardizing equation (25): this can now be done by just setting  $\partial = 0$  in its right-hand side. This yields, by identifying  $\text{st}(\omega_i) = \omega_i^-$  and  $\text{st}(\omega_i^\bullet) = \omega_i^+$ :

$$\omega_1^+ = \omega_2^+ = \frac{b_2(\omega_2^-)\omega_1^- + b_1(\omega_1^-)\omega_2^-}{b_1(\omega_1^-) + b_2(\omega_2^-)}, \quad (26)$$

that is, the reset value for the two velocities is the weighted arithmetic mean of  $\omega_1^-$  and  $\omega_2^-$ . Eq. (26) provides us with the reset values for the positions in the engaged mode, which is enough to restart the simulation in this mode. The actual impulsive values for the torques are useless and can be discarded. Note, however, that their sum explicitly remains equal to 0 during the instant of mode change. As a result, the momentum is preserved by the mode change, so that the obtained solution is satisfactory from a physical point of view.

As a final observation, instead of computing the exact standard part of  $\omega_i^\bullet$ , one could instead attempt to approximate it numerically by substituting  $\partial$  with a small (non-infinitesimal) step size  $\delta$  in System (24). This alternative approach is developed in Section 9.3. It has the advantage of not requiring the elimination of the impulsive variables.

We thus achieved the objectives stated at the end of Section 1.2.1. Figure 7 shows a simulation of the clutch model where the resets are computed following our approach developed hereabove. As one may expect, the reset value sits between the two values of  $\omega_1^-$  and  $\omega_2^-$  when  $\gamma : F \rightarrow T$  (at  $t = 5s$ ), and the transition is continuous at the second reset (at  $t = 10s$ ).

### 3. The Cup-and-Ball example

In this section, we study the example of Section 1.2.2 and address the objectives stated at the end of that section.

### 3.1. Rejecting the submitted model

We first explain why we reject the model proposed in Section 1.2.2, and how we do so. Let us consider again System (5).

*Nonstandard semantics.* Following Section 2.3 for the clutch, the derivatives  $x', y', x'', y''$  are replaced using the explicit first-order Euler scheme with an infinitesimal time step  $\partial$ :

$$x' \leftarrow \frac{x^\bullet - x}{\partial} \ ; \ x'' \leftarrow \frac{x^{\bullet 2} - 2x^\bullet + x}{\partial^2} \ . \quad (27)$$

In the sequel, when referring to model (5), substitution (27) will be implicit.

*Structural analysis.*  $x, y, x^\bullet, y^\bullet$  are the state variables, the value of which is known from previous nonstandard instants. The system has guard  $\gamma$ . In contrast to the clutch example,  $\gamma$  is no longer an input, but is rather evaluated based on system variables using equation ( $k_0$ ). For both modes, the leading variables are  $x^{\bullet 2}, y^{\bullet 2}, \lambda$ ; the first two of these arise from the second derivatives  $x'', y''$  by using substitution (27).

**Comment 5. (logico-numerical fixpoints)** We could consider the direct solving of System (5), seen as a mixed logico-numerical system of equations with dependent variables  $x^{\bullet 2}, y^{\bullet 2}, \lambda, \gamma$ . However, this system mixes guarded numerical equations with the evaluation of a predicate. A possible solution would consist in performing a relaxation, by iteratively updating the numerical variables based on the previous value for the guards, and then re-evaluating the guard based on the updated values of the numerical variables, hoping for a fixpoint to occur. Such fixpoint equation, however, can have zero, one, several, or infinitely many solutions. No characterization exists that could serve as a basis for a (graph-based) structural analysis. We thus decided not to try to solve such mixed logico-numerical systems.  $\square$

As a consequence, we are unable to evaluate guard  $\gamma$ , so that the mode the system is in cannot be determined. Should we give up? Not yet.

It could be that, based on information obtained by solving the subsystem consisting of the equations that are always enabled, guard  $\gamma$  can be evaluated. In this model, this subsystem is  $G = \{(e_1), (e_2)\}$ , with dependent variables  $x^{\bullet 2}, y^{\bullet 2}, \lambda$ . Unfortunately, subsystem  $G$  is underdetermined: it cannot be solved, and the evaluation of guard  $\gamma$  is still impossible. We thus reject model (5), because we are unable to evaluate guard  $\gamma$  on the basis of the subsystem that is enabled prior to evaluating this guard.

**Comment 6. (forbidding logico-numerical fixpoints based on syntax)** The reader may wonder why we did not simply reject the model right from the beginning, on the sole basis of the existence of the logico-numerical fixpoint equation. The rationale is that the method discussed above allows, in principle, more models to get accepted than just the ones without such equations. This design issue will be addressed more thoroughly in Comment 30 of Section 6.1.2.  $\square$

### 3.2. Correcting the model

To break the fixpoint equation defining  $\gamma$ , the idea is to replace the predicate  $s \leq 0$  defining the guard  $\gamma$  by its left-limit:  $\gamma = [s^- \leq 0]$ , where  $s^-(t) =_{\text{def}} \lim_{u \nearrow t} s(u)$ :

$$\left\{ \begin{array}{ll} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ \gamma = [s^- \leq 0] & (k_0) \\ \text{if } \gamma \text{ then } 0 = L^2 - (x^2 + y^2) & (k_1) \\ \quad \text{and } 0 = \lambda + s & (k_2) \\ \text{if not } \gamma \text{ then } 0 = \lambda & (k_3) \\ \quad \text{and } 0 = (L^2 - (x^2 + y^2)) - s & (k_4) \end{array} \right.$$

We map this model to nonstandard semantics by stating that predicate  $s \leq 0$  defines the value of the guard *at the next nonstandard instant*—therefore, an initial condition for guard  $\gamma$  is required; we initialize the system in free motion mode  $\gamma(0) = \text{F}$ . This yields the corrected model (28), where the modification is highlighted in red:

$$\left\{ \begin{array}{ll} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ \gamma^\bullet = [s \leq 0]; \gamma(0) = \text{F} & (k_0) \\ \text{if } \gamma \text{ then } 0 = L^2 - (x^2 + y^2) & (k_1) \\ \quad \text{and } 0 = \lambda + s & (k_2) \\ \text{if not } \gamma \text{ then } 0 = \lambda & (k_3) \\ \quad \text{and } 0 = (L^2 - (x^2 + y^2)) - s & (k_4) \end{array} \right. \quad (28)$$

This model is understood in the nonstandard semantics, meaning that the derivatives are expanded using (27). Therefore, the leading variables in all modes are  $\lambda, s, x^{\bullet 2}, y^{\bullet 2}$ .

### 3.3. Nonstandard structural analysis

This section focuses on the core difficulty, namely the mode change  $\gamma : \text{F} \rightarrow \text{T}$ , i.e., the event when the rope gets straight. Due to equation (k<sub>1</sub>), the mode  $\gamma = \text{T}$  (where the rope is straight) requires index reduction. We thus augment model (28) with the two latent equations shown in red:

$$\left\{ \begin{array}{ll} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ \gamma^\bullet = [s \leq 0]; \gamma(0) = \text{F} & (k_0) \\ \text{if } \gamma \text{ then } 0 = L^2 - (x^2 + y^2) & (k_1) \\ \quad \text{and } 0 = L^2 - (x^2 + y^2)^\bullet & (k_1^\bullet) \\ \quad \text{and } 0 = L^2 - (x^2 + y^2)^{\bullet 2} & (k_1^{\bullet 2}) \\ \quad \text{and } 0 = \lambda + s & (k_2) \\ \text{if not } \gamma \text{ then } 0 = \lambda & (k_3) \\ \quad \text{and } 0 = (L^2 - (x^2 + y^2)) - s & (k_4) \end{array} \right. \quad (29)$$

In this model, the two equations  $(k_1)$  and  $(k_1^\bullet)$  are in conflict with equations from the previous two instants, shown in blue in the following subsystem:

$$\left\{ \begin{array}{ll} 0 = \frac{x-2\dot{x}+\ddot{x}}{\partial^2} + \bullet^2\lambda \bullet^2x & (\bullet^2e_1) \\ 0 = \frac{y-2\dot{y}+\ddot{y}}{\partial^2} + \bullet^2\lambda \bullet^2y + g & (\bullet^2e_2) \\ 0 = \frac{x-2\dot{x}+\ddot{x}}{\partial^2} + \lambda \bullet x & (\bullet e_1) \\ 0 = \frac{y-2\dot{y}+\ddot{y}}{\partial^2} + \lambda \bullet y + g & (\bullet e_2) \\ 0 = L^2 - (x^2 + y^2) & (k_1) \\ 0 = L^2 - (x^2 + y^2)^\bullet & (k_1^\bullet) \end{array} \right.$$

Applying again Principle (1) of causality, we erase, in model (29), equations  $(k_1)$  and  $(k_1^\bullet)$  at the instant of mode change  $\bullet\gamma=\text{F}, \gamma=\text{T}$ . This yields the following system:

$$\text{at } \left[ \begin{array}{l} \bullet\gamma=\text{F} \\ \gamma=\text{T} \end{array} \right] : \left\{ \begin{array}{ll} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ 0 = L^2 - (x^2 + y^2)^\bullet & (k_1^{\bullet 2}) \\ 0 = \lambda + s & (k_2) \end{array} \right. \quad (30)$$

It uniquely determines all the leading variables from the state variables  $x, y$  and  $x^\bullet, y^\bullet$ . In turn, equations  $(k_1)$  and  $(k_1^\bullet)$ , which were erased from this model, are not satisfied.

Also, we erase, in model (29), the only equation  $(k_1)$  at the next instant, i.e., when  $\bullet^2\gamma=\text{F}, \bullet\gamma=\text{T}, \gamma=\text{T}$ . This yields the following system:

$$\text{at } \left[ \begin{array}{l} \bullet^2\gamma=\text{F} \\ \bullet\gamma=\text{T} \\ \gamma=\text{T} \end{array} \right] : \left\{ \begin{array}{ll} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ 0 = L^2 - (x^2 + y^2)^\bullet & (k_1^{\bullet 2}) \\ 0 = L^2 - (x^2 + y^2)^\bullet & (k_1^{\bullet 2}) \\ 0 = \lambda + s & (k_2) \end{array} \right. \quad (31)$$

This uniquely determines all the leading variables, given values for the state variables  $x, y$  and  $x^\bullet, y^\bullet$ . Note that  $(k_1^{\bullet 2})$  is a consistency equation that is satisfied by the state variables  $x^\bullet, y^\bullet$ . In turn, equation  $(k_1)$ , which was erased from this model, is not satisfied. At subsequent instants, equation erasure is no longer needed. This completes the nonstandard structural analysis of the mode change  $\gamma : \text{F} \rightarrow \text{T}$ , i.e., when the rope gets straight.

### 3.4. Standardization

Code generation for restarts consists in standardizing nonstandard systems (30) and (31). Remember, from our development of the clutch example, that standardizing systems of equations requires more care than standardizing numbers, due to impulsive behaviors and singularity issues that result. The simplest method from a conceptual point of view is to eliminate impulsive variables from the restart system, as they are of no use for restarting the new mode.

We focus on the standardization of the mode change  $\gamma : \text{F} \rightarrow \text{T}$ , i.e., when the rope gets straight. Our task is to standardize systems (30) and (31), by targeting discrete-time dynamics, for the two successive instants composing the restart phase. This will provide us with restart values for positions and velocities.

Due to the expansion of derivatives in equations  $(e_1, e_2, e_1^\bullet, e_2^\bullet)$ , tensions  $\lambda$  and  $\lambda^\bullet$  are both impulsive, hence so are  $s$  and  $s^\bullet$  by  $(k_2, k_2^\bullet)$ . We eliminate the impulsive variables by ignoring  $(k_2, k_2^\bullet)$ , combining  $(e_1)$  and  $(e_2)$  to eliminate  $\lambda$ , and  $(e_1^\bullet)$  and  $(e_2^\bullet)$  to eliminate  $\lambda^\bullet$ . This yields:

$$\text{at } \begin{bmatrix} \bullet\gamma=\text{F} \\ \gamma=\text{T} \end{bmatrix} : \begin{cases} 0 = y''x + gx - x''y \\ 0 = L^2 - (x^2 + y^2)^{\bullet 2} \end{cases} \quad (32)$$

$$\text{at } \begin{bmatrix} \bullet 2\gamma=\text{F} \\ \bullet\gamma=\text{T} \\ \gamma=\text{T} \end{bmatrix} : \begin{cases} 0 = y''x + gx - x''y \\ 0 = L^2 - (x^2 + y^2)^\bullet \\ 0 = L^2 - (x^2 + y^2)^{\bullet 2} \end{cases} \quad (33)$$

In System (32), we expand second derivatives using (27). Note that substitution (27) implies that we can as well expand the second derivatives as functions of first derivatives:

$$x'' \leftarrow \frac{x'^\bullet - x'}{\partial} . \quad (34)$$

Since we will use System (33) to evaluate restart values for the velocities, we will use expansion (34) in it. Consequently, System (32) has dependent variables  $x^{\bullet 2}, y^{\bullet 2}$ , whereas System (33) has dependent variables  $x'^\bullet, y'^\bullet$ . We are now ready to standardize the two systems.

*System (32) to define restart positions.* We expand second derivatives using (27):

$$\begin{cases} 0 = (y^{\bullet 2} - 2y^\bullet + y)x - (x^{\bullet 2} - 2x^\bullet + x)y + \partial^2 gx \\ 0 = L^2 - (x^2 + y^2)^{\bullet 2} \end{cases} \quad (35)$$

Regarding (35), one can safely set  $\partial = 0$  to standardize it, since the resulting system is still structurally regular—see Comment 4 regarding the clutch example. The standardization of (35) is thus:

$$\begin{cases} 0 = (y^{\bullet 2} - 2y^\bullet + y)x - (x^{\bullet 2} - 2x^\bullet + x)y \\ 0 = L^2 - (x^2 + y^2)^{\bullet 2} \end{cases} \quad (36)$$

In contrast, had we directly set  $\partial = 0$  in System (30) (without eliminating impulsive variable  $\lambda$ ), we would get a structurally singular system, an incorrect standardization.

In System (36), we can interpret  $x$  and  $x^\bullet$  as the left-limit  $x^-$  of state variable  $x$  in previous mode, and  $x^{\bullet 2}$  as the restart value  $x^+$  for the new mode. This yields

$$\begin{cases} 0 = (y^+ - y^-)x^- - (x^+ - x^-)y^- \\ 0 = L^2 - (x^2 + y^2)^+ \end{cases} \quad (37)$$

which determines the restart values for positions. Note that the constraint that the rope is straight is satisfied. Furthermore, if we remember that the rope is straight at mode change (this is the mode change condition), then  $0 = L^2 - (x^2 + y^2)^-$  also holds, which ensures that  $x^+ = x^-, y^+ = y^-$  is the unique solution of (37): positions are continuous.

System (33) to define restart velocities. We expand second derivatives using (34):

$$\begin{cases} 0 = (y'^{\bullet} - y')x - (x'^{\bullet} - x')y + \partial \cdot gx \\ 0 = L^2 - (x^2 + y^2)^{\bullet} \\ 0 = L^2 - (x^2 + y^2)^{\bullet 2} \end{cases} \quad (38)$$

By expanding  $x^{\bullet 2} = x^{\bullet} + \partial x'^{\bullet}$ , the right-hand side of the last equation rewrites

$$\begin{aligned} L^2 - (x^2 + y^2)^{\bullet 2} &= L^2 - (x^2 + y^2)^{\bullet} \\ &\quad + 2\partial(x^{\bullet}x'^{\bullet} + y^{\bullet}y'^{\bullet}) \\ &\quad + \partial^2((x'^{\bullet})^2 + (y'^{\bullet})^2) \\ &= 0 \quad (\text{by second equation of (38)}) \\ &\quad + 2\partial(x^{\bullet}x'^{\bullet} + y^{\bullet}y'^{\bullet}) \\ &\quad + O(\partial^2) \end{aligned} \quad (39)$$

Using this expansion of  $L^2 - (x^2 + y^2)^{\bullet 2}$ , setting  $\partial = 0$  in (38) yields

$$\begin{cases} 0 = (y'^{\bullet} - y')x - (x'^{\bullet} - x')y \\ 0 = x^{\bullet}x'^{\bullet} + y^{\bullet}y'^{\bullet} \end{cases} \quad (40)$$

where the dependent variables are now  $x'^{\bullet}, y'^{\bullet}$ , whereas other variables are state variables whose values are determined by previous time steps. Note that System (40) is structurally regular, hence it is the correct standardization of System (38). We are now ready to get restart values for velocities. In System (40), we perform the following substitutions, recalling that superscripts  $-$  and  $+$  denote left- and right-limits:

$$x = x^- ; x^{\bullet} = x^+ ; x' = x'^- \quad (41)$$

and the restart velocity is evaluated by

$$x'^+ = x'^{\bullet} . \quad (42)$$

After similar substitutions for  $y$ , we get

$$\begin{cases} 0 = (y'^+ - y'^-)x^- - (x'^+ - x'^-)y^- \\ 0 = x^+x'^+ + y^+y'^+ \end{cases} \quad (43)$$

System (43) determines  $x'^+$  and  $y'^+$ , which are the velocities for restart. The second equation guarantees that the velocity will be tangent to the constraint. With (37) and (43), we determine the restart conditions for positions and velocities. Invariants from the physics are satisfied.

Our reasoning so far produces a behavior in which the two modes (free motion and straight rope) gently alternate; the system always stays in one mode for some positive period of time before switching to the other mode. *This indeed amounts to assuming that the impact is totally inelastic at mode change, an assumption that was not explicit at all in (28).* So, what happened?

### 3.5. Handling transient modes

In fact, the straight rope mode was implicitly assumed to last for at least three non-standard successive instants, since we allowed ourselves to shift  $(k_1)$  twice.

Now, let us instead assume elastic impact. Hence, the cascade of mode changes  $\gamma : F \rightarrow T \rightarrow F$  is possible in three successive nonstandard instants, reflecting that the straight rope mode is *transient* (it is left immediately after being reached). We wish to adapt our reasoning to handle this transient mode.

Consider again model (28). We regard the instant of the cascade when  $\gamma = T$  occurs as the current instant. We cannot add latent equations by simply shifting  $(k_1)$ , since these shifted versions are not active in the mode  $\gamma = F$ . Instead, our line of reasoning goes as follows. Set

$$\begin{aligned} S(T) &= \{(e_1), (e_2), (k_1), (k_2)\} , \\ S(F) &= \{(e_1), (e_2), (k_3), (k_4)\} . \end{aligned}$$

Systems  $S^\bullet(T)$  and  $S^\bullet(F)$  are obtained by shifting the equations constituting  $S(T)$  and  $S(F)$ ; systems  $S^{\bullet 2}(T)$  and  $S^{\bullet 2}(F)$  are defined similarly.

Now, the idea is to consider the *differentiation array* originally proposed by Campbell and Gear [21], except that we take into account the trajectory  $T, F, F, \dots$  for guard  $\gamma$ . Using shifting instead of differentiation yields the following *array of shifted systems*, where size  $n$  is a parameter:

$$\mathcal{A}_n(S) \stackrel{\text{def}}{=} \begin{bmatrix} S(T) \\ S^\bullet(F) \\ S^{\bullet 2}(F) \\ \vdots \\ S^{\bullet n}(F) \end{bmatrix} \quad (44)$$

We search for an  $n$  such that  $\mathcal{A}_n(S)$  provides us with the due latent equations. Thus, the question is: can we find latent equations for  $S(T)$  by shifting appropriate equations from  $S(F)$ ? Unfortunately, the answer is no: although shifting twice  $(k_4)$  in System (28) produces one more equation involving the leading variables  $x^{\bullet 2}, y^{\bullet 2}$ , this equation also involves the new variable  $s^{\bullet 2}$ , which keeps the augmented system underdetermined; shifting other equations fails as well.

Therefore, the structural analysis rejects this model as being underdetermined at transient mode  $\gamma = T$ . The user is then asked to provide one more equation. For example, he/she could specify an impact law for the velocity  $y'$  by providing the equation

$$(y')^+ = -(1 - \alpha)(y')^- ,$$

where  $0 \leq \alpha < 1$  is a fixed damping coefficient. This is reinterpreted in the nonstandard domain as

$$y'^{\bullet} = -(1 - \alpha)y' , \quad (45)$$

yielding the following refined system for use at mode  $\gamma=T$  within the cascade  $\gamma:F \rightarrow T \rightarrow F$ :

$$\left\{ \begin{array}{ll} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ 0 = y'^{\bullet} + (1 - \alpha)y' & (\tau_1) \\ 0 = L^2 - (x^2 + y^2) & (k_1) \\ 0 = \lambda + s & (k_2) \end{array} \right. \quad (46)$$



We proceed again with the structural analysis. Variables  $x, y$  are the states, so that their values are set by the previous instants. Current equation ( $k_1$ ) creates a conflict with the past. Hence, we discard it from System (46), which leaves us with the following system:

$$\begin{cases} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ 0 = y'^{\bullet} + (1 - \alpha)y' & (\tau_1) \\ 0 = \lambda + s & (k_2) \end{cases} \quad (47)$$

Model (47) is structurally nonsingular, recalling that  $y''$  and  $y'^{\bullet}$  can be interchanged for the structural analysis. This refined model is therefore accepted. Standardization proceeds accordingly, with no new difficulty.

### 3.6. Consequences for the modeling language

Through the Cup-and-Ball example, we demonstrated the need for the following user-given information: is the current mode long or transient? We argue that this information cannot be derived from inspecting the system model.

In the clutch example, guard  $\gamma$  was a system input, and we implicitly assumed that each of the two modes was long. Had we composed the clutch model with a ‘feedback’ determining  $\gamma$  as a predicate over system variables, this implicit assumption would have become questionable. This situation indeed occurs with the Cup-and-Ball example, where the elastic/inelastic impact assumption is a side information that cannot be deduced from the equation syntax. This leads to the following important observation:

**Comment 7. (Long / Transient)** Long / Transient is an information regarding modes, that cannot be found by inspecting the equations of the system. It must be inferred from understanding the system physics. Therefore, it has to be provided as an extra information by the model designer.  $\square$

The natural way of addressing Comment 7 is to provide a different syntax for specifying long modes on the one hand, and events corresponding to transient modes on the other hand. In contrast, mode changes separating two successive long modes need not be specified. A candidate syntax is:

long mode	if predicate then $S$	(48)
transient mode	when predicate then $S$	

The ‘when’ statement of Modelica, which we reuse here, points to the event when ‘predicate’ switches from false to true. For this ‘when’ statement, we could furthermore restrict ‘predicate’ to be a zero-crossing condition, by which a real-valued expression crosses zero from below [15]. We devote the ‘if’ statement to long-lasting modes specified by ‘predicate’.

To illustrate this, our syntax offering if-equations and when-equations allows specifying the Cup-and-Ball example with elastic impact as shown in Figure 8.

## 4. Issues and systematic approach

From our discussion of the examples, we see that the key issue is the handling of mode changes, and particularly the reconciliation of the dynamics just before the change,

$$\left\{ \begin{array}{ll} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ \gamma = [s^- \leq 0]; \gamma(0) = F & (k_0) \\ \text{when } \gamma \text{ then } y' = -\alpha y'^- & (\tau_1) \\ \text{if not } \gamma \text{ then } 0 = \lambda & (k_3) \\ \text{and } 0 = (L^2 - (x^2 + y^2)) - s & (k_4) \end{array} \right.$$

Figure 8: Model of the Cup-and-Ball example with elastic impact specified by equation  $(\tau_1)$ , highlighted in red. This equation replaces the two equations  $(k_1), (k_2)$  specifying the straight rope mode in case of inelastic impact.

and the consistency conditions set by the new mode. How to properly address this in the original model of continuous real-time was an open question; mapping the system dynamics to the nonstandard domain clarified this task.

In the rest of this paper, we extend our method, from the toy examples to a systematic approach for mDAE system compilation. This approach, illustrated in Figure 9, decomposes into several steps.

- First, the mDAE system is transformed into a system of *multi-mode difference Algebraic Equations* (mdAE) using the nonstandard Euler expansion of derivatives  $x' \leftarrow \frac{x^\bullet - x}{\delta}$ . The resulting mdAE system is the *nonstandard semantics* of the original model.
- Second, the structural analysis of this nonstandard semantics is performed, resulting in a new mdAE system where latent equations are made explicit and conflicts at mode changes are solved.
- Finally, the standardization step recovers both the smooth dynamics in each mode, and the possibly discontinuous/impulsive restart conditions at mode changes. The resulting code can then be executed with the help of a solver.

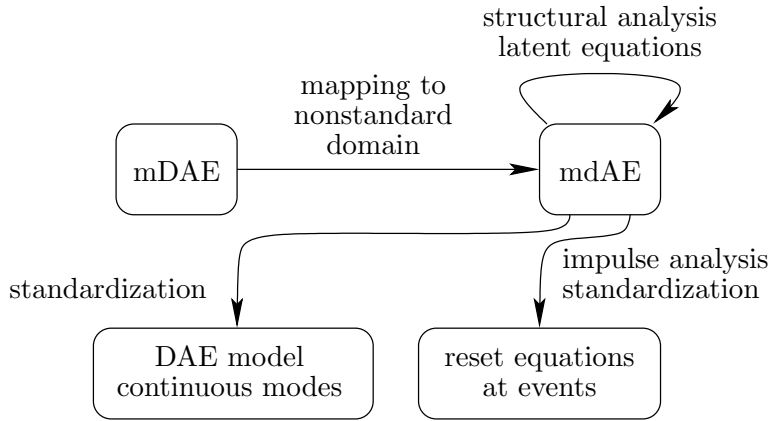


Figure 9: Our approach to the compilation of mDAE systems.

The intuitions developed on the examples are formalized in the sequel of the paper. We first recall the background on structural analysis and we extend it to our context. We then focus on standardization.

## 5. Background on Structural Analysis

We recall in this section the basics of structural analysis for algebraic systems of equations and provide in passing a few new lemmas for subsequent use. We then present the  $\Sigma$ -method, which is the structural analysis method we adopt for DAE systems.

**Notations 8.** In this paper,  $\mathbb{N} = \{0, 1, 2, \dots\}$  is the set of nonnegative natural integers and  $\mathbb{R}_+ = [0, \infty)$  is the set of nonnegative real numbers.  $\square$

We will sometimes need to handle variables, their valuations, and vectors thereof. To this end, the following conventions will be used (unless no confusion occurs from using more straightforward notations):

**Notations 9.** Lowercase letters ( $x, y$ , etc.) denote scalar real variables; capitals ( $X, Y$ , etc.) denote vectors of real variables; whenever convenient, we regard  $X$  as a set of variables and write  $x \in X$  to refer to an entry of  $X$ . We adopt similar conventions for functions ( $f, g$ , etc.) and vectors of functions ( $F, G$ , etc.). A value for a variable  $x$  is generically denoted by  $\nu_x$ , and similarly for  $X$  and  $\nu_X$ .  $\square$

### 5.1. Structural analysis of algebraic equations

In this section we focus on systems of algebraic equations, that is, equations involving no time and no dynamics.

#### 5.1.1. Structural nonsingularity of algebraic equations

Consider a system of smooth algebraic equations:

$$f_i(y_1, \dots, y_k, x_1, \dots, x_n) = 0, \quad i = 1, \dots, m \quad (49)$$

rewritten as  $F(Y, X) = 0$ , where  $Y$  and  $X$  denote the vectors  $(y_1, \dots, y_k)$  and  $(x_1, \dots, x_n)$ , respectively, and  $F$  is the vector  $(f_1, \dots, f_m)$ . The system has  $m$  equations,  $k$  input variables collected in vector  $Y$ , and  $n$  *dependent variables* (or unknowns) collected in vector  $X$ . Throughout this section, we assume that the  $f_i$ 's are all  $\mathcal{C}^1$  at least.

If the system is square, i.e., if  $m = n$ , the *Implicit Function Theorem* (see, e.g., Theorem 10.2.2 in [24]) states that, if  $(\nu_Y, \nu_X) \in \mathbb{R}^{k+n}$  is a value for the pair  $(Y, X)$  such that  $F(\nu_Y, \nu_X) = 0$  and the Jacobian matrix of  $F$  with respect to  $X$  evaluated at  $(\nu_Y, \nu_X)$  is nonsingular, then there exists, in an open neighborhood  $U$  of  $\nu_Y$ , a unique vector of functions  $G$  such that  $F(v, G(v)) = 0$  for all  $v \in U$ . In words, Eq. (49) uniquely determines  $X$  as a function of  $Y$ , locally around  $\nu_Y$ . Denote by  $\mathbf{J}_X F$  the above mentioned Jacobian matrix. Solving  $F = 0$  for  $X$ , given a value  $\nu_Y$  for  $Y$ , requires forming  $\mathbf{J}_X F(\nu_Y)$  as well as inverting it.

One could avoid considering  $\mathbf{J}_X F$  by focusing on *structural* nonsingularity, a weaker but much cheaper criterion previously introduced for handling sparse matrices in high performance computing.

### 5.1.2. Structural analysis

We begin with some background on graphs, for which a basic reference is [11]. Given a graph  $\mathcal{G} = (V, E)$ , where  $V$  and  $E$  are the sets of vertices and edges, a *matching*  $\mathcal{M}$  of  $\mathcal{G}$  is a set of edges of  $\mathcal{G}$  such that no two edges in  $\mathcal{M}$  are incident on a common vertex. Matching  $\mathcal{M}$  is called *complete* if it covers all the vertices of  $\mathcal{G}$ . An  $\mathcal{M}$ -*alternating path* is a path of  $\mathcal{G}$  whose edges are alternatively in  $\mathcal{M}$  and not in  $\mathcal{M}$ —we say simply “alternating path” when  $\mathcal{M}$  is understood from the context. A vertex is said to be *matched* if there is an edge in  $\mathcal{M}$  that is incident on the vertex, and *unmatched* otherwise.

Let  $F(Y, X)=0$  be a system of algebraic equations of the form (49); its *bipartite graph*  $\mathcal{G}_F$  is the graph having  $F \cup X$  as set of vertices and an edge  $(f_i, x_j)$  if and only if variable  $x_j$  occurs in function  $f_i$ .

*Square systems.* We first consider square systems, in which equations and dependent variables are in equal numbers.

**Definition 10 (structural nonsingularity).** *System  $F(Y, X)=0$  is called structurally nonsingular if its bipartite graph  $\mathcal{G}_F$  possesses a complete matching.*

A structurally nonsingular system is necessarily square, i.e., with equations and variables in equal numbers. The link to numerical regularity is formalized in the following lemma (implicitly used in [40, 37, 46, 48]):

**Lemma 11.** *Assume that  $F$  is  $\mathcal{C}^1$ . The following properties are equivalent:*

1. *System  $F(Y, X)=0$  is structurally nonsingular.*
2. *For every  $(\nu_Y, \nu_X)$  satisfying  $F(\nu_Y, \nu_X)=0$ , the Jacobian matrix  $\mathbf{J}_X F(\nu_Y, \nu_X)$  remains generically<sup>11</sup> nonsingular when its nonzero coefficients vary over some neighborhood.*

Let us detail Property 2. We can represent graph  $\mathcal{G}_F$  by the  $n \times n$ -incidence matrix  $M_F = (m_{ij})$ , having a 1 at  $(i, j)$  if  $(f_i, x_j)$  is an edge of  $\mathcal{G}_F$ , and a 0 otherwise. We order the entries  $\{a_{ij} \mid m_{ij} = 1\}$  of Jacobian  $\mathbf{J}_X F(\nu_Y, \nu_X)$ , for example by lexicographic order over the pair  $(i, j)$ . In this way, every  $n \times n$ -matrix whose pattern of nonzero coefficients is  $M_F$  identifies with a unique vector of  $\mathbb{R}^K$ , where  $K$  is the number of 1’s in  $M_F$ . We denote by  $\mathcal{J}_X F \in \mathbb{R}^K$  the image of the Jacobian  $\mathbf{J}_X F(\nu_Y, \nu_X)$  obtained via this correspondence. Property 2 says:

There exist an open neighborhood  $U$  of  $\mathcal{J}_X F$  in  $\mathbb{R}^K$ , and a subset  $V \subseteq U$  such that: (i) the set  $U \setminus V$  has zero Lebesgue measure, and (ii) every  $\mathcal{J} \in V$  yields a regular matrix. (50)

In the sequel, the so defined sets  $U$  and  $V$  will be denoted by

$$U_F(\nu_Y, \nu_X) \quad \text{and} \quad V_F(\nu_Y, \nu_X) \tag{51}$$

or simply  $U_F$  and  $V_F$  when no confusion can result.

*Practical use of structural nonsingularity.* If a square matrix is structurally singular, then it is singular. The converse is false: structural nonsingularity does not guarantee

---

<sup>11</sup>Generically means: outside a set of values of Lebesgue measure zero.

nonsingularity. Therefore, the practical use of Definition 10 is as follows: We first check if  $F=0$  is structurally nonsingular. If not, then we abort searching for a solution. Otherwise, we then proceed to computing a solution, which may or may not succeed depending on the actual numerical regularity of the Jacobian matrix  $\mathbf{J}_X F$ .  $\square$

Let  $F(Y, X)=0$  be structurally nonsingular and let  $\mathcal{M}$  be a complete matching for it. Using  $\mathcal{M}$ , graph  $\mathcal{G}_F$  can be directed as follows. Edge  $(f_i, x_j)$  is directed from  $f_i$  to  $x_j$  if  $(f_i, x_j) \in \mathcal{M}$ , from  $x_j$  to  $f_i$  otherwise. Denote by  $\vec{\mathcal{G}}_F^{\mathcal{M}}$  the resulting directed graph. The following result holds ([25], Chapter 6.10):

**Lemma 12.** *The strongly connected components of  $\vec{\mathcal{G}}_F^{\mathcal{M}}$  are independent of  $\mathcal{M}$ .*

Each strongly connected component defines a *block* of  $F$ , consisting of the set of all equations that are vertices of this component. Blocks are partially ordered and we denote by  $\preceq_F$  this order. Extending  $\preceq_F$  to a total order (via topological sorting) yields an ordering of equations that puts the Jacobian matrix  $\mathbf{J}_X F$  in *Block Triangular Form* (BTF).

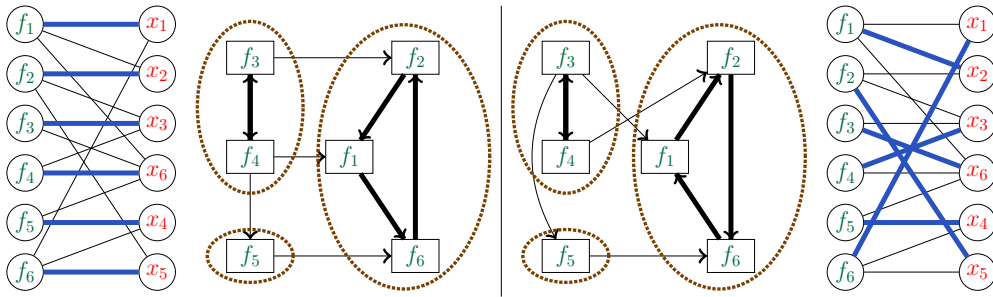


Figure 10: Illustrating Lemma 12 (we only show the projection of strongly connected components on function nodes).

This lemma is illustrated in Figure 10, inspired by [lecture notes by J-Y. L'Excellent and Bora Uçar, 2010](#). In this figure, a same bipartite graph  $\mathcal{G}$  is shown twice (left and right), with the equations sitting on the left-hand side in green, and the variables on the right-hand side in red. Two complete matchings  $\mathcal{M}_1$  (left) and  $\mathcal{M}_2$  (right) are shown in thick blue, with other edges of the bipartite graph being black. The restriction, to the equation vertices, of the two directed graphs  $\vec{\mathcal{G}}_F^{\mathcal{M}_1}$  (left) and  $\vec{\mathcal{G}}_F^{\mathcal{M}_2}$  (right) are shown on both sides. Although the two directed graphs  $\vec{\mathcal{G}}_F^{\mathcal{M}_1}$  and  $\vec{\mathcal{G}}_F^{\mathcal{M}_2}$  differ, the resulting block structures (encircled in dashed brown) are identical.

*Nonsquare systems.* In our development, we will also encounter nonsquare systems of algebraic equations. For general systems (with a number of variables not equal to the number of equations,  $n \neq m$  in (49)), call *block* a pair  $\beta = (\mathbf{f}, \mathbf{x})$  collecting a subset  $\mathbf{f}$  of the set of  $f_i$ 's and a subset  $\mathbf{x}$  of variables  $x_j$  involved in it.

**Definition 13 (Dulmage-Mendelsohn).** *For  $F=0$  a general system of algebraic equations, the Dulmage-Mendelsohn decomposition [46] of  $\mathcal{G}_F$  yields the partition of system  $F=0$  into the following three blocks, given some matching of maximal cardinality for  $\mathcal{G}_F$ :*

- block  $\beta_{\mathcal{D}}$  collects the variables and equations reachable via some alternating path from some unmatched equation;

- block  $\beta_{\mathcal{U}}$  collects the variables and equations reachable via some alternating path from some unmatched variable;
- block  $\beta_{\mathcal{E}}$  collects the remaining variables and equations.

Blocks  $\beta_{\mathcal{D}}, \beta_{\mathcal{U}}, \beta_{\mathcal{E}}$  are the  $\mathcal{D}$ overdetermined,  $\mathcal{U}$ nderdetermined, and  $\mathcal{E}$ nabled parts of  $F$ .

Statement 1 of the following lemma ensures that Definition 13 is meaningful:

**Lemma 14.**

1. The triple  $(\beta_{\mathcal{D}}, \beta_{\mathcal{E}}, \beta_{\mathcal{U}})$  defined by the Dulmage-Mendelsohn decomposition does not depend on the particular maximal matching used for its definition. We thus write  $(\beta_{\mathcal{U}}, \beta_{\mathcal{E}}, \beta_{\mathcal{D}}) = \text{DM}(F)$ .
2. System  $F$  is structurally nonsingular if and only if the overdetermined and underdetermined blocks  $\beta_{\mathcal{D}}$  and  $\beta_{\mathcal{U}}$  are both empty.

A refined DM decomposition exists, which in addition refines block  $\beta_{\mathcal{E}}$  into its indecomposable block triangular form. This is a direct consequence of applying Lemma 12 to  $\beta_{\mathcal{E}}$  in order to get its BTF.

The following obvious lemma will be useful in the sequel.

**Lemma 15.** Let  $\text{DM}(F) = (\beta_{\mathcal{U}}, \beta_{\mathcal{E}}, \beta_{\mathcal{D}})$  be the Dulmage-Mendelsohn decomposition of  $F=0$ . Then, no overdetermined block exists in the Dulmage-Mendelsohn decomposition of  $F \setminus \beta_{\mathcal{D}}$ .

**Comment 16.** Suppose that, instead of removing from  $F$  all equations belonging to  $\beta_{\mathcal{D}}$ , we only remove all unmatched equations with reference to the matching of maximal cardinality used for generating  $\text{DM}(F)$ , and call  $F'$  the remaining subsystem of  $F$ . Then, we still get that no overdetermined block exists in the Dulmage-Mendelsohn decomposition of  $F'$  and we have  $F' \supseteq (F \setminus \beta_{\mathcal{D}})$ , the inclusion being strict in general. However, this policy depends on the particular matching. In contrast, our policy is canonical since the DM decomposition is independent from the matching. Since we will rely on Lemma 15 for our compilation scheme in Section 6, adopting the finer policy would introduce a risk of nondeterminism due to the arbitrary choice of the particular complete matching selected. This is the reason for sticking with the policy stated in Lemma 15.  $\square$

### 5.1.3. Local versus non-local use of structural analysis

Structural analysis relies on the Implicit Function Theorem for its justification, as we have seen. Therefore, the classical use of structural analysis is local: let  $F(Y, X)=0$  be a structurally nonsingular system of equations with dependent variables  $X$ , and let  $\nu_X$  satisfy  $F(\nu_Y, \nu_X)=0$  for given values  $\nu_Y$  for  $Y$ ; then, the system remains (generically) nonsingular if  $\nu$  varies in a neighborhood of  $\nu_Y$ . This is the situation encountered in the running of ODE or DAE solvers, since a close guess of the solution is known from the previous time step.

We are, however, also interested in invoking the structural nonsingularity of system  $F(Y, X) = 0$  when no close guess is known. This is the situation encountered when handling mode changes, see our examples of Sections 2 and 3. We thus need another argument to justify the use of structural analysis in this case.

As a prerequisite we recall basic definitions regarding smooth manifolds (see, e.g., [22], Section 4.7). In the next lemma, we consider the system  $F(Y, X)=0$  for a fixed value of  $Y$ , and thus we omit  $Y$ .

**Lemma 17.** *Let  $k, m, n \in \mathbb{N}$  be such that  $m < n$  and set  $p = n - m$ . For  $\mathcal{S} \subset \mathbb{R}^n$  and  $x^* = (x_1^*, \dots, x_n^*) \in \mathcal{S}$ , the following properties are equivalent:*

1. *There exists an open neighborhood  $V$  of  $x^*$  and a  $\mathcal{C}^k$ -diffeomorphism  $F : V \rightarrow W \subset \mathbb{R}^n$  such that  $F(x^*) = 0$  and  $F(\mathcal{S} \cap V)$  is the intersection of  $W$  with the subspace of  $\mathbb{R}^n$  defined by the  $m$  equations  $w_{p+1}=0, \dots, w_n=0$ , where  $w_i$  denote the coordinates of points belonging to  $W$ .*
2. *There exists an open neighborhood  $V$  of  $x^*$ , an open neighborhood  $U$  of 0 in  $\mathbb{R}^p$  and a homeomorphism  $\psi : U \rightarrow \mathcal{S} \cap V$ , such that  $\psi$ , seen as a map with values in  $\mathbb{R}^n$ , is of class  $\mathcal{C}^k$  and has rank  $p$  at 0—i.e.,  $\psi'(0)$  has rank  $p$ .*

Statement 1 means that  $\mathcal{S}$  is, in a neighborhood of  $x^*$ , the solution set of the system of equations  $f_1(X)=0, \dots, f_m(X)=0$  in the  $n$ -tuple of dependent variables  $X$ , where  $F = (f_1, \dots, f_m)$ . Statement 2 expresses that  $\psi$  is a parameterization of  $\mathcal{S}$  in a neighborhood of  $x^*$ , of class  $\mathcal{C}^k$  and rank  $p$ —meaning that  $\mathcal{S}$  has dimension  $p$  in a neighborhood of  $x^*$ .

In order to apply Lemma 17 to the non-local use of structural analysis, we have to study the case of square systems of equations (i.e., let  $m = n$ ). Consider a system  $F(X)=0$ , where  $X=(x_1, \dots, x_n)$  is the  $n$ -tuple of dependent variables and  $F=(f_1, \dots, f_n)$  is an  $n$ -tuple of functions  $\mathbb{R}^n \rightarrow \mathbb{R}$  of class  $\mathcal{C}^k$ . Let  $\mathcal{S} \subseteq \mathbb{R}^n$  be the solution set of this system, that we assume is non-empty. To be able to apply Lemma 17, we augment  $X$  with one extra variable  $z$ , i.e., we set  $Z =_{\text{def}} X \cup \{z\}$  and we now regard our formerly square system  $F(X)=0$  as an augmented system  $G(Z)=0$  where  $G(X, z) =_{\text{def}} F(X)$ . Extended system  $G$  possesses  $n$  equations and  $n+1$  dependent variables, hence,  $p=1$ , and its solution set is equal to  $\mathcal{S} \times \mathbb{R}$ . Let  $x^* \in \mathcal{S}$  be a solution of  $F = 0$ . Then, we have  $G(x^*, z^*)=0$  for any  $z^* \in \mathbb{R}$ . Assume further that the Jacobian matrix  $\mathbf{J}_X F(x)$  is nonsingular at  $x^*$ . Then, it remains nonsingular in a neighborhood of  $x^*$ . Hence, the Jacobian matrix  $\mathbf{J}_X G(x, z)$  has rank  $n$  in a neighborhood  $V$  of  $(x^*, z^*)$  in  $\mathbb{R}^{n+1}$ . We can thus extend  $G$  by adding one more function in such a way that the resulting  $(n+1)$ -tuple  $\bar{G}$  is a  $\mathcal{C}^k$ -diffeomorphism, from  $V$  to an open set  $W$  of  $\mathbb{R}^{n+1}$ . The extended system  $\bar{G}=0$  satisfies Property 1 of Lemma 17. By Property 2 of Lemma 17,  $(\mathcal{S} \times \mathbb{R}) \cap V$  has dimension 1, implying that  $\mathcal{S} \cap \text{proj}_{\mathbb{R}^n}(V)$  has dimension  $< 1$ . This analysis is summarized in the following lemma:

**Lemma 18.** *Consider the square system  $F(X)=0$ , where  $X=(x_1, \dots, x_n)$  is the  $n$ -tuple of dependent variables and  $F=(f_1, \dots, f_n)$  is an  $n$ -tuple of functions  $\mathbb{R}^n \rightarrow \mathbb{R}$  of class  $\mathcal{C}^k, k \geq 1$ . Let  $\mathcal{S} \subseteq \mathbb{R}^n$  be the solution set of this system. Assume that system  $F=0$  has solutions and let  $x^*$  be such a solution. Assume further that the Jacobian matrix  $\mathbf{J}_X F(x)$  is nonsingular at  $x^*$ . Then, there exists an open neighborhood  $V$  of  $x^*$  in  $\mathbb{R}^n$ , such that  $\mathcal{S} \cap V$  has dimension  $< 1$ .*

The implications for the non-local use of structural analysis are as follows. Consider the square system  $F(X)=0$ , where  $X=(x_1, \dots, x_n)$  is the  $n$ -tuple of dependent variables and  $F=(f_1, \dots, f_n)$  is an  $n$ -tuple of functions  $\mathbb{R}^n \rightarrow \mathbb{R}$  of class  $\mathcal{C}^k$ . We first check the structural nonsingularity of  $F=0$ . If structural nonsingularity does not hold, then we abort solving the system. Otherwise, we proceed to solving the system and the following cases can occur:

1. the system possesses no solution;
2. its solution set is nonempty, of dimension  $< 1$ ; or

3. its solution set is nonempty, of dimension  $\geq 1$ .

Based on Lemma 18, one among cases 1 or 2 will hold, generically, whereas case 3 will hold exceptionally. Thus, structurally, we know that the system is not underdetermined. This is our justification of the non-local use of structural analysis. Note that the existence of a solution is not guaranteed. Furthermore, unlike for the local use, uniqueness is not guaranteed either.<sup>12</sup> Of course, subclasses of systems for which existence and uniqueness are guaranteed are of interest—the illustration examples we develop in this paper have this property.

#### 5.1.4. Equations with existential quantifiers

To our knowledge, the results of this section are new. For the handling of mode change events, we will also need to develop a structural analysis for the following class of (possibly nonsquare) algebraic systems of equations with existential quantifier:

$$\exists W : F(X, W, Y)=0 , \quad (52)$$

where  $(X, W)$  collects the dependent variables of system  $F(X, W, Y)=0$ . Our aim is to find structural conditions ensuring that (52) defines a partial function  $Y \rightarrow X$ , meaning that the value of tuple  $X$  is uniquely defined by the satisfaction of (52), given a value for tuple  $Y$ .

At this point, a fundamental difficulty arises. Whereas (52) is well defined as an abstract relation, it cannot in general be represented by a projected system of smooth algebraic equations of the form  $G(X, Y) = 0$ . Such a  $G$  can be associated to (52) only in subclasses of systems.<sup>13</sup> In general, no extension of the Implicit Function Theorem exists for systems of the form (52); thus, one cannot apply as such the arguments developed in Section 5.1.3. Therefore, we will reformulate our requirement differently.

Say that  $\nu_Y$  is *consistent* if the system of equations  $F(X, W, \nu_Y)=0$  possesses a solution for  $(X, W)$ . We thus consider the following property for the system  $F(X, W, Y)=0$ : for every consistent  $\nu_Y$ ,

$$\left. \begin{array}{l} F(\nu_X^1, \nu_W^1, \nu_Y) = 0 \\ F(\nu_X^2, \nu_W^2, \nu_Y) = 0 \end{array} \right\} \implies \nu_X^1 = \nu_X^2 , \quad (53)$$

expressing that  $X$  is independent of  $W$  given a consistent tuple of values for  $Y$ . To find structural criteria guaranteeing (53), we will consider the algebraic system  $F(X, W, Y)=0$  with  $X, W, Y$  as dependent variables (i.e., values for the entries of  $Y$  are no longer seen as being given).

**Definition 19.** *System (52) is structurally nonsingular if the following two conditions hold, almost everywhere when the nonzero coefficients of the Jacobian matrix  $\mathbf{J}_{X,W,Y}F$  vary over some neighborhood:<sup>14</sup> consistent values for  $Y$  exist and condition (53) holds.*

<sup>12</sup>Lemma 18 only states that there cannot be two arbitrarily close solutions of  $F = 0$ , as its conclusion  $\mathcal{S} \cap V = \{x^*\}$  involves an open neighborhood  $V$  of  $x^*$ .

<sup>13</sup>Examples are linear systems for which elimination is easy, and polynomial systems for which it is doable—but expensive—using Gröbner bases.

<sup>14</sup>The precise meaning of this statement is the same as in Lemma 11, see the discussion thereafter.



The structural nonsingularity of (52) can be checked by using the DM decomposition of  $F$  as follows. Let  $(\beta_{\mathcal{U}}, \beta_{\mathcal{E}}, \beta_{\mathcal{D}}) = \text{DM}(F)$  be the DM decomposition of  $F(X, W, Y)=0$ , with  $(X, W, Y)$  as dependent variables. We further assume that the regular block  $\beta_{\mathcal{E}}$  is expressed in its block triangular form, see Lemma 12 and comments thereafter. Let  $\mathbf{B}$  be the set of all indecomposable blocks of  $\beta_{\mathcal{E}}$  and  $\preceq$  be the partial order on  $\mathbf{B}$  following Lemma 12. The following holds:

**Lemma 20.** *System (52) is structurally nonsingular if and only if:*

1. Block  $\beta_{\mathcal{D}}$  is empty;
2. Block  $\beta_{\mathcal{U}}$  involves no variable belonging to  $X$ ;
3. For every  $\beta \in \mathbf{B}$  containing some variable from  $X$ , then,  $\beta$  contains no variable from  $W$ , and for every  $\beta' \prec \beta$  and every directed edge  $(z', z) \in \vec{\mathcal{G}}_F^{\mathcal{M}}$  such that  $z' \in \beta'$ ,  $z \in \beta$ , and  $\mathcal{M}$  is an arbitrary complete matching for  $\mathcal{G}_F$ , then  $z' \notin W$ .

Conditions 1 and 2 speak by themselves. Regarding Condition 3, the intuition is the following. Every directed path of  $\vec{\mathcal{G}}_F^{\mathcal{M}}$ , originating from  $W$  and terminating in  $X$ , of minimal length, must traverse  $Y$ . Consequently,  $W$  influences  $X$  “through”  $Y$  only.

*Proof.* “If” part. By Condition 1, there exist consistent values for  $Y$ . By condition 2, no variable of  $X$  belongs to the underdetermined part of  $\text{DM}(F)$ . It remains to show that condition (53) holds. By condition (3), each indecomposable block  $\beta \in \mathbf{B}$  involving a variable of  $X$  has the form  $G(Z, U)=0$ , where

1.  $Z$  is the  $n$ -tuple of dependent variables of  $G = (g_1, \dots, g_n)$ ,
2. no variable of  $W$  belongs to  $Z$ , and
3.  $U \subseteq X \cup Y$  is a subset of the dependent variables of the blocks  $\beta'$  immediately preceding  $\beta$ .

Since  $G=0$  is structurally regular, fixing a value for  $U$  entirely determines all the dependent variables of block  $\beta$ . This proves the “if” part.

“Only if” part. We prove it by contradiction. If condition 1 does not hold, then the existence of consistent values for  $Y$  is not structurally guaranteed. If condition 2 does not hold, then the variables of  $X$  that belong to  $\beta_{\mathcal{U}}$  are not determined, even for given values of  $W, Y$ . If condition 3 does not hold, then two cases can occur:

- Some  $x \in X$  and  $w \in W$  are involved in a same indecomposable block  $\beta \in \mathbf{B}$ . Then, condition (53) will not hold for  $x$  due to the structural coupling with  $w$  through block  $\beta$ .
- There exists an indecomposable block  $\beta \in \mathbf{B}$  of the form  $G(Z, U)=0$ , where
  1.  $Z$  is the  $n$ -tuple of dependent variables of  $G = (g_1, \dots, g_n)$ , and  $Z$  contains a variable  $x \in X$ ,
  2. no variable of  $W$  belongs to  $Z$ , and
  3. some variable  $w \in W$  belongs to  $U$ , the set of the dependent variables of the blocks  $\beta'$  immediately preceding  $\beta$ .

Hence,  $w$  influences  $x$ , structurally, which again violates (53). This finishes the proof.  $\square$

We complement Lemma 20 with the algorithm `ExistQuantifEqn` (Algorithm 5), which requires a system  $F$  of the form (52). If condition 1 of Lemma 20 fails to be satisfied, then  $b_{\mathcal{D}} \leftarrow \text{F}$  is returned, indicating overdetermination. If conditions 2 or 3 of Lemma 20 fail to be satisfied, then  $b_{\mathcal{U}} \leftarrow \text{F}$  is returned, indicating underdetermination. Otherwise, `ExistQuantifEqn` succeeds and returns the value  $\text{T}$  for both Booleans, together with the decomposition  $F_{\Sigma} \cup \overline{F}_{\Sigma}$  of  $\beta_{\mathcal{E}}$ . In this decomposition:

- Subsystem  $F_{\Sigma}$  collects the indecomposable blocks involving variables belonging to  $X$ , so that  $F_{\Sigma}$  determines  $X$  as a function of  $\nu_Y$  when  $\nu_Y$  is consistent;
- Subsystem  $\overline{F}_{\Sigma}$  collects the consistency conditions, whose dependent variables belong to  $W \cup Y$ .

---

**Algorithm 5** `ExistQuantifEqn`

---

**Require:**  $F$ ; **return**  $(b_{\mathcal{D}}, b_{\mathcal{U}}, F_{\Sigma}, \overline{F}_{\Sigma})$

```

1:  $(\beta_{\mathcal{U}}, \beta_{\mathcal{E}}, \beta_{\mathcal{D}}) = \text{DM}(F)$ 
2: if condition 1 of Lemma 20 holds then
3:    $b_{\mathcal{D}} \leftarrow \text{T}$ ;
4:   if 2 and 3 of Lemma 20 hold then
5:      $b_{\mathcal{U}} \leftarrow \text{T}$ ;
6:     partition  $\beta_{\mathcal{E}} = F_{\Sigma} \cup \overline{F}_{\Sigma}$ 
7:   else
8:      $b_{\mathcal{U}} \leftarrow \text{F}$ 
9:   else
10:   $b_{\mathcal{D}} \leftarrow \text{F}$ 

```

---

Our background on the structural analysis of algebraic equations is now complete. In the next section, we recall the background on the structural analysis of (single-mode) DAE systems.

### 5.2. The $\Sigma$ -method for DAE systems

The DAE systems we consider are “square”, i.e., have the form

$$f_j(\text{the } x_i\text{'s and their derivatives}) = 0 \quad (54)$$

where  $x_1, \dots, x_n$  denote the dependent variables and  $f_1 = 0, \dots, f_n = 0$  denote the equations. Call *leading variables* of System (54) the  $d_i$ -th derivatives<sup>15</sup>  $x_i^{/d_i}$  for  $i = 1, \dots, n$ , where  $d_i$  is the maximal differentiation degree of variable  $x_i$  throughout  $f_1 = 0, \dots, f_n = 0$ . The problem addressed by the structural analysis of DAE systems of the form (54) is the following. Regard (54) as a system of algebraic equations with the leading variables as unknowns. If this system is structurally nonsingular, then, given a value for all the

---

<sup>15</sup>The notation  $x^{/k}$  is adopted throughout this paper, instead of the more classical  $x^{(k)}$ , for the  $k$ -th derivative of  $x$ . See Notations 24.

$x_i^k$  for  $i = 1, \dots, n$  and  $k = 0, \dots, d_i - 1$ , a unique value for the leading variables can be computed, structurally; hence, System (54) is like an ODE. If this is not the case, finding additional *latent equations* by differentiating suitably selected equations from (54) will bring the system to an ODE-like form, while not changing its set of solutions. Performing this is known as *index reduction*.

In our simple examples, finding latent equations was easy. In general, this is difficult and algorithms were proposed in the literature for doing it efficiently. Among them, the Pantelides algorithm [40] is the historical solution. We decided, however, to base our subsequent developments on the beautiful method proposed in [48] by J. Pryce, called the  $\Sigma$ -method. The  $\Sigma$ -method also covers the construction of the Block Triangular Form and addresses numerical issues, which we do not discuss here.

*Weighted bipartite graphs.* We consider System (54), which is entirely characterized by its set of dependent variables  $X$  (whose generic element is denoted by  $x$ ) and its set of equations  $F=0$  (whose generic element is written  $f=0$ ). We attach to (54) the bipartite graph  $\mathcal{G}$  having  $F \cup X$  as set of vertices and having an edge  $(f, x)$  if and only if  $x$  occurs in function  $f$ , regardless of its differentiation degree. Recall that a matching is *complete* iff it involves all equations of  $F$ .

So far,  $\mathcal{G}$  is agnostic with respect to differentiations. To account for this, we further equip  $\mathcal{G}$  with *weights*: to each edge  $(f, x) \in \mathcal{G}$  is associated a nonnegative integer  $d_{fx}$ , equal to the maximal differentiation degree of variable  $x$  in function  $f$ . This yields a *weight* for any matching  $\mathcal{M}$  of  $\mathcal{G}$  by the formula  $w(\mathcal{M}) = \sum_{(f,x) \in \mathcal{M}} d_{fx}$ . Suppose we have a solution to the following problem:

**Problem 21.** Find a complete matching  $\mathcal{M}$  for  $\mathcal{G}$  and integer offsets  $\{c_f \mid f \in F\}$  and  $\{d_x \mid x \in X\}$ , satisfying the following conditions, where the  $d_{fx}$  are the weights as before:

$$\begin{aligned} d_x - c_f &\geq d_{fx} \quad \text{with equality if } (f, x) \in \mathcal{M} \\ c_f &\geq 0 \end{aligned} \tag{55}$$

Then, differentiating  $c_f$  times each function  $f$  yields a DAE system  $F_\Sigma=0$  having the following properties. Inequality  $d_x \geq c_f + d_{fx}$  holds for each  $x \in X$ , and  $d_x = c_f + d_{fx}$  holds for the unique  $f$  such that  $(f, x) \in \mathcal{M}$ . Hence, the leading variables of DAE system  $F_\Sigma=0$  are the  $d_x$ -th derivatives  $x^{d_x}$ . Consequently, system  $F_\Sigma=0$ , now seen as a system of algebraic equations having  $x^{d_x}$  as dependent variables, is structurally nonsingular by Definition 10. Hence,  $F_\Sigma=0$  is like an ODE. The integer  $k = \max_{f \in F} c_f$  is called the *index* of the system.<sup>16</sup>

**Definition 22.** For  $F$  a DAE system, the solution to Problem 21 yields the DAE system  $F_\Sigma = \{f^{c_f} \mid f \in F\}$  together with the system of consistency constraints  $\overline{F}_\Sigma = \{f^{k'} \mid f \in F, 0 \leq k' < c_f\}$ .

Knowing the offsets also allows transforming  $F$  into a system of index 1, by not performing the final round of differentiations. There are infinitely many solutions to Problem 21 with unknowns  $c_f$  and  $d_x$ , since, for example, adding the same  $\ell \in \mathbb{N}_{\geq 0}$  to all  $c_f$ 's and  $d_x$ 's yields another solution. We thus seek for a *smallest* solution, elementwise. Hence,

<sup>16</sup>We should rather say the *differentiation index* as, once again, other notions of index exist for DAEs [21] that are not relevant to our work.

Problem 21 is the fundamental problem we must solve, and we seek a smallest solution for it.

The beautiful idea of J. Pryce is to propose a linear program encoding Problem 21. As a preliminary step, we claim that a bruteforce encoding of Problem 21 is the following: Find integers  $\xi_{fx}, d_x, c_f$  such that

$$\begin{aligned}
& \left. \begin{array}{l} \sum_{f:(f,x) \in \mathcal{G}} \xi_{fx} = 1 \\ \sum_{x:(f,x) \in \mathcal{G}} \xi_{fx} = 1 \\ \xi_{fx} \geq 0 \end{array} \right\} \text{complete matching} \\
& \left. \begin{array}{l} d_x - c_f - d_{fx} \geq 0 \\ c_f \geq 0 \end{array} \right\} \text{encodes “}\geq\text{” in (55)} \\
& \sum_{(f,x) \in \mathcal{G}} \xi_{fx} (d_x - c_f - d_{fx}) = 0 \text{ encodes “}=\text{” in (55)}
\end{aligned} \tag{56}$$

We now justify our claim. Focus on the first block. Since the  $\xi$ 's are integers, they can only take values in  $\{0, 1\}$  and one can define a bipartite graph by stating that edge  $(f, x)$  exists iff  $\xi_{fx} = 1$ . The first two equations formalize that this graph is a matching, which is complete since all equations are involved. The second block is a direct encoding of (55) if we ignore the additional statement “with equality iff”. The latter is encoded by the last constraint (since having the sum equal to zero requires that all the terms be equal to zero). Constraint problem (56) does not account for our wish for a “smallest” solution: this will be handled separately.

Following the characterization of solutions of linear programs via *complementary slackness conditions*, every solution of problem (56) is a solution of the following dual linear programs (LP), where the  $\xi_{fx}$  are real:

$$\begin{aligned}
\text{primal:} \quad & \begin{array}{l} \text{maximize} \quad \sum_{(f,x) \in \mathcal{G}} d_{fx} \xi_{fx} \\ \text{subject to} \quad \sum_{f:(f,x) \in \mathcal{G}} \xi_{fx} = 1 \\ \text{and} \quad \sum_{x:(f,x) \in \mathcal{G}} \xi_{fx} \geq 1 \\ \text{and} \quad \xi_{fx} \geq 0 \end{array}
\end{aligned} \tag{57}$$

$$\begin{aligned}
\text{dual:} \quad & \begin{array}{l} \text{minimize} \quad \sum_x d_x - \sum_f c_f \\ \text{subject to} \quad d_x - c_f \geq d_{fx} \\ \text{and} \quad c_f \geq 0 \end{array}
\end{aligned} \tag{58}$$

In these two problems,  $f$  ranges over  $F$ ,  $x$  ranges over  $X$ , and  $(f, x)$  ranges over  $\mathcal{G}$ . Also, we have relaxed the integer LP to a real LP, as the two possess identical solutions in this case. Note that LP (57) encodes the search for a complete matching of maximum weight for  $\mathcal{G}$ . By the principle of complementary slackness in linear programming,

$$\begin{aligned}
& \text{for respective optima of problems (57) and (58),} \\
& \xi_{fx} > 0 \text{ if and only if } d_x - c_f = d_{fx},
\end{aligned} \tag{59}$$

which is exactly the last constraint of (56). Using this translation into linear programs (57) and (58), it is proved in [48], Thm 3.6, that, if a solution exists to Problem 21, then a unique elementwise smallest solution exists. Based on the above analysis, the following algorithm was proposed in [48] for solving (57,58) and was proved to provide the smallest solution for (58). Let  $\mathcal{G}$  be a bipartite graph with weights  $\{d_{fx} | (f, x) \in \mathcal{G}\}$ :

1. Solve LP (57), which gives a complete matching  $\mathcal{M}$  of maximum weight for  $\mathcal{G}$ ; any method for solving LP can be used;
2. Apply the following iteration until a fixpoint is reached (in finitely many steps), from the initial values  $c_f = 0$ :
  - (a)  $\forall x : d_x \leftarrow \max\{d_{fx} + c_f \mid (f, x) \in \mathcal{G}\}$  ;
  - (b)  $\forall f : c_f \leftarrow d_x - d_{fx}$  where  $(f, x) \in \mathcal{M}$  .

The reason for using the special iterative algorithm for solving the dual LP (58) is that a standard LP-solving algorithm will return an arbitrary solution, not necessarily the smallest one. The following lemma, that will be crucially put to use in Section 9.1, is an obvious consequence of the linearity of problems (57) and (58):

**Lemma 23.** *Let  $\mathcal{G}$  be a given bipartite graph and let two families of weights  $(d_{fx}^1)_{(f,x) \in \mathcal{G}}$  and  $(d_{fx}^2)_{(f,x) \in \mathcal{G}}$  be related by  $d_{fx}^2 = M \times d_{fx}^1$  for every  $(f, x) \in \mathcal{G}$ , where  $M$  is a fixed positive integer. Then, the offsets of the corresponding  $\Sigma$ -method are also related in the same way:  $d_x^2 = M \times d_x^1$  for every variable  $x$ , and  $c_f^2 = M \times c_f^1$  for every function  $f$ .*

A necessary and sufficient condition for Problem 21 to have a solution is that the set of complete matchings for  $\mathcal{G}$  is non-empty. We thus consider the function

$$b \leftarrow \text{ExistsMatching}(\mathcal{G}), \quad (60)$$

which returns a Boolean  $b$  such that  $b = \top$  iff there exists a complete matching for the bipartite graph  $\mathcal{G}$  (regardless of its associated weights). In this case, step 1 of the above algorithm is guaranteed successful (since the set of complete matchings for  $\mathcal{G}$  is nonempty and finite). In the sequel, we encode the above algorithm as the function

$$(\mathbf{c}, \mathbf{d}) \leftarrow \text{FindOffsets}(\mathcal{G}, \mathbf{D}), \quad \text{where} \quad \begin{cases} \mathbf{c} : & F \rightarrow \mathbb{N} \\ \mathbf{d} : & X \rightarrow \mathbb{N} \\ \mathbf{D} : & F \times X \rightarrow \mathbb{N} \end{cases} \quad (61)$$

Having  $b = \top$  returned by  $\text{ExistsMatching}()$  is a prerequisite to calling  $\text{FindOffsets}()$ , which is then guaranteed successful.

---

#### Algorithm 6 SigmaMethod

---

**Require:**  $F$ ; **return**  $(b, \overline{F}_\Sigma, F_\Sigma)$

- 1:  $b \leftarrow \text{ExistsMatching}(\mathcal{G})$ .
  - 2: **if**  $b$  **then**
  - 3:      $(\mathbf{c}, \mathbf{d}) \leftarrow \text{FindOffsets}(\mathcal{G}, \mathbf{D})$
  - 4:      $(F_\Sigma, \overline{F}_\Sigma) \leftarrow \text{LatentEqns}(\mathbf{c}, F)$
- 

We encapsulate these two functions as Algorithm 6 (SigmaMethod). It requires a DAE system  $F=0$ .  $\mathcal{G}$  is the weighted bipartite graph of  $F$  with weights  $\mathbf{D}$ . SigmaMethod returns a Boolean  $b$  and, if  $b = \top$ , a pair  $(\overline{F}_\Sigma, F_\Sigma)$  defining an index-0 system  $F_\Sigma=0$  and a (possibly underdetermined) system  $\overline{F}_\Sigma=0$  of consistency constraints following Definition 22. In Line 4,  $F_\Sigma$  collects the  $f^{c_f}$ , for  $f \in F$ , and  $\overline{F}_\Sigma$  collects the  $f'^k$ , for  $f \in F$  and  $0 \leq k < c_f$ . This completes our background material on structural analysis.

## 6. Structural analysis of multimode DAE systems

In this section, we extend structural analysis, from (single-mode) DAE systems, to multimode DAE systems. As a prerequisite, we precisely define the class of multimode DAE systems considered. Since the structural analysis will actually be applied to the nonstandard semantics, we also need to define multimode *difference Algebraic Systems* (dAE), operating in discrete time.

### 6.1. Defining multimode DAE/dAE systems

In this section, the class of systems of multimode Differential/difference Algebraic Equations considered in this paper is formally defined. By analogy with DAE's, we use the acronym dAE to mean *difference algebraic equations*. Notation  $x^\bullet$  denotes the shift of  $x$ , as per Definition 2.

**Notations 24.** Let  $\mathcal{X}$  be an underlying set of variables. For  $x \in \mathcal{X}$  and  $m \in \mathbb{N}$ , the  $m$ -differentiation and  $m$ -shift of  $x$  are denoted by  $x'^m$  and  $x^{\bullet m}$ , respectively. Let  $\mathcal{X}'^m$  and  $\mathcal{X}^{\bullet m}$  denote the set of all  $x'^m$  and  $x^{\bullet m}$ , for  $x$  ranging over the set  $\mathcal{X}$  of variables. Let

$$\mathcal{X}^{(\prime)} =_{\text{def}} \bigcup_{m \in \mathbb{N}} \mathcal{X}'^m \text{ and } \mathcal{X}^{(\bullet)} =_{\text{def}} \bigcup_{m \in \mathbb{N}} \mathcal{X}^{\bullet m} . \quad (62)$$

For  $X \subset \mathcal{X}$ , sets  $X'^m, X^{\bullet m}, X^{(\prime)}, X^{(\bullet)}$  are defined in a similar way. □

The following assumption will be in force in the remainder of this paper:

**Assumption 1.** *All the functions and predicates over  $\mathcal{X}^{(\prime)}$  are assumed to involve only finitely many of the  $x'^m$  for  $x \in \mathcal{X}$  and  $m \in \mathbb{N}$ . The same holds regarding  $\mathcal{X}^{(\bullet)}$ .*

#### 6.1.1. Syntax and meaning

From the discussion in Section 3.6 regarding the Cup-and-Ball example, we know we need to type the modes as either *long* or *transient*, since the two situations require different index reduction techniques. This justifies the consideration of two different kinds of guarded equations.

**Definition 25 (mDAE/mdAE syntax).** Multimode DAE systems (mDAEs), *respectively* multimode dAE systems (mdAEs), *are defined by the following syntax:*

$$\begin{aligned} S &::= s \mid S; s \\ s &::= g \mid e \\ g &::= \gamma = b \end{aligned} \quad (63)$$

$$e ::= \text{if } \gamma \text{ then } f=0 \quad (64)$$

$$\quad \mid \text{when } \gamma \text{ then } f=0 \quad (65)$$

mDAE  $S$  is a finite set of  $s$ ;  $s$  is either a *guard evaluation*  $g$  (63), or a *guarded equation*  $e$ , of two different kinds (64) or (65);  $f$  is a smooth scalar function over  $\mathcal{X}^{(\prime)}$ , and  $b$  is a Boolean expression of predicates over  $\mathcal{X}^{(\prime)}$ . The same holds for the syntax of mdAE, with  $\mathcal{X}^{(\bullet)}$  replacing  $\mathcal{X}^{(\prime)}$ .

**Notations 26.** Let  $e$  be a guarded equation. Its *body*  $f=0$  and *guard*  $\gamma$  are denoted by  $f(e)$  and  $\gamma(e)$ . For  $S$  an mDAE or mdAE system, we denote by  $X(S)$ ,  $E(S)$ , and  $\Gamma(S)$  its sets of numerical variables, guarded equations, and guards, respectively. The mention of  $S$  is omitted when it is clear from the context. Finally, for  $e : \text{if/when } \gamma \text{ then } f=0$  a guarded equation following Definition 25, we define

$$e^\bullet/e' : \text{if/when } \gamma \text{ then } f^\bullet/f'=0 \quad (66)$$

where  $f^\bullet$  follows Definition 2, and, for  $E$  a set of equations, Notations 24 are adapted so that  $E'^m, E^{\bullet m}, E^{(\prime)}, E^{(\bullet)}$  can be considered. We insist that, in (66), the guard  $\gamma$  is left unchanged.  $\square$

**Meaning.** The meanings of **if**-equation (64) and **when**-equation (65) differ:

- In the **if**-equation (64), equation  $f=0$  is enabled if and only if the guard  $\gamma$  holds. Otherwise, this equation is disabled.
- In the **when**-equation (65), equation  $f=0$  is enabled exactly at the events when guard  $\gamma$  switches from F to T. Otherwise, this equation is disabled.  $\square$

For mdAE systems, the formalization of the second statement is straightforward: equation  $f=0$  is enabled when

$$\gamma^\uparrow \text{ holds, where } \gamma^\uparrow =_{\text{def}} (\text{not } \bullet\gamma) \text{ and } \gamma . \quad (67)$$

For mDAE systems: equation  $f=0$  is enabled when

$$\gamma^\uparrow \text{ holds, where } \gamma^\uparrow =_{\text{def}} (\text{not } \gamma^-) \text{ and } \gamma \quad (68)$$

where  $\gamma^-$  denotes the left-limsup of  $\gamma$  at the current instant, defined by:

$$\begin{aligned} \gamma^-(t) &=_{\text{def}} \limsup_{s \nearrow t} \gamma(s) \\ &= \text{if } \exists s_n \nearrow t : \gamma(s_n)=\text{T} \text{ then T else F} \end{aligned}$$

The left-limsup is always defined and coincides with the left limit when the latter exists.

**Comment 27.** Why bother with the two types of equations, since **when**-equations are mapped to **if**-equations by mapping  $\gamma$  to  $\gamma^\uparrow$ ? The motivation was given in Section 3.6. We need to type the modes as either *long* or *transient*, since the two situations require different index reduction techniques. This is precisely what our **when**-equations guarantee, since, by construction, switching from F to T occurs in a single nonstandard instant. *Hence, by construction, when-equations are active at transient modes only.* Actually, this distinction between **if**-equations and **when**-equations already exists in object-oriented DAE-based languages such as Modelica.  $\square$

Regarding **if**-equations, we state the following assumption:

**Assumption 2.** *In if-equations, guard  $\gamma$  holds true for a (standard) positive duration.*

It is the responsibility of the programmer to make sure that **if**-equations are properly used so that Assumption 2 holds.

### 6.1.2. Long versus transient modes

For our theoretical development and the presentation of algorithms, it will be convenient to simplify our syntax by replacing **when**-equations by **if**-equations using (67) and Assumption 2. Our syntax of Definition 25 with arbitrary guards is thus replaced with the following (equivalent) syntax:

**Definition 28 (mDAE/mdAE revisited).** Multimode DAE systems (mDAEs), respectively multimode dAE systems (mdAEs), are defined by the following syntax:

$$\begin{aligned}
S &::= s \mid S; s \\
s &::= g \mid e \\
g &::= \gamma = b \\
e &::= \text{if } \gamma \text{ then } f=0
\end{aligned} \tag{69}$$

where  $b$  is a long/transient-typed Boolean expression of predicates over  $\mathcal{X}^{(l)}$  (respectively  $\mathcal{X}^{(\bullet)}$ ).

From now on, and unless otherwise specified, we will use the modified syntax of Definition 28. The notion of mode is then easily introduced.

**Definition 29 (mode).** For  $S$  a mDAE or mdAE, a mode is a valuation, in  $\{\text{F}, \text{T}\}$ , of all of its guards  $\gamma \in \Gamma(S)$ . The mode is called **transient** if at least one guard typed **transient** takes the value **T** in it; otherwise, it is called **long**. Modes will be generically denoted by the symbol  $\mu$ .

A mode enables a subset of the equations  $f=0$  and disables the other ones.

*Nonstandard semantics.* An mDAE  $S$  is transformed to a (nonstandard) mdAE  ${}^*S$  through the following syntactic transformation:

$${}^*S \stackrel{\text{def}}{=} S[x' \text{ is replaced with } \frac{x \bullet - x}{\partial}] \tag{70}$$

Since we perform the structural analysis on the nonstandard semantics, we must develop a structural analysis for mdAE systems, operating in discrete time. Therefore, in the rest of this section, we work in nonstandard semantics.  $\square$

**Comment 30. (Progressive evaluation of the modes)** It may be that some guards have their value unknown when the execution of the current instant starts; this was, for example, the case with the Cup-and-Ball example in its first form given in Section 1.2.2. We can use the known enabled equations to evaluate part of the variables, hoping that some of the yet unknown guards get their value defined. This yields more equations getting enabled/disabled, allowing for more guards being evaluated. If this process leads to the evaluation of all guards, then our execution scheme succeeds; otherwise, we report a failure due to a logico-numerical fixpoint equation, as was the case in Section 1.2.2.

This progressive evaluation of guards becomes problematic when both long and transient modes are considered. Indeed, different structural analyses are used for each case (see Section 3.3 for the long mode case of the Cup-and-Ball example, versus Section 3.5 for the transient mode case). Unfortunately, we do not know which one should be used if the mode is only partially known, since the long/transient qualification is only known when all guards have been evaluated. Consequently, we have to choose between the following two incompatible features for our compilation method:



1. either we support the progressive evaluation of modes, but are unable to handle transient modes;
2. or we support both long and transient modes, but we do not support the progressive evaluation of modes.

Since many physically meaningful models involve transient modes, our opinion is that the second option is much more useful in practice, hence it is adopted in our work.  $\square$

Comment 6, regarding the Cup-and-Ball example, discussed what we miss by favoring the second option. Indeed, delaying the effect of guards was the appropriate correction to remove logico-numerical fixpoint problems in the Cup-and-Ball example. Remember that the so introduced delay is infinitesimal, so this does not really matter.<sup>17</sup>

Therefore, throughout the rest of this work, Definition 28 is modified to restrict ourselves *by the syntax* to systems in which all guards can be evaluated at the initialization of each nonstandard instant:<sup>18</sup>

**Definition 31 (mdAE modified).** Multimode dAE systems (mdAEs), are defined by the following syntax:

$$\begin{aligned}
S &::= s \mid S; s \\
s &::= g \mid e \\
g &::= \gamma^\bullet = b \\
e &::= \text{if } \gamma \text{ then } f=0
\end{aligned} \tag{71}$$

where  $b$  is a long/transient-typed Boolean expression of predicates over  $\mathcal{X}^{(\bullet)}$ .

**Notations 32.** Referring to (71), we write  $b(\gamma)$  to refer to the expression defining  $\gamma^\bullet$ .  $\square$

From now on, by mdAE, we will mean a system of the form given in Definition 31.

### 6.2. Structural analysis: intuition

In this section, we give an intuition about how to extend the structural analysis from dAE systems to mdAE systems. To simplify our explanation, *we consider the restricted case in which both modes (previous and current) are long*. Also, it will be convenient to assume, for the two modes before and after the change, that we have two index-1 models, of the form

$$S_i : \begin{cases} X'_i &= F_i(X_i, Y_i) \\ 0 &= H_i(X_i) \end{cases}, \text{ for } i = 0, 1,$$

where indices 0 and 1 refer to the system before and after the change. Observe that system state  $X_i$  may collect different variables before and after change, and so does the algebraic tuple  $Y_i$ . Of course, the dynamics  $F_i, H_i$  are subject to changes as well. Let

---

<sup>17</sup>This is unlike Synchronous Languages for discrete time systems, in which the restriction resulting from adopting Definition 31 would be significant.

<sup>18</sup>With this new syntax, the original forms for the Cup-and-Ball and Westinghouse examples (given in the introduction) could not be expressed.

$\mathbf{J}_i =_{\text{def}} \mathbf{J}_{X_i} H_i$  denote the Jacobian of  $H_i$  w.r.t.  $X_i$ . Hence, the system augmented with its latent equations

$$S_{i\Sigma} : \begin{cases} X_i' &= F_i(X_i, Y_i) \\ 0 &= H_i(X_i) \\ 0 &= \mathbf{J}_i(X_i)X_i' \end{cases}, \text{ for } i = 0, 1,$$

seen as a system of algebraic equations with  $X_i', Y_i$  as dependent variables, is structurally nonsingular. Equivalently, we may replace, in this system, the differentiation of the algebraic constraint,  $\frac{d}{dt}H_i(X_i)$ , by its shifting  $H_i(X_i^\bullet)$  (their bipartite graphs are equal):

$$S_{i\Sigma} : \begin{cases} X_i' &= F_i(X_i, Y_i) \\ 0 &= H_i(X_i) \\ 0 &= H_i(X_i^\bullet) \end{cases}, \text{ for } i = 0, 1. \quad (72)$$

Form (72) can be seen as the result of applying Pantelides' or Pryce's structural analysis in each mode, before and after change.

**Comment 33 (key idea).** *Our key idea is to handle the mode change as a task of reconciling the possible conflict between the predictions generated by the previous mode and the consistency conditions associated to the new mode.* Indeed, at the instant of mode change, there is a possible conflict between:

- on the one hand, the value for  $X_0$  predicted by model  $S_{0\Sigma}$  which sets the value of the derivative  $X_0'$  just before the change;
- on the other hand, the consistency constraints on the pair  $X_1, Y_1$  resulting from model  $S_{1\Sigma}$  after the change.  $\square$

To perform this reconciliation, we move to the nonstandard semantics by *syntactically replacing*, in (72), the derivatives  $x'$  by their first-order Euler approximations with infinitesimal time step  $\partial$ :

$$x'(t) \quad \underbrace{\longleftarrow}_{\text{is replaced by}} \quad \frac{x(t + \partial) - x(t)}{\partial} = \frac{x^\bullet(t) - x(t)}{\partial}.$$

In addition, we replace, in (72), the derivative  $\frac{d}{dt}H_i(X_i) = \mathbf{J}_i(X_i)X_i'$  by the time shift  $H_i(X_i^\bullet)$ . This substitution is legitimate, since the bipartite graph of the system is left unchanged by performing this replacement.

Performing the above substitutions in (72) and sampling it on the nonstandard discrete time basis  $\mathbb{T} = \{0, \partial, 2\partial, \dots\}$  yields a discrete-time two-modes system. At the instant of mode change, we inherit the following system of algebraic equations:

$$S = \underbrace{\bullet S_{0\Sigma}}_{\text{previous instant}} \cup \underbrace{S_{1\Sigma}}_{\text{current instant}} \quad (73)$$

where

$$\bullet S_{0\Sigma} : \begin{cases} \frac{X_0 - \bullet X_0}{\partial} &= F_0(\bullet X_0, \bullet Y_0) \\ 0 &= H_0(\bullet X_0) \\ 0 &= H_0(X_0) \end{cases}, \quad S_{1\Sigma} : \begin{cases} \frac{X_1^\bullet - X_1}{\partial} &= F_1(X_1, Y_1) & (e_0) \\ 0 &= H_1(X_1) & (e_1) \\ 0 &= H_1(X_1^\bullet) & (e_1^\bullet) \end{cases} \quad (74)$$

and  $\cup$  refers to the union of sets of equations. For each state variable  $x$  that is shared between the two systems, i.e.,  $x \in X_0 \cap X_1$ , the value for  $x$  must be identical in  $\bullet S_{0\Sigma}$  and  $S_{1\Sigma}$ . To avoid duplicates of a same equation while taking the union of previous and current systems in  $S$ , we identify previous and current equations that are identical in their syntax, if any. The dependent variables of  $S$  are the  $X_0$  inherited from  $S_{0\Sigma}$  and the  $X_1, Y_1, X_1^\bullet$  inherited from  $S_{1\Sigma}$ . When seen as sets of variables,  $X_0$  on the one hand, and  $X_1, Y_1, X_1^\bullet$  on the other hand, possess in general a nonempty intersection. Being the result of a structural analysis, system  $S_{1\Sigma}$  is:

- structurally nonsingular, when seen as a system of algebraic equations in its leading variables  $Y_1, X_1^\bullet$ , for given values for  $X_1$ ; and
- possibly structurally underdetermined, when seen as a system of algebraic equations in all its variables  $X_1, Y_1, X_1^\bullet$  (ensuring consistent initialization).

Focus on  $S_{1\Sigma}$  in decomposition (73). Denote by  $\Phi$  the subsystem of  $S_{1\Sigma}$  collecting all its equations that also belong to  $\bullet S_{0\Sigma}$ , and are thus already solved.

The remaining system  $S_{1\Sigma} \setminus \Phi$  must then be solved for its dependent variables, which consist of the subset of  $X_1 \cup Y_1 \cup X_1^\bullet$  collecting all variables not yet determined by  $\bullet S_{0\Sigma}$ .

To prepare for this, we submit  $S_{1\Sigma} \setminus \Phi$  to the Dulmage-Mendelsohn (DM) decomposition, see Definition 13:  $(\beta_{\mathfrak{U}}, \beta_{\mathfrak{E}}, \beta_{\mathfrak{D}}) \leftarrow \text{DM}(S_{1\Sigma} \setminus \Phi)$ . If  $\beta_{\mathfrak{U}} = \beta_{\mathfrak{D}} = \emptyset$ , then  $S_{1\Sigma} \setminus \Phi$  is structurally nonsingular by Lemma 14 and we solve it. Otherwise, the following can occur:

*Case  $\beta_{\mathfrak{U}} \neq \emptyset$ .* This points to an underspecified subsystem at the instant of mode change. This information is returned to the designer for a correcting action.

*Case  $\beta_{\mathfrak{U}} = \emptyset$  and  $\beta_{\mathfrak{D}} \neq \emptyset$ .* This points to a conflict between  $\bullet S_{0\Sigma}$  and  $S_{1\Sigma}$  and is handled differently. We solve this conflict by applying Principle 1 of causality. Accordingly, the equations belonging to  $\beta_{\mathfrak{D}}$  are removed from the consistency equations  $(e_1)$  of the new system  $S_{1\Sigma}$ . We thus keep the system  $S_{1\Sigma}^\downarrow =_{\text{def}} \{(e_0), (e_1) \setminus \beta_{\mathfrak{D}}, (e_1^\bullet)\}$ . By Lemma 15, since the system  $\{(e_0), (e_1^\bullet)\}$  is structurally nonsingular, performing the DM decomposition on system  $S_{1\Sigma}^\downarrow$  returns a triple with an empty overdetermined block. Hence, the so obtained reduced system  $S_{1\Sigma}^\downarrow$  is structurally regular and ready to be solved.

**Comment 34.** If the two systems  $S_{i\Sigma}, i = 0, 1$  are identical (there is no mode change), then  $X_1 = X_0, Y_1 = Y_0, F_1 = F_0$ , and  $H_1 = H_0$ , so that  $\Phi$  consists of the consistency subsystem  $0 = H_1(X_1)$ , see the definition of  $S_{1\Sigma}$  in (74). Hence,  $S_{1\Sigma} \setminus \Phi$  coincides in this case with the structurally nonsingular system obtained by performing index reduction (see Definition 22), and we are left with the usual structural analysis for DAE.

The reconciliation may take several nonstandard instants following the mode change, where less and less equations from  $S_{1\Sigma}$  are erased when applying DM, until we reach the steady regime. The interpretation of the erasure of some equations from  $S_{1\Sigma}$  is that their satisfaction gets postponed for a few instants. Recall that we are in nonstandard semantics, so that all of this will occur in a total zero duration, in standard real time.

In the rest of this section, we formalize this intuition and generalize it.

### 6.3. Structural Analysis of long modes

In this section we develop the structural analysis of mdAE systems for the subclass of systems possessing only long modes. Within each long mode, we simply reuse the existing

structural analysis based on the  $\Sigma$ -method. The only difficulty sits in the mode changes, for which we use conflict reconciliation as explained above. Throughout this section, we work in the nonstandard semantics as given by Equation (70).

### 6.3.1. Constructive Semantics

We reuse the techniques of *constructive semantics*, first introduced in the context of reactive synchronous programming languages [12, 9, 10, 47], for the purpose of grounding compilation on solid mathematical foundations. For synchronous languages, the execution of a program proceeds through successive *reactions*, by which discrete time progresses from the current instant to the next one. The compilation of a program consists in generating the code for executing a reaction. This task is formalized through the notion of *constructive semantics*, which relies on the following two pillars:

1. The specification of the set of *atomic actions*, which are effective, non-interruptible, state transformations. For synchronous languages, atomic actions consist of: (a) the evaluation of a single expression, and (b) control flow actions. Executing an instance of an atomic action is referred to as performing a *micro-step*.
2. The *causality analysis* of the reaction, which is a partial order on the set of micro-steps, abstracting the dependencies between them. A scheduling of the micro-steps is correct if and only if it complies with the causality analysis.

The *constructive semantics* consists in decomposing a reaction as the symbolic execution of a sequence of micro-steps subject to causality constraints. Intermediate stages of this execution are represented by the ‘status’ of each variable, belonging to the abstract alphabet {not evaluated/evaluated}. The assignment of a status to all variables is called a *configuration*. Such executions are monotonic, in that the so produced sequence of configurations is increasing with respect to the product order derived from the order *not evaluated* < *evaluated* for each variable. The constructive semantics is complete if it terminates by having all variables evaluated. Failure occurs if an incomplete sequence cannot be extended.

This approach is suited to mdAE systems, since it will provide the formal support for correctly chaining the different actions composing the structural analysis, and describing the resulting increase in knowledge on all involved equations, variables and guards. For mdAE systems, the *atomic actions*, however, consist of:

- (a) the evaluation of a guard;
- (b) manipulations of systems of equations: adding latent equations or deleting conflicting equations at mode changes; and
- (c) solving a block of algebraic equations.

Actions (a) and (b) participate in the structural analysis, whereas (c) is the runtime solving of systems of equations to evaluate the leading variables. To keep this evaluation symbolic, we abstract away the values of numerical variables and regard the values of guards as oracles, so that all the possible valuations must be explored.

*Statuses and Contexts.* The different statuses of the guards, variables and equations of an mdAE are the following: an evaluated guard may be *true* or *false*; a variable may be *not evaluated*, or *evaluated*; an equation may be *not evaluated*, *disabled* (its guard is

false), or *solved*. This is summed up in Table 1. There is no need to consider the status *not evaluated* for a guard, since guards can be evaluated at the initialization of each nonstandard instant, by Definition 31.

	I	U	F	T
guard	irrelevant		evaluated to F	evaluated to T
variable	irrelevant	not evaluated		evaluated
equation	irrelevant	not evaluated	disabled	solved

Table 1: The partially ordered domain  $D = \{I, U, F, T\}$  and its interpretation for guards, variables, and equations.

Unlike for single-mode DAEs, the set of equations defining the dynamics of an mdAE is mode-dependent. This information should therefore be traced as part of the constructive semantics of mdAE systems. To capture this, we tag as *irrelevant* in the considered mode all the variables, guards, and equations that are not used to define the dynamics of the mdAE system in this mode. Formally, the domain of statuses is the partially ordered set  $(D, <)$  defined by

$$D = \{I, U, F, T\} \quad \text{and} \quad I < U < F, T . \quad (75)$$

The interpretation of this domain for variables, equations and guards, is summarized in Table 1 and detailed next:

- The minimal element I is used to represent the fact that a variable, a guard, or an equation is *irrelevant*.
- Value U means that a variable or equation has *not been evaluated yet*. At the beginning of a time step, only state variables are known, and all other variables are set to U, reflecting that their numerical values are not known yet.
- Maximal element T has different meanings, depending on whether it applies to a variable, a guard or an equation. In the case of a variable, it means that the numerical value of the variable has been computed, whatever it could be. For a guard, it means that the guard has been evaluated to true. For an equation, it means that the equation has been solved.
- Maximal element F also has different meanings, depending on whether it applies to a guard or an equation. In the context of a guard, it means that the guard has been evaluated to false. When it applies to an equation, it means that the equation is disabled, i.e., its guard is false. This value does not apply to variables.

We call *status* the assignment, to each guard, variable, and equation, of a value over  $D$ . It is formally defined as follows:

**Definition 35 (status).** *Let  $S$  be an mdAE system. Its set  $\mathbb{V}$  of S-variables<sup>19</sup> is*

$$\mathbb{V} \stackrel{\text{def}}{=} \Gamma \cup X^{(\bullet)} \cup E^{(\bullet)} \quad (76)$$

<sup>19</sup>The prefix “S” is reminiscent of constructive Semantics.

(using Notations 24 and 26). A status  $\sigma$  of  $S$  is a valuation in  $D$  of the  $S$ -variables, that is, a mapping  $\sigma : \mathbb{V} \rightarrow D$ . A status is called *finite* if it is equal to  $\mathbf{I}$ , except for a finite number of  $S$ -variables. The set of statuses is partially ordered by the product order:  $\sigma_1 \leq \sigma_2$  if and only if for all  $v \in \mathbb{V}$ ,  $\sigma_1(v) \leq \sigma_2(v)$ .

We denote by  $\mathbf{I} \nearrow \mathbf{U}$ ,  $\mathbf{U} \nearrow \mathbf{F}$ , and  $\mathbf{U} \nearrow \mathbf{T}$  the increasing changes of status for a given  $S$ -variable.

*Enabled Equations and Leading Variables.* Let  $\sigma$  be a status. Equation  $e^{\bullet m}$  is *enabled in*  $\sigma$  (respectively *disabled in*  $\sigma$ ) if and only if  $\sigma(\gamma) = \mathbf{T}$  (respectively  $\sigma(\gamma) = \mathbf{F}$ ), where  $\gamma$  is the guard of  $e^{\bullet m}$ , see (66). Denote by

$$\text{Enab}(\sigma) \text{ (respectively } \text{Disab}(\sigma))$$

the set of equations that are enabled (respectively disabled) in  $\sigma$ . Recall that, for any finite status  $\sigma$ , these sets are finite. The set  $\text{Enab}(\sigma)$  defines the DAE system that is enabled at status  $\sigma$ . The *leading variables* in status  $\sigma$  are defined, with reference to this enabled DAE system, following the introduction of Section 5.2. We define the set

$$\text{Undef}(\sigma) \stackrel{\text{def}}{=} \{v \in \mathbb{V} \mid \sigma(v) \leq \mathbf{U}\}$$

collecting the  $S$ -variables that are either irrelevant or not evaluated in status  $\sigma$ . We are now ready to introduce micro-steps, and runs as finite sequences of micro-steps.

**Definition 36 (micro-step).** A micro-step is a transformation, via some atomic action, of a status  $\sigma$  to a status  $\sigma'$  by updating the values of a finite subset of  $S$ -variables, from  $\mathbf{I}$  to  $\mathbf{U}$ , or from  $\mathbf{U}$  to  $\mathbf{T}$  or  $\mathbf{F}$ .

A run is a finite sequence of micro-steps. It is complete when every equation is either irrelevant, disabled, or solved, and all the enabled leading variables are evaluated. Formally (the reader is referred to Definition 35 for the notions used):

**Definition 37 (runs).** Given a finite initial status  $\sigma_0$ , a run of mdAE system  $S$  is a finite increasing sequence of statuses

$$\sigma_0 < \sigma_1 < \dots < \sigma_k < \sigma_{k+1} < \dots < \sigma_{\text{end}} \tag{77}$$

such that, for every  $k < \text{end}$ ,  $\sigma_{k+1}$  is obtained from  $\sigma_k$  by performing a micro-step. The final status  $\sigma_{\text{end}}$  is called *complete* if, in  $\sigma_{\text{end}}$ , all equations  $e_i$  have either the value  $\mathbf{T}$  or  $\mathbf{F}$  and no leading variable has the value  $\mathbf{U}$ . A run  $\sigma_0 < \dots < \sigma_{\text{end}}$  of mdAE system  $S$  is complete if and only if its final status  $\sigma_{\text{end}}$  is complete.

*Handling conflicts between past and present.* To handle the conflicts between the previous and current instants when a mode change occurs, we will need to keep track of the possible sources of conflict between previous and current instants. If equation  $e^{\bullet(m+1)}$ ,  $m \geq 0$ , was solved at the previous instant (meaning that  $\sigma(e^{\bullet(m+1)}) = \mathbf{T}$  for the final status of that instant), then  $e^{\bullet m}$  is a possible source of conflict at the current instant. The set of all such equations is called the *context* at the current instant. All the variables involved in the context are already evaluated at the initial status of the current run (value  $\mathbf{T}$ ).

If, however, equation  $e^{\bullet m} \in \Delta$  is also enabled at status  $\sigma$  of the current instant, then  $e^{\bullet m}$  is no longer a source of conflict, but rather a satisfied consistency equation. We call it a *fact* at  $\sigma$ . Formally:

**Definition 38 (contexts and facts).** Let  $S$  be an mdAE system. Its context  $\Delta$  at the current instant, and set of facts  $\Phi(\sigma)$  at the current status, are defined as follows:

$$\begin{aligned}\Delta &= \{e^{\bullet m} \mid \text{the final status of } e^{\bullet(m+1)} \text{ at previous instant is } \top\} \\ \Phi(\sigma) &= \Delta \cap \text{Enab}(\sigma)\end{aligned}$$

By abuse of notation, we will also denote by  $\Delta$  and  $\Phi$  the sets of functions  $f(e)$ , for  $e$  ranging over  $\Delta$  and  $\Phi$ .

To prepare for the algorithm computing a run for given initial status  $\sigma_0$  and context  $\Delta$ , we now collect a few useful auxiliary algorithms.

### 6.3.2. Auxiliary algorithms

**Tick algorithm** (Algorithm 7). When a run is complete, the system can proceed to the next time step by executing the algorithm Tick. Tick requires a complete (final) status  $\sigma$  and returns a pair  $(\sigma^\circ, \Delta^\circ)$  consisting of an initial status  $\sigma^\circ$  and a context  $\Delta^\circ$ , for use as initial conditions for the next run.

---

#### Algorithm 7 Tick

---

**Require:**  $\sigma$ , a complete status; **return** pair  $(\sigma^\circ, \Delta^\circ)$ .

1:  $\text{Tick}(\sigma) =_{\text{def}} (\sigma^\circ, \Delta^\circ)$ , where:

$$\begin{aligned}\sigma^\circ(\gamma) &= \text{if } \gamma \text{ is involved in } S \text{ then } \sigma(b(\gamma)) \text{ else } \text{I} \\ \sigma^\circ(x^{\bullet m}) &= \text{if } \sigma(x^{\bullet(m+1)}) = \top \text{ then } \top \text{ else} \\ &\quad \text{if } x^{\bullet m} \text{ is a variable of } S \text{ then } \text{U} \text{ else } \text{I} \\ \sigma^\circ(e^{\bullet m}) &= \text{if } e^{\bullet m} \text{ is a variable of } S \text{ then } \text{U} \text{ else } \text{I} \\ \Delta^\circ &= \{e^{\bullet m} \mid e \in E, m \in \mathbb{N} : \sigma(e^{\bullet(m+1)}) = \top\}\end{aligned}$$


---

Performing Tick has the effect of shifting backward defined variables, and setting all other S-variables  $v \in \mathbb{V}$  to either U, if  $v$  is relevant to the mdAE  $S$ , or I (irrelevant) otherwise. Guards are set either to I or to the value  $\sigma(b(\gamma)) \in \{\text{F}, \top\}$ , using Notations 32 and the definition (71) for  $\gamma^\bullet$ . Thus, the value of all guards at the next instant is set. The new context is defined to be the set of equations that are known to be satisfied in the next instant. Tick is not increasing and it does not have to be, since it applies when moving to the next time step.

**IndexReduc algorithm** (Algorithm 8). Given a system  $G$  with its set of leading variables determined by the status  $\sigma$ , the IndexReduc algorithm is a renaming of the  $\Sigma$ -method (Algorithm 6), together with a refined interpretation of its result. *It is understood that we use the discrete-time adaptation of it, in which forward shift replaces differentiation but everything else remains unchanged.*

Recall that this algorithm requires a square system  $G=0$  and returns a success/failure flag  $b$ , as well as, in case of success, a pair consisting of  $G_\Sigma=0$ , a structurally nonsingular system determining some leading variables, and the system  $\overline{G}_\Sigma=0$  collecting the consistency equations. As a consequence, the status of every latent equation added by the  $\Sigma$ -method is updated as  $\text{I} \nearrow \text{U}$ .

---

**Algorithm 8** IndexReduc

---

**Require:**  $(G, \sigma)$ ; **return**  $(\sigma, b, G_\Sigma, \overline{G}_\Sigma)$ 

- 1:  $(b, G_\Sigma, \overline{G}_\Sigma) \leftarrow \text{SigmaMethod}(G)$  with dependent variables set by  $\sigma$ ; increase  $\sigma$
  - 2: **if**  $b$  **then return**  $(G_\Sigma, \overline{G}_\Sigma)$
- 

*SolveConflict algorithm* (Algorithm 9). This algorithm requires a system  $K=0$  and a status  $\sigma$  for all variables involved in  $K$ . System  $K$  is submitted to the DM decomposition (Line 1); in doing so, the dependent variables are those with status  $\cup$ . Possible conflicts are collected in  $\beta_\mathcal{D}$ . If the latter is nonempty, corresponding equations are deleted (Line 3). The status is updated as  $\cup \nearrow \mathbb{F}$  for every removed equation. The so reduced system is again submitted to the DM decomposition (Line 4), which, by Lemma 15, will return an empty overdetermined block: conflicts are removed. If block  $\beta_\mathcal{U}$  is nonempty, then system  $K$  is underdetermined. Therefore,  $b = \mathbb{F}$  is returned (Line 5). Otherwise,  $b = \mathbb{T}$  and  $H = \beta_\mathcal{E}$  are returned (Line 6).

---

**Algorithm 9** SolveConflict

---

**Require:**  $(K, \sigma)$ ; **return**  $(\sigma, b, H)$ 

- 1:  $(\beta_\mathcal{U}, \beta_\mathcal{E}, \beta_\mathcal{D}) \leftarrow \text{DM}(K)$  with dependent variables set by  $\sigma$
  - 2: **if**  $\beta_\mathcal{D} \neq \emptyset$  **then**
  - 3:    $K \leftarrow K \setminus \beta_\mathcal{D}$ ; increase  $\sigma$
  - 4:    $(\beta_\mathcal{U}, \beta_\mathcal{E}, \emptyset) \leftarrow \text{DM}(K)$  with dependent variables set by  $\sigma$
  - 5:  $b \leftarrow [\beta_\mathcal{U} = \emptyset]$
  - 6: **if**  $b$  **then**  $H \leftarrow \beta_\mathcal{E}$
- 

*Eval algorithm.*

$$\sigma \leftarrow \text{Eval}(H, \sigma) \tag{78}$$

This algorithm requires a status  $\sigma$  and a structurally regular system  $H$ , enabled at  $\sigma$ . It evaluates  $H$  for its dependent variables, which sets the values of the leading variables. The corresponding change in status is  $\cup \nearrow \mathbb{T}$  for both the variables that are evaluated and the equations that are solved.

### 6.3.3. Executing a nonstandard instant

*ExecRun* requires a finite status  $\sigma$  and a finite context  $\Delta$  to serve as initial conditions for the run. It returns, either a documented *Fail*, or a final value for status  $\sigma$ , as well as values for  $\sigma^\circ$  and  $\Delta^\circ$  to serve as initial conditions for the next instant. A line-by-line description of *ExecRun* is given next.

**Line 1** Not all the equations belonging to the context are active in the current instant. The active ones are collected in the set  $\Phi$  of *facts*.

**Line 2** System  $G$  is set to the enabled guarded equations that are not evaluated yet in status  $\sigma$ .

**Line 3** We submit  $G$  to index reduction (*IndexReduc* algorithm, Algorithm 8).



---

**Algorithm 10** ExecRun

---

**Require:**  $(\sigma, \Delta)$ ; **return**  $(Fail, (\sigma^\circ, \Delta^\circ))$ 

```
1:  $\Phi \leftarrow \{e \in \Delta \mid \sigma(\gamma(e)) = \text{T}\}$ 
2:  $G \leftarrow [Enab(\sigma) \cap Undef(\sigma)]$ 
3:  $(\sigma, b, G_\Sigma, \overline{G}_\Sigma) \leftarrow \text{IndexReduc}(G, \sigma)$ ; increase  $\sigma$ 
4: if  $\neg b$  then return  $Fail(\sigma)$ 
5: else
6:    $(\sigma, b, H) \leftarrow \text{SolveConflict}((G_\Sigma \cup \overline{G}_\Sigma) \setminus \Phi, \sigma)$ ; increase  $\sigma$ 
7:   if  $\neg b$  then return  $Fail(\sigma)$ 
8:   else
9:      $\sigma \leftarrow \text{Eval}(\sigma, H)$ ; increase  $\sigma$ 
10:     $(\sigma^\circ, \Delta^\circ) \leftarrow \text{Tick}(\sigma)$ 
```

---

**Line 4** If IndexReduc returns  $b=\text{F}$ , then ExecRun returns  $Fail(\sigma)$  and stops.

**Line 5** Otherwise, ExecRun returns a pair consisting of  $G_\Sigma=0$ , a structurally nonsingular system determining some leading variables, and the system  $\overline{G}_\Sigma=0$  collecting consistency equations. In this case, the algorithm can further progress.

**Line 6** Taking the set  $\Phi$  of facts into account, the possible conflict is solved using the SolveConflict algorithm (Algorithm 9). When applying SolveConflict, the dependent variables of  $(G_\Sigma \cup \overline{G}_\Sigma) \setminus \Phi$  include all the variables (both leading and non-leading) of the current system that have status  $\sigma = \text{U}$ , i.e., whose value was not set by executing the previous nonstandard instant.

**Line 7** If SolveConflict fails, ExecRun returns  $Fail(\sigma)$  and stops.

**Line 9** Otherwise, SolveConflict returns a structurally regular system  $H$ , which is solved using Eval, and the so reached status  $\sigma$  is complete.

**Line 10** Since  $\sigma$  is complete, Tick (Algorithm 7) is performed, so that the values of guards, the status of all S-variables and the context are known for the next nonstandard instant.

#### 6.3.4. Important properties

*Evaluating guards.* Since ExecRun (Algorithm 10) performs a *symbolic* interpretation of the model, all the possible modes (valuations for the guards) must be hypothesized and explored at Line 10.<sup>20</sup> The program is accepted if and only if, for all possible modes, no *Fail* is returned by ExecRun. In case of success, all reachable pairs  $(\sigma, \Delta)$  of statuses and contexts have been explored.

*ExecRun as a labeled automaton.* We can see ExecRun as a labeled finite state automaton having nodes labeled with the different encountered status-context pairs, and transitions labeled with the micro-steps. In Figure 11, we depict this graph for the clutch example. The Tick transitions indicate a move to the next nonstandard instant.

---

<sup>20</sup> Assertions may be used to restrict the set of possible modes, by stating some physical knowledge the programmer may have, e.g., ‘this sequence of modes is not possible’. This is a minor change to our approach, not developed here.

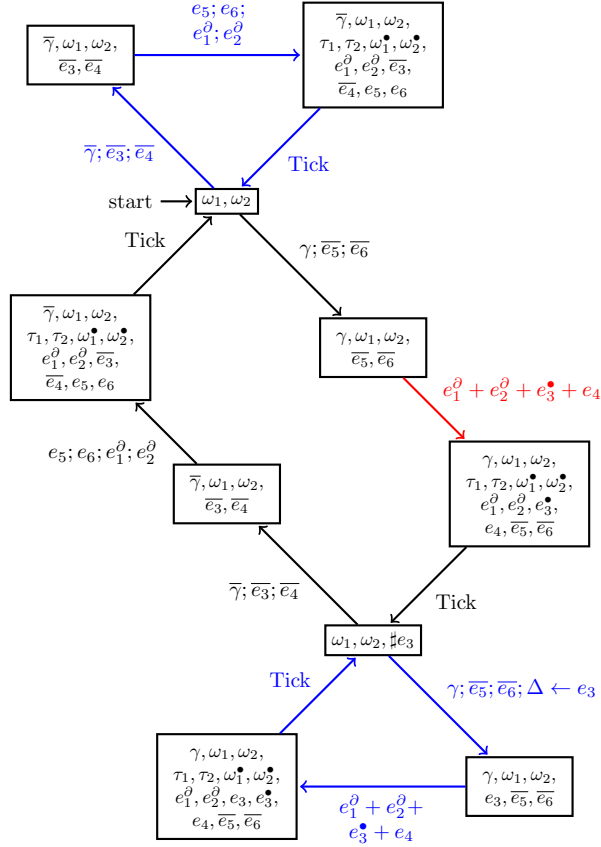


Figure 11: ExecRun as a labeled automaton for the Simple Clutch. Notations: For all statuses (shown in boxes),  $v$  (respectively  $\bar{v}$ ) means  $v = \text{T}$  (respectively  $v = \text{F}$ ), and not mentioning  $v$  means  $v = \text{U}$ .  $\#e$  means that  $e_f$  belongs to context  $\Delta$ . The assignment  $\Delta \leftarrow e_3$  refers to Line 10 of Algorithm 10. Blue (respectively black) transitions belong to a continuous-time (respectively discrete-time) dynamics. The red transition is impulsive. A semicolon is the sequential composition of computations, and the  $+$  symbol defines blocks of equations.

*Rejecting models.* Firstly, a system exhibiting logico-numerical fixpoint equations is rejected via syntactic checks, since it does not agree with Definition 31. Then, a failure is found in the following cases:

- **SigmaMethod** (Algorithm 6) itself fails, indicating that no matching is found for the bipartite graph  $\mathcal{G}_G$  of the running dAE system  $G=0$ , a failure criterion for the  $\Sigma$ -method. Failure message indicates either over- or under-determination.
- The conflict between the predictions from the past and the consistency conditions from the current mode cannot be solved, i.e.,  $\beta_{\mathfrak{M}} \neq \emptyset$  at Line 5 of **SolveConflict**, revealing some missing restart equation(s) for the new mode.

*Detecting continuous modes.* Elementary cycles of **ExecRun** capture runs with stationary valuations of the guards and define the continuous dynamics in each mode. Other runs capture mode changes and their reset actions. During runs with constant valuations of the guards, we have  $(G_{\Sigma} \cup \overline{G}_{\Sigma}) \setminus \Phi = G_{\Sigma}$ . Hence, Line 1 returns  $\beta_{\mathfrak{M}} = \beta_{\mathfrak{D}} = \emptyset$  and, thus,  $\beta_{\mathfrak{E}} = G_{\Sigma}$ : we recover the usual structural analysis for DAE systems.

*Solving conflicts at mode changes takes only finitely many nonstandard instants.* This follows from the fact that, at a mode change, the maximal shifting degree of context  $\Delta$  is finite and the maximal shifting degree of the conflicting subsystem decreases by one at each subsequent time step. See Section 3.3 on the Cup-and-Ball example as an illustration of this.

*Determinism of the algorithm.* All the algorithms called throughout the execution of Algorithm 10 possess a unique return (at Line 3 and Line 6). Since Algorithm 10 involves no choice by itself, it is deterministic. Comment 16 following Lemma 15 finds its justification here: our design choice in Lemma 15 was essential in guaranteeing the determinism of our compilation method.

### 6.3.5. Generic Blocks within a long mode

The results of this section will be used in Section 8, and more specifically in Section 8.4. The successive application of Line 3 and Line 6 of Algorithm 10 yields blocks of a specific form, that is investigated next.

**Notations 39.** For  $f : \mathbb{R}^K \rightarrow \mathbb{R}$  a smooth numerical function and  $n \geq 0$  an integer, set

$$\begin{aligned}
 B_{\bullet}(f, n) &=_{\text{def}} \begin{bmatrix} B_{\bullet}^{\text{h}}(f, n) \\ B_{\bullet}^{\text{t}}(f, n) \end{bmatrix} \text{ where} \\
 B_{\bullet}^{\text{h}}(f, n) &=_{\text{def}} \begin{bmatrix} f \\ f^{\bullet} \\ \vdots \\ f^{\bullet(n-1)} \end{bmatrix} \text{ and} \\
 B_{\bullet}^{\text{t}}(f, n) &=_{\text{def}} f^{\bullet n} .
 \end{aligned} \tag{79}$$

$B_{\bullet}(f, n)$  is called an  $(f, n)$ -block,  $B_{\bullet}^{\text{h}}(f, n)$  is its *head* and  $B_{\bullet}^{\text{t}}(f, n)$  is its *tail*. By convention,  $B_{\bullet}(f, n)$  is empty for  $n < 0$  and  $B_{\bullet}^{\text{h}}(f, n)$  is empty for  $n = 0$ . For  $\{f_1, \dots, f_k, \dots, f_K\}$  a

tuple of  $K$  smooth functions in  $K$  dependent variables, define

$$F = \begin{bmatrix} f_1 \\ \vdots \\ f_K \end{bmatrix}, \quad N = \begin{bmatrix} n_1 \\ \vdots \\ n_K \end{bmatrix}, \quad B_\bullet(F, N) = \begin{bmatrix} B_\bullet(f_1, n_1) \\ \vdots \\ B_\bullet(f_K, n_K) \end{bmatrix} \quad (80)$$

and the definitions of  $B_\bullet^h(F, N)$  and  $B_\bullet^t(F, N)$  follow accordingly.  $\square$

**Lemma 40.** *Assume that the system is within a long mode. Then, successively applying Line 3 and Line 6 of Algorithm 10 to the system  $G=0$  returns, either ‘Fail’, or a structurally nonsingular system  $H=B_\bullet^t(G, N)$  of the form (80).*

*Proof.* Consider Line 3 of Algorithm 10. Since the system is in a long mode, SigmaMethod returns, either *Fail*, or  $(G_\Sigma, \bar{G}_\Sigma) = (B_\bullet^t(G, N), B_\bullet^h(G, N))$ , where  $G$  is the subset of enabled equations at status  $\sigma$ .

Move to Line 6. Since the system is in a long mode, we have  $\bar{G}_\Sigma = \Phi$ , expressing that consistency conditions are guaranteed by previous nonstandard instants. Hence, the Dulmage-Mendelsohn decomposition performed at Line 6 simply returns  $\beta_{\mathfrak{E}} = G_\Sigma = B_\bullet^t(G, N)$ .  $\square$

#### 6.4. Structural Analysis: general case

So far, we restricted ourselves to systems having only long modes. In this section, we consider the general case, where transient modes can also occur, possibly in cascades. The Cup-and-Ball example of Section 3 can be kept in mind while reading this section.

Our main task is to revisit index reduction. So far, this was performed by invoking the  $\Sigma$ -method, which requires the new mode to last for long enough to allow for an unlimited number of time shifts to be applied. If the current mode is transient, this no longer holds and the continuation of the system dynamics is determined by its successive future modes—see the analysis of the Cup-and-Ball example in Section 3.5.

Whereas the  $\Sigma$ -method applies to time-invariant systems only, the method of *differentiation arrays* (see Section 3.1 of [21]) can be adapted to time-varying systems.

##### 6.4.1. Index reduction with Difference Arrays

The index reduction algorithm proposed in this section applies to time-varying DAE systems. It is, however, less computationally efficient than the  $\Sigma$ -method (when the latter can be applied).

##### *Difference arrays for time-varying systems*

In the following,  $\mathbb{T}$  denotes a discrete-time index set possessing a minimal element, denoted by 0.

**Definition 41.** *A time-varying dAE system over  $\mathbb{T}$  is a pair  $(X, \mathcal{F})$ , where  $X$  is a tuple of numerical variables taking its values in an Euclidian space  $D$ , and  $\mathcal{F} : \mathbb{T} \rightarrow \mathbb{E}$  maps every instant  $t \in \mathbb{T}$  to a dAE system  $F_t = 0$  whose set of dependent variables is contained in  $X$ . A solution of  $(X, \mathcal{F})$  is a trajectory  $\mathbf{x} : \mathbb{T} \rightarrow D$  such that, for every  $t \in \mathbb{T}$ ,  $\mathbf{x}(t)$  satisfies  $F_t = 0$ .*

**Example 42.** Consider  $\mathbb{T} = \mathbb{N}$ ,  $X = (u, v, w)$ , and  $\mathcal{F}$  is defined by

$$t=0 : 0 = v^\bullet - 3u \quad ; \quad \text{and } \forall t > 0 : \begin{cases} 0 = u^\bullet - v \\ 0 = v - uw \\ 0 = w^\bullet + u \end{cases} .$$

The solutions of this system are trajectories  $t \mapsto (u(t), v(t), w(t))$  such that

$$0 = v(1) - 3u(0) \quad ; \quad \forall t > 0 : \begin{cases} 0 = u(t+1) - v(t) \\ 0 = v(t) - u(t)w(t) \\ 0 = w(t+1) + u(t) \end{cases} .$$

Note that solutions are only partially specified at  $t = 0$  since  $w$  is not involved in the initial conditions.  $\square$

Given a time-varying dAE system  $(X, \mathcal{F})$ , we consider, for each  $k \in \mathbb{T}$ , the following dAE array:

$$\mathcal{A}_k \stackrel{\text{def}}{=} \begin{bmatrix} F_0 \\ F_1 \\ \vdots \\ F_k \end{bmatrix} . \quad (81)$$

This array allows us to perform the ‘index reduction’ of  $F_0$ , by adding ‘latent equations’ taken from the future systems  $F_1, F_2$ , etc., generally different from  $F_0^\bullet, F_0^{\bullet 2}$ , etc. To this end, the set of all variables of array  $\mathcal{A}_k$  decomposes as  $Y \cup X \cup W$ , where:

$$\begin{aligned} Y &\text{ collects the variables in } F_0 \text{ that participate in the consistency;} \\ X &\text{ collects the leading variables of } F_0; \\ W &\text{ collects all the remaining variables of the array.} \end{aligned} \quad (82)$$

In its Section 3.1 on ‘standard indices’, Campbell and Gear’s landmark paper [21] states that, for time-invariant systems ( $F_t \equiv F_0$ ), the index is the smallest  $k$  such that  $\exists W : \mathcal{A}_k(X, W, Y) = 0$  uniquely defines  $X$  as a function of  $Y$ , provided that  $Y$  is consistent.

Therefore, we consider the following structural translation of the above problem (compare with Problem 21 in Section 5.2):

**Problem 43.** *Find the smallest integer  $k$  such that the system  $\exists W : \mathcal{A}_k(X, W, Y) = 0$  is structurally nonsingular in the sense of Definition 19.*

For time-invariant systems ( $F_t \equiv F$  for all  $t$ ), Problem 43 is monotonic in that, if this problem has a solution for  $k$ , then it also has a solution for  $k' > k$ . However, this no longer holds for time-varying systems, since it can be that system  $F_{k+1} = 0$  is in conflict with  $\mathcal{A}_k$ .

#### *Handling finite cascades of transient modes*

Now, we specialize the structural analysis of difference arrays to the handling of finite cascades of transient modes. That is, we assume the following finite chain of modes:

$$\underbrace{\mu_0, \mu_1 \dots, \mu_K}_{\text{finite cascade of transient modes}} \quad , \quad \underbrace{\mu_\infty}_{\text{final long mode}} \quad (83)$$

We adopt the following principles in developing our structural analysis:

---

**Algorithm 11** DiffArray
 

---

**Require:**  $(\mathcal{A}, \sigma)$  as in (84); **return**  $(\sigma, b, F_\Sigma, \overline{F}_\Sigma)$

```

1:  $k \leftarrow -1$ ;
2: if  $k < K$  then
3:    $k \leftarrow k+1$ 
4:    $\mathcal{A}_k \leftarrow (F_0, \dots, F_k)$ 
5:   arrange the variables of  $\mathcal{A}_k$  as  $(Y_k \cup X_k \cup W_k)$ 
6:    $(b_\mathfrak{D}, b_\mathfrak{U}, F_\Sigma, \overline{F}_\Sigma) \leftarrow \text{ExistQuantifEqn}(\mathcal{A}_k)$ 
7:   if  $\neg b_\mathfrak{D}$  then return  $b \leftarrow \mathbb{F}$ 
8:   else
9:     if  $\neg b_\mathfrak{U}$  then go to Line 2
10:    else return  $(b \leftarrow \mathbb{T}, F_\Sigma, \overline{F}_\Sigma)$ ; increase  $\sigma$ 
11: else consider  $\mathcal{A}$  defined in (84)
12:    $\mathcal{A} \leftarrow \mathcal{A}$  where  $\overline{F}_{\infty, \Sigma} \leftarrow (\overline{F}_{\infty, \Sigma} \setminus \mathcal{A}_K)$ 
13:    $(\beta_\mathfrak{U}, \beta_\mathfrak{E}, \beta_\mathfrak{D}) = \text{DM}(\mathcal{A})$ 
14:   if  $\beta_\mathfrak{D} \neq \emptyset$  then
15:      $\mathcal{A} \leftarrow \mathcal{A}$  where  $\overline{F}_{\infty, \Sigma} \leftarrow \overline{F}_{\infty, \Sigma} \setminus \beta_\mathfrak{D}$ ;
16:     go to Line 13
17:   else
18:     if cond. 2 and 3 of Lemma 20 hold then
19:       partition  $\beta_\mathfrak{E} = F_\Sigma \cup \overline{F}_\Sigma$ 
20:       return  $(b \leftarrow \mathbb{T}, F_\Sigma, \overline{F}_\Sigma)$ ; increase  $\sigma$ 
21:     else return  $b \leftarrow \mathbb{F}$ 

```

---

- For the final long mode  $\mu_\infty$ , it is legitimate to postpone the satisfaction of some of the consistency equations. That is, we reuse the conflict solving approach of the algorithm SolveConflict (Algorithm 9).
- For the transient modes  $\mu_0, \mu_1, \dots, \mu_K$ , however, this is not legitimate, so that we require the corresponding equations to be satisfied and refuse to erase any of them.

We associate to the chain (83) the array

$$\mathcal{A} =_{\text{def}} \begin{bmatrix} F_0 \\ F_1 \\ \vdots \\ F_K \\ F_{\infty, \Sigma} \cup \overline{F}_{\infty, \Sigma} \end{bmatrix} \quad (84)$$

with a decomposition of its variables as  $Y \cup X \cup W$  following (82).

The last row of this array originates from applying index reduction on the long mode  $\mu_\infty$ .

Based on the results of Section 5.1.4, in particular Lemma 20, the structural analysis for time-varying systems is described in Algorithm 11 (DiffArray). This algorithm requires as inputs an array  $\mathcal{A}$  as in (84) and a status; it returns, either a *Fail* information ( $b = \mathbb{F}$ ), or a pair  $(F_\Sigma, \overline{F}_\Sigma)$ , where  $F_\Sigma = 0$  is a structurally nonsingular system extracted from  $\mathcal{A}$

and determining the leading variables of  $F_0$ , and  $\overline{F}_\Sigma=0$  is a solvable system of consistency equations. The algorithm proceeds according to two phases.

In a first phase (Line 1–Line 10), the subchain of transient modes is explored, with increasing array size  $k = 0, \dots, K$ . Using status  $\sigma$ , the variables of  $\mathcal{A}_k$  are sorted at Line 5 following (82). At Line 6, array  $\mathcal{A}_k$  is submitted to `ExistQuantifEqn` (Algorithm 5), which returns the two Boolean flags  $b_\Sigma$ , indicating whether the overdetermined block is empty, and  $b_\mathcal{U}$ , indicating whether the underdetermined block is empty. If  $b_\Sigma = \text{F}$ , which indicates an overdetermined system, a failure indication is returned at Line 7. If  $b_\Sigma = \text{T}$  and  $b_\mathcal{U} = \text{F}$  is returned, the array is increased (Line 3) and reprocessed the same way, until  $k = K$  is reached. Otherwise, the first phase returns the pair  $(F_\Sigma, \overline{F}_\Sigma)$ . The status is updated as  $\text{I} \nearrow \text{U}$  for each latent equation added, and the algorithm successfully terminates.

If  $k=K$  is reached with no success, we move to the second phase, starting with the removal of duplicates between  $\mathcal{A}_K$  and  $\overline{F}_{\infty, \Sigma}$  (Line 12). Then, we apply the DM decomposition at Line 13 to identify possible conflicts and underdetermined parts. We handle conflicts as in the `SolveConflict` algorithm (Algorithm 9), by deleting conflicting equations from  $\overline{F}_{\infty, \Sigma}$  (Line 15) and returning to Line 13; this loop is visited only once thanks to Lemma 15. When reaching Line 17, we know that condition 1 of Lemma 20 is satisfied, so that we can conclude by using the successful case of Lemma 20 at Lines 18 and subsequent. The status is updated as  $\text{I} \nearrow \text{U}$  for each latent equation added (Line 20).

#### 6.4.2. Executing a nonstandard instant

We first revisit the algorithm `IndexReduc`. If the current mode is long, then `IndexReduc` is implemented as Algorithm 8. Otherwise, we will implement it by hypothesizing a future mode trajectory and calling the algorithm `DiffArray` (Algorithm 11). `IndexReduc`, revisited, is specified as Algorithm 12.

---

#### Algorithm 12 `IndexReduc`, revisited

---

**Require:**  $(G, \sigma, \text{type})$ ; **return**  $(\sigma, b, G_\Sigma, \overline{G}_\Sigma)$

1: **if**  $\text{type}=\text{long}$  **then**

2:      $(b, G_\Sigma, \overline{G}_\Sigma) \leftarrow \text{SigmaMethod}(G)$  with dependent variables set by  $\sigma$ ; increase  $\sigma$

3: **else**

4:     hypothesize  $\mathcal{F} = (G_0, G_1, \dots)$  with  $G_0=G$ ; let  $\mathcal{A}$  be the corresponding array

5:      $(b, G_\Sigma, \overline{G}_\Sigma) \leftarrow \text{DiffArray}(\mathcal{A}, \sigma)$ ; increase  $\sigma$

6: **if**  $b$  **then return**  $(G_\Sigma, \overline{G}_\Sigma)$

---

In this algorithm, Line 4 requires hypothesizing a continuation for the current mode. This, of course, is problematic if all modes can be visited in this continuation, particularly transient modes. It is therefore essential to be able, at compile time, to restrict the number of possible continuations for the current mode. This question is investigated in Section 6.4.3.

Having revisited index reduction, we can now revisit `ExecRun`, see Algorithm 13. Only two lines are modified with reference to the original algorithm `ExecRun`:

**Line 1** We add the evaluation of the *type* (long or transient) of the current mode.

---

**Algorithm 13** ExecRun , revisited

---

**Require:**  $(\sigma, \Delta)$ ; **return**  $(Fail, \sigma, (\sigma^\circ, \Delta^\circ))$

- 1:  $type \leftarrow \text{EvalModeType}(\sigma)$
- 2:  $\Phi \leftarrow \{e \in \Delta \mid \sigma(\gamma(e)) = \text{T}\}$
- 3:  $G \leftarrow [Enab(\sigma) \cap Undef(\sigma)]$
- 4:  $(\sigma, b, G_\Sigma, \overline{G}_\Sigma) \leftarrow \text{IndexReduc}(G, \sigma, type)$ ; increase  $\sigma$
- 5: **if**  $\neg b$  **then return**  $Fail(\sigma)$
- 6: **else**
- 7:    $(\sigma, b, H) \leftarrow \text{SolveConflict}((G_\Sigma \cup \overline{G}_\Sigma) \setminus \Phi, \sigma)$ ; increase  $\sigma$
- 8:   **if**  $\neg b$  **then return**  $Fail(\sigma)$
- 9:   **else**
- 10:      $\sigma \leftarrow \text{Eval}(\sigma, H)$ ; increase  $\sigma$
- 11:      $(\sigma^\circ, \Delta^\circ) \leftarrow \text{Tick}(\sigma)$

---

**Line 4** We call instead the `IndexReduc` corresponding to Algorithm 12, which selects the appropriate method for index reduction, depending on the *type* of the current mode.

#### 6.4.3. Continuations of transient modes

We need to investigate this issue for transient modes only. Therefore, we begin with assumptions regarding transient modes. Recall that we follow Notations 24, hence, by *numerical variable* of an mDAE system  $S$ , we mean a variable of the form  $x^{/m}$  ( $m$ -th derivative of  $x$ ) for some algebraic or state variable  $x$  of the considered system. The *degree* of  $x^{/m}$  is  $m$ . The same holds regarding mdAE system  $S$ , with  $x^{\bullet m}$  ( $m$ -shifted version of  $x$ ) replacing  $x^{/m}$ . We repeatedly use these notations in the sequel.

*Assumptions regarding transient modes.* According to Definition 29, a mode  $\mu$  is transient if at least one of its T-valued guards is transient. Guards are Boolean expressions of predicates over numerical variables. We will restrict ourselves to the class of nonstandard mdAE models  $S$  with set  $X$  of numerical variables, satisfying the following assumption regarding transient predicates over numerical variables:

**Assumption 3.** *Let  $\gamma_1, \gamma_2 \in \Gamma(S)$  be any two distinct transient predicates over numerical variables of mdAE model  $S$ . Then, for  $\nu_X^1$  and  $\nu_X^2$  any two valuations of the tuple  $X$  such that  $\nu_X^1 - \nu_X^2$  is infinitesimal,  $\gamma_1(\nu_X^1) = \gamma_2(\nu_X^2) = \text{T}$  cannot occur.*

*Discussion.* Assumption 3 formalizes that two different transient predicates never take the value T simultaneously, “by chance”. For example, for  $x, y \in X$ , the two predicates  $\gamma_1 = \text{when } x \geq 0$  and  $\gamma_2 = \text{when } y \leq 5$  specify the events when  $x$  becomes  $\geq 0$  and when  $y$  becomes  $\leq 5$ . The assumption says that these two events will never occur simultaneously. For this case, Assumption 3 is very reasonable. Of course, if  $\gamma_1 = \text{when } x^3 \geq 8$  and  $\gamma_2 = \text{when } x \geq 2$ , the two guards are considered different (their syntax differ) but generate simultaneous events. It is not the duty of the compiler to discover if two guards happen to be mathematically identical whereas they differ in their syntax. It is rather the responsibility of the programmer to use identical syntax for different occurrences of a same predicate.  $\square$



Using Assumption 3 we can in particular bound the number of transient modes that can occur at a given instant  $t \in \mathbb{T}$ . For  $\gamma_o \in \Gamma(S)$  a transient predicate, let

$$\Gamma(S, \gamma_o) \subseteq \Gamma(S) \quad (85)$$

be the subset of all guards  $\gamma$  of  $S$  that are Boolean expressions of predicates such that  $\gamma = \top$  requires  $\gamma_o = \top$ . An example would be  $\gamma = \gamma_o \wedge \gamma'$ . The following obvious lemma holds:

**Lemma 44.** *Under Assumption 3, the number of transient modes that can possibly occur at a given instant  $t \in \mathbb{T}$  is bounded by*

$$M(S) = \max_{\gamma_o} |2^{\Gamma(S, \gamma_o)}|$$

where  $\gamma_o$  ranges over the set of all transient predicate guards of  $S$ , and  $|A|$  denotes the cardinal of set  $A$ .

In practice, this bound is generally much smaller than the number of possible modes for  $S$ , which is bounded by  $|2^{\Gamma(S)}|$ .

*Assumptions regarding continuations of transient modes.* Let  $\mu_0 = \mu(\mathbf{t})$  be the current mode at time  $\mathbf{t} \in \mathbb{T}$ , and let  $\mu_0, \mu_1, \mu_2, \dots$  be any continuation for  $\mu_0$  at instants  $\mathbf{t}, \mathbf{t}^\bullet, \mathbf{t}^{\bullet^2}, \dots$ . Such continuation can have one of the following two forms:

$$\text{finite cascade} : \mu_0, \mu_1, \dots, \mu_{m-1}, \mu_\infty \quad (86)$$

$$\text{infinite cascade} : \mu_0, \mu_1, \dots, \mu_{m-1}, \dots \quad (87)$$

where the final mode  $\mu_\infty$  in (86) is long (and thus repeated), and all other listed modes, including the "...", are transient. We will state assumptions that statically rule out (87) and limit the number of possible continuations in (86). For  $\mu$  a transient mode, let

$$X(\mu) \subseteq X$$

be the subset of numerical variables  $x$  such that  $x^\bullet - x$  is possibly non-infinitesimal at this mode. In other words,  $X(\mu)$  is the set of variables that are reset by this transient mode. We then define, for  $M$  a set of transient modes:

$$X(M) = \bigcup_{\mu \in M} X(\mu).$$

This set can be identified using structural algorithms.

**Example 45.** The statement

$$x' = f(x, u) \text{ reset } 0 \text{ when } x \geq 1$$

expands, according to Definition 25, as

$$\left\{ \begin{array}{ll} \text{when } x \geq 1 & \text{then } x^+ = 0 \\ & \text{else } x' = f(x, u) \end{array} \right.$$

whose nonstandard semantics is

$$\begin{cases} \text{when } x \geq 1 & \text{then } x^\bullet = 0 & (e_1) \\ & \text{else } \frac{x^\bullet - x}{\delta} = f(x, u) & (e_2) \end{cases} \quad (88)$$

Since  $x^\bullet - x$  is not guaranteed infinitesimal by  $(e_1)$ ,  $x^\bullet - x$  is non-infinitesimal at the instant of zero-crossing (i.e., the first time when  $x \geq 1$  occurs). Thus,  $X(\mu_{zc}) = \{x\}$ , where  $\mu_{zc}$  is the mode when the zero-crossing occurs. This can be found at compile time, based on the syntax of System (88).  $\square$

**Example 46.** Consider the Cup-and-Ball example with elastic impact, see Section 3.5, in its original specification prior to adding the restart condition (45). Let  $\mu_{\text{imp}}$  be the mode when the elastic impact occurs. Then,  $X(\mu_{\text{imp}}) = \emptyset$ . Thus, we know that the only continuation is:  $\mu_{\text{imp}}, \mu_{\text{free}}$ , in which  $\mu_{\text{free}}$  is the mode in which the rope is not straight and, thus, free motion occurs. The same holds after adding restart condition (45).  $\square$

In the sequel, we denote by  $M$  the set of all possible modes of system  $S$ . For  $Y \subset X$  a subset of variables, let

$$M(Y) \subseteq M$$

be the set of all transient modes that can possibly be visited as a consequence of making noninfinitesimal changes to the variables of  $Y$ —this set can be identified at compile time using structural algorithms.

In bounding the possible continuations of a transient mode  $\mu_0$ , we consider the following chain of sets of possibly reachable transient modes:

$$\begin{aligned} \text{Continuations } (\mu_0) &= M_0, M_1, M_2, \dots \\ \text{where : } M_0 &= \{\mu_0\} \\ \text{and } \forall k > 0 : M_k &= M(X(M_{k-1})) \end{aligned}$$

Once again thanks to structural algorithms, the following assumption can be checked at compile time on the considered system  $S$ :

**Assumption 4.** *For every transient mode  $\mu$ , the chain  $X_k, k \geq 0$  in  $\text{Continuations}(\mu)$  converges to  $\emptyset$  in finitely many steps.*

**Lemma 47.** *Assumption 5 implies Assumption 4,*

where Assumption 5 states:

**Assumption 5.** *For every subset  $Y$  of reset variables,*

$$d^\circ(X(M(Y))) > d^\circ(Y)$$

where  $d^\circ(Z)$  is the smallest (shifting or differentiation) degree among all variables of  $Z$ .

Assumption 5 holds if, when a transient mode is reached due to a predicate on positions  $x, y, \dots$ , then only velocities  $x', y', \dots$  or accelerations  $x'', y'', \dots$  can be reset. This kind of property is often encountered in classes of physical systems, e.g., contact mechanics. This finishes our study of structural analysis for mDAE or mdAE systems. Referring to Figure 9 in Section 4, we are ready to move to the two downgoing arrows, consisting in the standardization of the nonstandard model resulting from the structural analysis, together with the impulse analysis. The rest of the paper is devoted to this.

## 7. Background on nonstandard analysis

The development of the clutch example in Section 2 consisted of three main steps: (1) the mapping of the original clutch model to its nonstandard semantics, (2) the symbolic execution of this nonstandard semantics (Section 2.4), and (3) the back-standardization of this symbolic execution (Section 2.5). All of this resulted in the actual code for the simulation results reported in Figure 7. Formally justifying back-standardization requires a non-trivial use of nonstandard analysis, way beyond the quick and superficial introduction we gave in Section 2.3. In this and the next sections, we provide the mathematical background supporting back-standardization procedures in full generality.

*Nonstandard analysis* was proposed by Abraham Robinson in the 1960s to allow for the explicit manipulation of “infinitesimals” in analysis [49, 23]. Robinson’s approach was axiomatic, by adding three new axioms to the basic Zermelo-Fraenkel (ZFC) framework. An alternative presentation was later proposed by Lindstrøm [35]. Its interest is that it does not require any fancy axiomatic material but only makes use of Zorn’s lemma, equivalent to the axiom of choice in the ZF set theory. The proposed construction bears some resemblance to the construction of  $\mathbb{R}$  as the set of equivalence classes of Cauchy sequences in  $\mathbb{Q}$  modulo the equivalence relation  $\sim$  defined by  $(u_n) \sim (v_n)$  iff  $\lim_{n \rightarrow \infty} (u_n - v_n) = 0$ .

The important point for us is that nonstandard analysis allows the use of the nonstandard discretization of continuous dynamics “as if” it was operational and with infinitesimal discretization error. Iwasaki et al. [33] first proposed using nonstandard analysis to discuss the nature of time in hybrid systems. Bliudze and Krob [14, 13] have also used nonstandard analysis as a mathematical support for defining a system theory for hybrid systems. In their study of mathematical foundations of hybrid systems modelers of the ODE class, Benveniste et al. [5] used extensively the nonstandard semantics of hybrid systems, with practical consequences for how to design a specification formalism for multimode ODE systems [15, 4].<sup>21</sup> See also [51, 39, 29] for similar studies using the axiomatic approach to nonstandard analysis. The following presentation is borrowed from [5].

### 7.1. Intuitive introduction

We begin with an intuitive introduction to the construction of the nonstandard reals. The goal is to augment  $\mathbb{R} \cup \{\pm\infty\}$  by adding, to each  $x$  in the set, a set of elements that are “infinitesimally close” to it. We will call the resulting set  ${}^*\mathbb{R}$ . Another requirement is that all operations and relations defined on  $\mathbb{R}$  should extend to  ${}^*\mathbb{R}$ .

A first idea is to represent such additional numbers as convergent real sequences. For example, elements infinitesimally close to the real number zero are the sequences  $u_n = 1/n$ ,  $v_n = 1/\sqrt{n}$  and  $w_n = 1/n^2$ . Observe that the above three sequences can be ordered:  $v_n > u_n > w_n > 0$  where 0 denotes the constant zero sequence. Of course, infinitely large elements (close to  $+\infty$ ) can also be considered, e.g., sequences  $x_n = n$ ,  $y_n = \sqrt{n}$ , and  $z_n = n^2$ .

Unfortunately, this way of defining  ${}^*\mathbb{R}$  does not yield a total order since two sequences converging to zero cannot always be compared: if  $u_n$  and  $u'_n$  are two such sequences, the three sets  $\{n \mid u_n > u'_n\}$ ,  $\{n \mid u_n = u'_n\}$ , and  $\{n \mid u_n < u'_n\}$  may even all be infinite. The

---

<sup>21</sup>Corresponding to the Simulink class of modeling languages, with no DAE.

beautiful idea of Lindstrøm is to enforce that *exactly one of the above sets is important and the other two can be ignored*.

The key step in Lindstrøm's construction consists in fixing once and for all a finitely additive positive measure  $\mu$  over the set  $\mathbb{N}$  of integers with the following properties:<sup>22</sup>

$$\begin{aligned}\mu &: 2^{\mathbb{N}} \rightarrow \{0, 1\} \\ \mu(X) &= 0 \text{ whenever } X \text{ is finite} \\ \mu(\mathbb{N}) &= 1\end{aligned}$$

Once  $\mu$  is fixed, one can compare any two sequences  $u_n$  and  $u'_n$  as follows. Exactly one of the three sets  $\{n \mid u_n > u'_n\}$ ,  $\{n \mid u_n = u'_n\}$ , or  $\{n \mid u_n < u'_n\}$  has  $\mu$ -measure 1 whereas the other two must have  $\mu$ -measure 0. Thus, we say that

$$\begin{aligned}u > u' & \text{ if } \mu(\{n \mid u_n > u'_n\}) = 1 \\ u = u' & \text{ if } \mu(\{n \mid u_n = u'_n\}) = 1 \\ u < u' & \text{ if } \mu(\{n \mid u_n < u'_n\}) = 1\end{aligned}$$

respectively. Indeed, the same trick works for many other relations and operations on nonstandard real numbers, as we shall see. We now proceed with a more formal presentation.

## 7.2. Nonstandard domains

For  $I$  an arbitrary set, a *filter*  $\mathcal{F}$  over  $I$  is a family of subsets of  $I$  such that: (1) the empty set does not belong to  $\mathcal{F}$ ; (2)  $P, Q \in \mathcal{F}$  implies  $P \cap Q \in \mathcal{F}$ ; and (3)  $P \in \mathcal{F}$  and  $P \subset Q \subseteq I$  implies  $Q \in \mathcal{F}$ . Consequently,  $\mathcal{F}$  cannot contain both a set  $P$  and its complement  $P^c$ . A filter that contains one of the two for any subset  $P \subseteq I$  is called an *ultra-filter*. At this point we recall Zorn's lemma, known to be equivalent to the axiom of choice:

**Lemma 48 (Zorn's lemma).** *Any partially ordered set  $(X, \leq)$  such that any chain in  $X$  possesses an upper bound has a maximal element.*

A filter  $\mathcal{F}$  over  $I$  is an ultra-filter if and only if it is maximal with respect to set inclusion. By Zorn's lemma, any filter  $\mathcal{F}$  over  $I$  can be extended to an ultra-filter over  $I$ . Now, if  $I$  is infinite, the family of sets  $\mathcal{F} = \{P \subseteq I \mid P^c \text{ is finite}\}$  is a *free filter*, meaning it contains no finite set. It can thus be extended to a free ultra-filter over  $I$ . Hence:

**Lemma 49.** *Any infinite set has a free ultra-filter.*

Every free ultra-filter  $\mathcal{F}$  over  $I$  uniquely defines, by setting  $\mu(P) = 1$  if  $P \in \mathcal{F}$  and 0 otherwise, a finitely additive measure<sup>23</sup>  $\mu : 2^I \mapsto \{0, 1\}$ , which satisfies

$$\mu(I) = 1 \text{ and, if } P \text{ is finite, then } \mu(P) = 0. \tag{89}$$

<sup>22</sup>The existence of such a measure is non-trivial and is explained later.

<sup>23</sup>Observe that, as a consequence,  $\mu$  cannot be sigma-additive (in contrast to probability measures or Radon measures) in that it is *not* true that  $\mu(\bigcup_n A_n) = \sum_n \mu(A_n)$  holds for an infinite denumerable sequence  $A_n$  of pairwise disjoint subsets of  $\mathbb{N}$ .

Now, fix an infinite set  $I$  and a finitely additive measure  $\mu$  over  $I$  as above. Let  $\mathbb{X}$  be a set and consider the Cartesian product  $\mathbb{X}^I = (x_i)_{i \in I}$ . Define

$$(x_i) \sim (x'_i) \quad (90)$$

if and only if  $\mu\{i \in I \mid x_i \neq x'_i\} = 0$ . Relation  $\sim$  is an equivalence relation whose equivalence classes are denoted by  $[x_i]$  and we define

$${}^*\mathbb{X} = \mathbb{X}^I / \sim . \quad (91)$$

$\mathbb{X}$  is naturally embedded into  ${}^*\mathbb{X}$  by mapping every  $x \in \mathbb{X}$  to the constant tuple such that  $x_i = x$  for every  $i \in I$ ; we denote it by  $[x]$  or, simply, by  $x$ , if no confusion can result. Any algebraic structure over  $\mathbb{X}$  (group, ring, field) carries over to  ${}^*\mathbb{X}$  by almost pointwise extension. In particular, if  $[x_i] \neq 0$ , meaning that  $\mu\{i \mid x_i = 0\} = 0$  we can define its inverse  $[x_i]^{-1}$  by taking  $y_i = x_i^{-1}$  if  $x_i \neq 0$  and  $y_i = 0$  otherwise. This construction yields  $\mu\{i \mid y_i x_i = 1\} = 1$ , whence  $[y_i][x_i] = 1$  in  ${}^*\mathbb{X}$ . The existence of an inverse for any non-zero element of a ring is indeed stated by the formula:  $\forall x (x = 0 \vee \exists y (xy = 1))$ .<sup>24</sup>

### 7.3. Nonstandard reals and integers

The above general construction applies to  $\mathbb{X} = \mathbb{R}$  and  $I = \mathbb{N}$ . The result, denoted by  ${}^*\mathbb{R}$ , is a field (according to the transfer principle). By the same principle,  ${}^*\mathbb{R}$  is totally ordered by  $[u_n] \leq [v_n]$  iff  $\mu\{n \mid u_n > v_n\} = 0$ . Call *infinitesimal* any nonstandard real number whose absolute value is smaller than any positive real number:  $[|x_n|] \leq [\varepsilon]$  for any  $\varepsilon \in \mathbb{R}, \varepsilon > 0$ . For  $x$  and  $y$  two nonstandard real numbers, write

$$x \approx y$$

if  $x - y$  is infinitesimal. Call *infinite* any nonstandard real number whose absolute value is larger than any positive real number:  $[|x_n|] \geq [K]$  for any finite  $K \in \mathbb{R}$ . Call *finite* a nonstandard real number that is not infinite. Call *bounded* a subset  $A \subseteq {}^*\mathbb{R}$  such that there exists a finite standard positive real  $K$  such that  $|x| \leq K$  for every  $x \in A$ .

**Lemma 50 (standard part).** *For any finite  $[x_n] \in {}^*\mathbb{R}$ , there exists a unique standard real number  $x \in \mathbb{R}$ , such that  $[x] - [x_n]$  is infinitesimal. We call  $x$  the standard part of  $[x_n]$  and denote it by  $\text{st}([x_n])$ . Infinite nonstandard reals have no standard part in  $\mathbb{R}$ .*

*Proof.* To prove this, let  $x = \sup\{u \in \mathbb{R} \mid [u] \leq [x_n]\}$ . Since  $[x_n]$  is finite,  $x$  exists and we only need to show that  $[x_n] - x$  is infinitesimal. If not, then there exists  $y \in \mathbb{R}, y > 0$  such that either  $[x] < [x_n] - [y]$  or  $[x] > [x_n] + [y]$ , which both contradict the definition of  $x$ . The uniqueness of  $x$  is clear, thus we can define  $\text{st}([x_n]) = x$ .  $\square$

It is also of interest to apply the general construction (91) to  $\mathbb{X} = I = \mathbb{N}$ , which results in the set  ${}^*\mathbb{N}$  of *nonstandard natural numbers*;  ${}^*\mathbb{Z}$  is defined similarly, from  $\mathbb{X} = \mathbb{Z}$  and  $I = \mathbb{N}$ . The nonstandard set  ${}^*\mathbb{N}$  differs from  $\mathbb{N}$  by the addition of *infinite natural numbers*, which are equivalence classes of sequences of integers whose essential limit is  $+\infty$ .

<sup>24</sup>More generally, the *Transfer Principle* states that every first-order formula is true over  ${}^*\mathbb{X}$  if and only if it is true over  $\mathbb{X}$ . Recall that a first-order formula is a formula involving quantifiers over variables but not functions.

#### 7.4. Internal functions and sets

Any sequence  $(f_n)$  of functions  $f_n : \mathbb{R} \rightarrow \mathbb{R}$  pointwise defines a function  $[f_n] : {}^*\mathbb{R} \rightarrow {}^*\mathbb{R}$  by setting

$$[f_n]([x_n]) =_{\text{def}} [f_n(x_n)]. \quad (92)$$

To justify this definition based on a representant of the equivalence class, note that  $[x_n] = [x'_n]$  means that the set of integers  $n$  such that  $x_n \neq x'_n$  is neglectible, and we have  $f_n(x_n) = f_n(x'_n)$  for every  $n$  outside of this set. A function  ${}^*\mathbb{R} \rightarrow {}^*\mathbb{R}$  obtained in this way is called *internal*. Properties of, and operations on, ordinary functions extend pointwise to internal functions of  ${}^*\mathbb{R} \rightarrow {}^*\mathbb{R}$ . The same notions apply to sets by considering their characteristic functions.

An internal set  $A = [A_n]$  is called *hyperfinite* if  $\mu\{n \mid A_n \text{ finite}\} = 1$ ; the *cardinal*  $|A|$  of  $A$  is defined as  $[|A_n|]$ , where  $|A_n|$  denotes the cardinal of the finite set  $A_n$  in the usual sense. For  $A$  a bounded set of hyperreals, we define its *shadow*  $\text{Sh}(A) \subseteq \mathbb{R}$  as follows:

$$\text{Sh}(A) = \{\text{st}(x) \in \mathbb{R} \mid \exists y \in A \text{ such that } y \approx x\}. \quad (93)$$

The *nonstandard lifting* of  $f : \mathbb{R} \rightarrow \mathbb{R}$  is the internal function  ${}^*f = [f, f, f, \dots]$ . Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a standard function that is differentiable at  $x \in \mathbb{R}$ , and  ${}^*f = [f, f, f, \dots]$  its nonstandard lifting. Then, we have, for every  $x \in \mathbb{R}$ :

$$f'(x) = \text{st}\left(\frac{{}^*f(y) - {}^*f(x)}{y - x}\right) \quad (94)$$

for every  $y \in {}^*\mathbb{R}$  such that  $y \approx [x]$ —the proof is similar to that of Lemma 50.

Now, consider an infinite number  $K \in {}^*\mathbb{N}$  and the set

$$T = \left\{ 0, \frac{1}{K}, \frac{2}{K}, \frac{3}{K}, \dots, \frac{K-1}{K}, 1 \right\}. \quad (95)$$

By definition, if  $K = [K_n]$ , then  $T = [T_n]$  with

$$T_n = \left\{ 0, \frac{1}{K_n}, \frac{2}{K_n}, \frac{3}{K_n}, \dots, \frac{K_n-1}{K_n}, 1 \right\}$$

hence  $|T| = [|T_n|] = [K_n + 1] = K + 1$ . Note that the shadow of  $T$  (see (93)) is the standard interval  $[0, 1]$ : for any given real number  $a$  in  $[0, 1]$ ,  $\frac{[aK]}{K}$  standardizes to  $a$ .

In our nonstandard semantics, we will be using the following time basis:

$$\mathbb{T} =_{\text{def}} \{k\partial \mid k \in {}^*\mathbb{N}\} \quad (96)$$

where  $\partial > 0$  is some fixed infinitesimal time. This definition is reminiscent of Equation (13), in Section 2.3. The shadow of  $\mathbb{T}$  is  $\mathbb{R}_+$ , the set of nonnegative real numbers, and that, informally speaking,  $\mathbb{T}$  is both “discrete” (every element has a previous and a next element, except for 0) and “continuous” (as the shadow of  $\mathbb{T}$  is  $\mathbb{R}_+$ ).

**Notations 51.** Generic elements of  $\mathbb{T}$  will be denoted by the special symbol  $\mathfrak{t}$  or simply  $t$  when no confusion can result.

### 7.5. Integrals

Now, consider an internal function  $f = [f_n]$  and a hyperfinite set  $A = [A_n]$ . The *sum* of  $f$  over  $A$  can be defined by

$$\sum_{a \in A} f(a) =_{\text{def}} \left[ \sum_{a \in A_n} f_n(a) \right].$$

If  $T$  is as above, and  $f : \mathbb{R} \rightarrow \mathbb{R}$  is a standard function, we obtain

$$\sum_{t \in T} \frac{1}{|T|} {}^*f(t) = \left[ \sum_{t \in T_n} \frac{1}{|T_n|} f(t) \right]. \quad (97)$$

Now, the continuity of  $f$  implies the convergence of the Riemann sums:  $\sum_{t \in T_n} \frac{1}{|T_n|} f(t) \rightarrow \int_0^1 f(t) dt$ . Hence,

$$\int_0^1 f(t) dt = \text{st} \left( \sum_{t \in T} \frac{1}{|T|} {}^*f(t) \right). \quad (98)$$

Under the same assumptions, for any  $t \in [0, 1]$ ,

$$\int_0^t f(u) du = \text{st} \left( \sum_{u \in T, u \leq t} \frac{1}{|T|} {}^*f(u) \right). \quad (99)$$

### 7.6. ODE

Consider the following ODE:

$$x' = f(x, t), \quad x(0) = x_0. \quad (100)$$

Assume (100) possesses a solution  $x : [0, 1] \rightarrow \mathbb{R}$  such that the function  $t \mapsto f(x(t), t)$  is continuous. Rewriting (100) in its equivalent integral form  $x(t) = x_0 + \int_0^t f(x(u), u) du$  and using (99) yields

$$x(t) = \text{st}({}^*x(t)) \text{ where } {}^*x(t) =_{\text{def}} x_0 + \sum_{u \in T, u \leq t} \frac{1}{|T|} {}^*f(x(u), u). \quad (101)$$

One can rewrite the positive infinitesimal quantity  $1/|T|$  as  $\partial$ , so that

$$T = \{t_k = k\partial \mid k = 0, \dots, |T|\}.$$

Then, after substituting in (101), one gets that the piecewise-constant right-continuous function  ${}^*x(t), t \in {}^*\mathbb{R}, 0 \leq t \leq 1$  satisfies the following difference equation, for  $k=0, \dots, |T|-1$ ,

$$\begin{aligned} {}^*x(t_{k+1}) &= {}^*x(t_k) + \partial \cdot {}^*f({}^*x(t_k), t_k) \\ {}^*x(t_0) &= x_0 \end{aligned} \quad (102)$$

By (101), the following holds:

**Theorem 52.** *The solution  $x(t), t \in \mathbb{R}$ , of ODE (100), and  ${}^*x(t), t \in \mathbb{R}_+$ , the piecewise constant interpolation of the trajectory of system (102), are related by  $x = \text{st}({}^*x)$ .*

The same line of arguments applies for systems of ODEs possessing a unique continuous solution, i.e.,  $x$  and  $f$  take their values in  $\mathbb{R}^m$  in (100). Formula (102) can be seen as a *nonstandard semantics* for ODE (100). This semantics seems to depend on the choice of infinitesimal step parameter  $\partial$ . Property (101), though, expresses that

$$\text{all of these nonstandard semantics are equivalent from the standard viewpoint, regardless of the choice made for } \partial. \quad (103)$$

The paper [35] goes beyond our short exposure by including a direct proof of the Peano theorem on the existence of solutions of (100) when  $f$  is continuous and bounded.

### 7.7. Infinitesimal calculus

We conclude this background material with a summary of the section “Infinitesimal calculus” from [35]. We formulate the essential characterizations of properties of the infinitesimal calculus in nonstandard terms.  $f$  denotes a standard function and  ${}^*f = [f, f, f, \dots]$  is its nonstandard lifting. The reader is referred to that text for the proofs.

#### Theorem 53.

1. *The function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is continuous at  $a \in \mathbb{R}$  if and only if  ${}^*f(x) \approx f(a)$  for all  $x \approx a$ .*
2. *The function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is uniformly continuous on set  $A \subseteq \mathbb{R}$  if and only if  ${}^*f(x) \approx {}^*f(y)$  for all  $x, y \in {}^*A$  such that  $x \approx y$ .*
3. *The function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is differentiable at  $a \in \mathbb{R}$  if and only if there exists a number  $b \in \mathbb{R}$  such that*

$$\frac{{}^*f(x) - {}^*f(a)}{x - a} \approx b, \text{ for all } x \in \mathbb{R}, x \approx a, x \neq a.$$

*Moreover, if such a number  $b$  exists, then it equals  $f'(a)$ .*

## 8. The toolkit supporting standardization

By building on the background of Section 7, we develop here the specific material we need for our developments. First, we develop the compile-time impulse analysis. Second, we formally justify our use of structural analysis in the nonstandard domain. Third, referring to Comment 4 in support of the clutch example, we formally justify our method for standardizing systems of equations. We consider systems of equations of the form

$$0 = \mathbf{H}(X', X^\bullet, V, X), \quad (104)$$

where  $Z =_{\text{def}} (X^\bullet, V)$  collects the dependent variables. Such systems are also written as

$$0 = H(Z, X, \partial) \quad (105)$$

after expanding  $X'$  as  $\frac{X^\bullet - X}{\partial}$ . In System (104),  $V$  collects the algebraic variables,  $X$  collects the state variables, and  $\frac{X^\bullet - X}{\partial}$  is the nonstandard semantics of  $X'$ . Writing (104)



involves both  $X'$  and  $X^\bullet$  but not  $\partial$ . As (105) is obtained from (104) by performing the expansion, it involves both  $X^\bullet$  and  $\partial$  but not  $X'$ .

As an illustration, with reference to the Cup-and-Ball example, System (29) is of the form (104)—we repeat it here for convenience:

$$\left\{ \begin{array}{ll} 0 = x'' + \lambda x & (e_1) \\ 0 = y'' + \lambda y + g & (e_2) \\ \gamma^\bullet = [s \leq 0]; \gamma(0) = \text{F} & (k_0) \\ \text{if } \gamma \text{ then } 0 = L^2 - (x^2 + y^2) & (k_1) \\ \quad \text{and } 0 = L^2 - (x^2 + y^2)^\bullet & (k_1^\bullet) \\ \quad \text{and } 0 = L^2 - (x^2 + y^2)^{\bullet 2} & (k_1^{\bullet 2}) \\ \quad \text{and } 0 = \lambda + s & (k_2) \\ \text{if not } \gamma \text{ then } 0 = \lambda & (k_3) \\ \quad \text{and } 0 = (L^2 - (x^2 + y^2)) - s & (k_4) \end{array} \right.$$

Its counterpart following form (105) is obtained by expanding the second derivatives by using Euler scheme:

$$\left\{ \begin{array}{ll} 0 = \frac{x^{\bullet 2} - 2x^\bullet + x}{\partial^2} + \lambda x & (e_1) \\ 0 = \frac{y^{\bullet 2} - 2y^\bullet + y}{\partial^2} + \lambda y + g & (e_2) \\ \gamma^\bullet = [s \leq 0]; \gamma(0) = \text{F} & (k_0) \\ \text{if } \gamma \text{ then } 0 = L^2 - (x^2 + y^2) & (k_1) \\ \quad \text{and } 0 = L^2 - (x^2 + y^2)^\bullet & (k_1^\bullet) \\ \quad \text{and } 0 = L^2 - (x^2 + y^2)^{\bullet 2} & (k_1^{\bullet 2}) \\ \quad \text{and } 0 = \lambda + s & (k_2) \\ \text{if not } \gamma \text{ then } 0 = \lambda & (k_3) \\ \quad \text{and } 0 = (L^2 - (x^2 + y^2)) - s & (k_4) \end{array} \right.$$

**Comment 54.** System (104) is the generic form of systems we need to solve at Line 10 of Algorithm 13, when we eliminate higher order derivatives and shifts by introducing auxiliary variables.  $\square$

We consider the following assumptions, which hold for the nonstandard semantics of any standard mDAE system:

**Assumption 6.**

1. The functions  $\mathbf{H}(\cdot)$  and  $H(\cdot)$  are standard (they result from applying Algorithm 13 to our given standard mDAE model);
2. The state variables  $X$  are standard and finite (their values were assigned by the standardized dynamics of the previous mode).

8.1. Impulse Analysis

We refer the reader to Section 2.5 about the standardization of the clutch model at mode changes, and particularly System (24) in it. As a prerequisite to standardization, we had to identify which variables could become impulsive at mode changes (the two torques in this example), and then, we had to eliminate them.

To identify impulsive and non-impulsive variables in a systematic way, we develop the *impulse analysis*, which consists in abstracting hyperreals with their “magnitude order” compared to the infinitesimal  $\partial$ . Since our original DAE models are all standard, only the occurrence of the infinitesimal  $\partial$  will have a role in this matter. The impulse analysis will be useful as a preparation step for standardization.

**Definition 55 (impulse order and analysis).**

1. Given a system of equations such as (104) or (105), say that a dependent variable  $x$  has impulse order (or simply order)  $\mathfrak{o} \in \mathbb{R}$ , if the solution of the considered system is such that  $x\partial^{\mathfrak{o}}$  is provably a finite non-zero (standard) real number. We denote by  $\llbracket x \rrbracket$  the impulse order of  $x$ . By convention, the constant 0 has impulse order  $-\infty$ .
2. Say that  $x$  is impulsive if  $\llbracket x \rrbracket > 0$ .
3. The impulse analysis of a system of equations such as (104) or (105) is the system of constraints satisfied by the impulse orders of the dependent variables of the system.

8.1.1. The rules of impulse analysis

Figures 12 and 13 display the rules defining the translation of a system of equations of the form (104) or (105) into its impulse analysis, for the restricted class where only rational expressions are involved.

$$\begin{aligned} e & ::= 0 \mid c \mid \partial \mid x \mid e^c \mid e + e \mid e \times e \\ E & ::= e = e \mid E \text{ and } E \end{aligned}$$

Figure 12: Syntax:  $E$  is a system of one or several equations  $e = e$ . An expression  $e$  is 0, a nonzero (standard) real constant  $c$ , the infinitesimal  $\partial$ , a variable  $x$ , the monomial  $e^c$ , a sum, or a product.

$$\begin{array}{ll} \text{(R1)} & \llbracket 0 \rrbracket = -\infty \\ \text{(R2)} & \llbracket c \rrbracket = 0 \\ \text{(R3)} & \llbracket \partial \rrbracket = -1 \\ \text{(R4)} & \llbracket e^c \rrbracket = c \llbracket e \rrbracket \\ \text{(R5)} & \llbracket e_1 \times e_2 \rrbracket = \llbracket e_1 \rrbracket + \llbracket e_2 \rrbracket \\ \text{(R6)} & \llbracket e_1 + e_2 \rrbracket \leq \max\{\llbracket e_1 \rrbracket, \llbracket e_2 \rrbracket\} \end{array} \quad \frac{\begin{array}{l} E \vdash e = e' \\ \llbracket E \rrbracket \vdash \llbracket e \rrbracket = \llbracket e' \rrbracket \end{array}}{\text{(R7)}} \quad \frac{\begin{array}{l} E \vdash x = y + e \text{ or } \\ E \vdash 0 = y - x + e \end{array} \text{ and } E \not\vdash y = x - e}{E \vdash E \text{ and } y = x - e} \quad \text{(R8)}$$

Figure 13: Rules: The left column displays the impulse order of the primitive expressions. Rule (R7) indicates that  $\llbracket e \rrbracket = \llbracket e' \rrbracket$  is an equation of the impulse analysis  $\llbracket E \rrbracket$  if  $e = e'$  is an equation of  $E$ ; Rule (R8) indicates that, if  $E$  involves the equation  $x = y + e$  but not the equation  $y = x - e$ , then we augment  $E$  with the latter, i.e., we saturate  $E$  with the rule  $x = y + e \implies y = x - e$ .

Figure 12 describes the syntax of a mini-language specifying such systems of equations. The left column of Figure 13 gives the rules for mapping expressions to their corresponding impulse orders. All the rules are self-explanatory, with the exception of Rule (R6), where the inequality deserves some explanation. The sum  $e_1 + e_2$ , the dominant terms in the expansion of the  $e_i$ 's as power series over  $\partial$  may compensate each other for their impulse orders. For an example of this, see equation  $(e_1^\partial)$  in System (24): rewriting this equation

as  $a_1(\omega_1) = \frac{\omega_1^{\bullet} - \omega_1}{\partial} - b_1(\omega_1)\tau_1$ , we see a case of strict inequality for (R6) since  $a_1(\omega_1)$  has order zero, whereas it is equal to the difference of two terms of order one, see Section 8.1.4 for details.

We will use Rule (R6) in the following way, thereby reinforcing it. Consider an equation

$$e : z = x + y.$$

We can rewrite  $e$  in the following equivalent ways:  $0 = x + y - z$ ,  $x = z - y$ , or  $y = z - x$ . To each of them we apply the max rule. This yields, for the impulse analysis of equation  $e$ , the following system of constraints:

$$\text{impulse analysis of } e : \begin{cases} \llbracket z \rrbracket \leq \max\{\llbracket x \rrbracket, \llbracket y \rrbracket\} & ; & \llbracket 0 \rrbracket \leq \max\{\llbracket x \rrbracket, \llbracket y \rrbracket, \llbracket z \rrbracket\} \\ \llbracket x \rrbracket \leq \max\{\llbracket z \rrbracket, \llbracket y \rrbracket\} & ; & \llbracket y \rrbracket \leq \max\{\llbracket x \rrbracket, \llbracket z \rrbracket\} \end{cases} \quad (106)$$

Note that the constraint  $\llbracket 0 \rrbracket \leq \dots$  is vacuously satisfied since  $\llbracket 0 \rrbracket = -\infty$ . Also note that, among the three nontrivial inequalities of (106), at least two must be saturated. We will use impulse analysis (106) for handling sums of terms. This reinforcement of the max rule is formalized by Rule (R8) of Figure 13, which mechanizes the association, to equation  $e$  of (106), of its different rewritings.

Using the rules of Figures 12 and 13 in the numerical expressions, we map any system of rational equations of the form (104) or (105) into a system of constraints over impulse orders.

To cover functions beyond polynomials, we need to extend  $\mathbb{R} \cup \{-\infty\}$  with  $+\infty$ . In this extension, we take the convention that  $-\infty + \infty = -\infty$ , justified by the equality  $0 \times x = 0$  for any nonstandard  $x$ . For functions  $f(x) = \sum_{k=0}^{\infty} a_k x^k$  that can be represented as absolutely converging power series, we then get

$$\llbracket f(x) \rrbracket = \llbracket \sum_{k=0}^{\infty} a_k x^k \rrbracket = \llbracket x \rrbracket \cdot \text{sup}(A), \text{ where } A = \{k \mid a_k \neq 0\} \quad (107)$$

is the support of the series and  $\text{sup}(A)$  is the supremum of set  $A$ . In particular, if  $\llbracket x \rrbracket > 0$  and if the support of the series is infinite, we get  $\llbracket f(x) \rrbracket = +\infty$ .

### 8.1.2. Impulse analysis of systems of equations for restarts

Here we particularize the impulse analysis to systems of equations of the form (104), where the only reason for  $\partial$  to occur is the expansion of derivatives using the Euler scheme:

$$0 = \mathbf{H} \left( \frac{X^{\bullet} - X}{\partial}, X^{\bullet}, V, X \right)$$

The dependent variables are  $X^{\bullet}, V$ . It will be convenient to introduce the auxiliary variables

$$U =_{\text{def}} X^{\bullet} - X,$$

so that the systems we consider take the following form:

$$\text{dependent variables } X^{\bullet}, V, U \text{ in } : \begin{cases} 0 & = & \mathbf{H} \left( \frac{U}{\partial}, X^{\bullet}, V, X \right) \\ U & = & X^{\bullet} - X \end{cases} \quad (108)$$

The following condition for System (108) can be assumed, based on physical considerations (restart values for an ODE or a DAE cannot be impulsive in a physically meaningful model):

**Assumption 7.** *Since  $X$  is a state, both  $X$  (a known value) and  $X^\bullet$  must be finite.*

First, the impulse orders  $\llbracket X \rrbracket$  are all known, from previous nonstandard instants. Next, from Assumption 7 we deduce the inequalities

$$\llbracket X^\bullet \rrbracket \leq 0 \quad \text{and} \quad \llbracket U \rrbracket \leq 0. \quad (109)$$

The impulse orders  $\llbracket V \rrbracket$  are a priori unknown. We have, however, more prior information, thanks to the structural analysis. As part of both the simpler algorithm ExecRun (Algorithm 10) and revisited algorithm ExecRun (Algorithm 13), the auxiliary algorithm SolveConflict (Algorithm 9) is called. At the considered mode change, we thus know which consistency equation(s) of the new mode was/were conflicting with the dynamics of the previous mode. Formally, call  $G=0$  the subsystem collecting all the equations that were erased by algorithm SolveConflict at Line 6 of Algorithm 10 or Line 7 of Algorithm 13. As a result,  $G=0$  is no longer satisfied at the considered mode change, and thus,  $G$  defines a tuple  $R$  of finite nonzero variables called *residuals*, by setting

$$R = G. \quad (110)$$

An example of residual in the clutch is  $r = \omega_1 - \omega_2$ , which is both finite and nonzero at mode change  $\gamma : F \rightarrow T$ ; this information was found by the structural analysis. Finally, the system of equations that we need to solve collects all the above items, namely:

$$\text{dependent variables } X^\bullet, V, U, R \text{ in} \quad : \quad \begin{cases} 0 & = \mathbf{H}(\frac{U}{\partial}, X^\bullet, V, X) \\ U & = X^\bullet - X \\ R & = G \end{cases} \quad (111)$$

$$\text{prior information on impulse orders} \quad : \quad \begin{cases} \llbracket \frac{1}{\partial} \rrbracket & = 1 \\ \llbracket X^\bullet \rrbracket & \leq 0 \\ \llbracket U \rrbracket & \leq 0 \\ \llbracket R \rrbracket & = 0 \end{cases} \quad (112)$$

*Decomposing algebraic dependent variables into impulsive and non-impulsive ones.* In the following, we assume that the vector functions  $\mathbf{H}$  or  $H$  can be represented as absolutely convergent power series in their arguments. We can then apply the rules of Figures 12 and 13, complemented with (107), to derive the system of equations that impulse orders must satisfy. In particular, this allows us to partition the set  $V$  of algebraic variables as

$$V = W^i \cup W \quad (113)$$

where  $W^i$  collects the *impulsive variables*, having impulse order  $> 0$ , and  $W$  collects the *non-impulsive variables*. In Section 9.3, we will propose a numerical scheme for computing restarts that only requires the splitting of algebraic dependent variables into impulsive and non-impulsive ones. For now, let us develop the impulse analysis for the two transitions  $\gamma : T \rightarrow F$  and  $\gamma : F \rightarrow T$  of System (15) assuming linear  $f_i$  as in Eq. (23).

8.1.3. *Example: mode change  $\gamma : \mathbb{T} \rightarrow \mathbb{F}$  of the clutch*

We develop this example here based on the rules of Figures 12 and 13. To better illustrate this mechanical reasoning, we forbid ourselves algebraic manipulations that otherwise could help simplifying a manual analysis. We repeat here the corresponding system of equations with its expansion (108)—there is no residual for this mode change:

$$\left\{ \begin{array}{ll} \frac{u_1}{\partial} = a_1(\omega_1) + b_1(\omega_1)\tau_1 & (e_1^\partial) \\ \frac{u_2}{\partial} = a_2(\omega_2) + b_2(\omega_2)\tau_2 & (e_2^\partial) \\ \tau_1 = 0 & (e_5) \\ \tau_2 = 0 & (e_6) \\ u_1 = \omega_1^\bullet - \omega_1 \\ u_2 = \omega_2^\bullet - \omega_2 \end{array} \right. \quad (114)$$

We then saturate System (114) using Rule (R8) of Figure 13 (added equations are in blue):

$$\left\{ \begin{array}{l} \frac{u_1}{\partial} = a_1(\omega_1) + b_1(\omega_1)\tau_1 ; a_1(\omega_1) = \frac{u_1}{\partial} - b_1(\omega_1)\tau_1 ; b_1(\omega_1)\tau_1 = \frac{u_1}{\partial} - a_1(\omega_1) \\ \frac{u_2}{\partial} = a_2(\omega_2) + b_2(\omega_2)\tau_2 ; a_2(\omega_2) = \frac{u_2}{\partial} - b_2(\omega_2)\tau_2 ; b_2(\omega_2)\tau_2 = \frac{u_2}{\partial} - a_2(\omega_2) \\ \tau_1 = 0 \\ \tau_2 = 0 \\ u_1 = \omega_1^\bullet - \omega_1 ; u_1 + \omega_1 = \omega_1^\bullet ; u_1 - \omega_1^\bullet = -\omega_1 \\ u_2 = \omega_2^\bullet - \omega_2 ; u_2 + \omega_2 = \omega_2^\bullet ; u_2 - \omega_2^\bullet = -\omega_2 \end{array} \right. \quad (115)$$

and we consider the prior information (112):

$$\llbracket \omega_i^\bullet \rrbracket \leq 0 ; \llbracket u_i \rrbracket \leq 0 ; \llbracket \frac{1}{\partial} \rrbracket = 1 \quad (116)$$

We now apply the rules of Figure 13 to System (115), then use (116):

$$\left\{ \begin{array}{l} 1 + \llbracket u_1 \rrbracket \leq 0 ; 0 \leq 1 + \llbracket u_1 \rrbracket ; \text{vacuous} \\ 1 + \llbracket u_2 \rrbracket \leq 0 ; 0 \leq 1 + \llbracket u_2 \rrbracket ; \text{vacuous} \\ \llbracket \tau_1 \rrbracket = -\infty \\ \llbracket \tau_2 \rrbracket = -\infty \\ \llbracket u_1 \rrbracket \leq 0 ; \text{vacuous} ; \text{vacuous} \\ \llbracket u_2 \rrbracket \leq 0 ; \text{vacuous} ; \text{vacuous} \end{array} \right. \quad (117)$$

where “vacuous” indicates that the equation sitting in the corresponding position in System (115) generates a vacuously satisfied constraint. System (117) shows that  $\llbracket u_i \rrbracket = -1$ , meaning that  $u_i$  is continuous. Note that System (117) uniquely determines the impulse orders for all dependent variables.

8.1.4. *Example: mode change  $\gamma : \mathbb{F} \rightarrow \mathbb{T}$  of the clutch*

We repeat here the corresponding system of equations with its expansion (108), and we include the residual  $r \stackrel{\text{def}}{=} \omega_1 - \omega_2$  (from structural analysis, we know it is nonzero),

highlighted in red:

$$\left\{ \begin{array}{ll} \frac{u_1}{\partial} = a_1(\omega_1) + b_1(\omega_1)\tau_1 & (e_1^\partial) \\ \frac{u_2}{\partial} = a_2(\omega_2) + b_2(\omega_2)\tau_2 & (e_2^\partial) \\ \omega_1^\bullet - \omega_2^\bullet = 0 & (e_3^\bullet) \\ \tau_1 + \tau_2 = 0 & (e_4) \\ u_1 = \omega_1^\bullet - \omega_1 \\ u_2 = \omega_2^\bullet - \omega_2 \\ r = \omega_1 - \omega_2 \end{array} \right. \quad (118)$$

We then saturate System (118) using Rule (R8) of Figure 13 (added equations are in blue):

$$\left\{ \begin{array}{l} \frac{u_1}{\partial} = a_1(\omega_1) + b_1(\omega_1)\tau_1 ; a_1(\omega_1) = \frac{u_1}{\partial} - b_1(\omega_1)\tau_1 ; b_1(\omega_1)\tau_1 = \frac{u_1}{\partial} - a_1(\omega_1) \\ \frac{u_2}{\partial} = a_2(\omega_2) + b_2(\omega_2)\tau_2 ; a_2(\omega_2) = \frac{u_2}{\partial} - b_2(\omega_2)\tau_2 ; b_2(\omega_2)\tau_2 = \frac{u_2}{\partial} - a_2(\omega_2) \\ \omega_1^\bullet - \omega_2^\bullet = 0 ; \omega_1^\bullet = \omega_2^\bullet \\ \tau_1 + \tau_2 = 0 ; \tau_1 = -\tau_2 \\ u_1 = \omega_1^\bullet - \omega_1 ; u_1 + \omega_1 = \omega_1^\bullet ; u_1 - \omega_1^\bullet = -\omega_1 \\ u_2 = \omega_2^\bullet - \omega_2 ; u_2 + \omega_2 = \omega_2^\bullet ; u_2 - \omega_2^\bullet = -\omega_2 \\ r = \omega_1 - \omega_2 \end{array} \right. \quad (119)$$

We also consider the prior information (112):

$$\llbracket \omega_i^\bullet \rrbracket \leq 0 ; \llbracket u_i \rrbracket \leq 0 ; \llbracket r \rrbracket = 0 ; \llbracket \frac{1}{\partial} \rrbracket = 1 . \quad (120)$$

We now apply the rules of Figure 13 to System (119), then use (120):

$$\left\{ \begin{array}{l} 1 \leq \llbracket \tau_1 \rrbracket ; \text{vacuous} ; \llbracket \tau_1 \rrbracket \leq 1 \\ 1 \leq \llbracket \tau_2 \rrbracket ; \text{vacuous} ; \llbracket \tau_2 \rrbracket \leq 1 \\ \text{vacuous} ; \llbracket \omega_1^\bullet \rrbracket = \llbracket \omega_2^\bullet \rrbracket \\ \text{vacuous} ; \llbracket \tau_1 \rrbracket = \llbracket \tau_2 \rrbracket \\ \text{vacuous} ; \text{vacuous} ; 0 \leq \max\{\llbracket u_1 \rrbracket, \llbracket \omega_1^\bullet \rrbracket\} \\ \text{vacuous} ; \text{vacuous} ; 0 \leq \max\{\llbracket u_2 \rrbracket, \llbracket \omega_2^\bullet \rrbracket\} \\ 0 = \llbracket r \rrbracket = \llbracket \omega_1 \rrbracket = \llbracket \omega_2 \rrbracket \\ \llbracket \omega_i^\bullet \rrbracket \leq 0 , \llbracket u_i \rrbracket \leq 0 , i = 1, 2 \end{array} \right. \quad (121)$$

which proves  $\llbracket \tau_1 \rrbracket = \llbracket \tau_2 \rrbracket = 1$ , expressing that the two torques are impulsive. Note, on the other hand, that (121) does not fully determine the impulse orders of  $\omega_i^\bullet$  and  $u_i$ , since one may, e.g., have  $\llbracket \omega_i^\bullet \rrbracket < 0$  and  $\llbracket u_i \rrbracket = 0$ .

#### 8.1.5. Using impulse analysis in code generation

Code generation for restarts consists in standardizing System (104) or System (105). Recall that standardizing systems of equations requires more care than standardizing numbers, due to impulsive behaviors and singularity issues that result.

We can exploit the impulse analysis through three different approaches, which we illustrate by using the clutch example and particularly the mode change  $\gamma : \text{F} \rightarrow \text{T}$  when the clutch gets engaged. The system of equations for the corresponding restart is (24), which we recall for convenience:

$$\left\{ \begin{array}{ll} \omega_1^\bullet = \omega_1 + \partial.(a_1(\omega_1) + b_1(\omega_1)\tau_1) & (e_1^\partial) \\ \omega_2^\bullet = \omega_2 + \partial.(a_2(\omega_2) + b_2(\omega_2)\tau_2) & (e_2^\partial) \\ \omega_1^\bullet - \omega_2^\bullet = 0 & (e_3^\bullet) \\ \tau_1 + \tau_2 = 0 & (e_4) \end{array} \right.$$

As we have shown, the two torques  $\tau_1$  and  $\tau_2$  are impulsive of order 1.

*Eliminating impulsive variables.* When this is practical, the simplest method is to eliminate impulsive variables from the restart system, namely the two torques in System (24). The details were given in Section 2.5.2 and the resulting reduced system of equations is

$$\omega_1^\bullet = \omega_2^\bullet = \frac{b_2(\omega_2)\omega_1 + b_1(\omega_1)\omega_2}{b_1(\omega_1) + b_2(\omega_2)} + \partial \frac{a_1(\omega_1)b_2(\omega_2) + a_2(\omega_2)b_1(\omega_1)}{b_1(\omega_1) + b_2(\omega_2)}, \quad (122)$$

whose standardization is simply achieved by substituting  $\partial \leftarrow 0$ . No numerical information is provided on the torques with this method.

This is a satisfactory solution when elimination of impulsive variables is practical. In our example, they entered linearly in the restart system, so that elimination was straightforward. When this is not the case, elimination becomes costly or even impossible. We thus need to look for alternatives.

*Rescaling impulsive variables.* Since  $\llbracket \tau_i \rrbracket = 1$ , we apply the change of variable  $\tau_i = \partial \hat{\tau}_i$ , which transforms System (118) into:

$$\begin{cases} \omega_1^\bullet - \omega_1 = a_1(\omega_1) + b_1(\omega_1)\hat{\tau}_1 & (e_1^\partial) \\ \omega_2^\bullet - \omega_2 = a_2(\omega_2) + b_2(\omega_2)\hat{\tau}_2 & (e_2^\partial) \\ \omega_1^\bullet - \omega_2^\bullet = 0 & (e_3^\bullet) \\ \hat{\tau}_1 + \hat{\tau}_2 = 0 & (e_4) \end{cases} \quad (123)$$

which yields again the solution (122) for  $\omega_1^\bullet = \omega_2^\bullet$ . However, the rescaled value  $\hat{\tau}_1$  can now be computed. Rescaling impulsive variables is simpler than eliminating them. Unfortunately, this method does not work in full generality since impulse orders can be infinite as we have shown in (107). The last method addresses such cases, at the price of a possibly poor numerical conditioning.

*Bruteforce solving of the restart system.* Consider once more System (24), replace in it the infinitesimal time step  $\partial$  by a small real positive time step  $\delta > 0$ , and expand the derivatives as before. This yields

$$\begin{cases} \omega_1^\bullet = \omega_1 + \delta.(a_1(\omega_1) + b_1(\omega_1)\tau_1) & (e_1^\delta) \\ \omega_2^\bullet = \omega_2 + \delta.(a_2(\omega_2) + b_2(\omega_2)\tau_2) & (e_2^\delta) \\ \omega_1^\bullet - \omega_2^\bullet = 0 & (e_3^\bullet) \\ \tau_1 + \tau_2 = 0 & (e_4) \end{cases} \quad (124)$$

Then, it will be proved in Theorem 82 of Section 9.3 that *solving System (124) for its dependent variables and then discarding the values found for the impulsive variables yields a converging approximation for the states  $\omega_1^+$  and  $\omega_2^+$  at restart*. Of course, numerical conditioning for System (124) is likely to be less good than for System (123), so that rescaling is recommended when impulse orders are finite.

*Qualitative information obtained from impulse analysis.* The dynamical profile of impulsive variables yields interesting information about the nature of the impulsive behaviors. For the clutch example, we get the following profile for  $\llbracket \tau_i \rrbracket : \dots, 0, 1, 0, \dots$ , where the 1 occurs at the event  $\gamma : \mathbb{F} \rightarrow \mathbb{T}$ . This indicates a Dirac-like impulse. As a second example

(not developed here), the impulse analysis of the Cup-and-Ball example with inelastic impact yields the following profile for the impulse order of the tension  $[\lambda] : \dots, 0, 1, 0, \dots$ , where the 1 occurs at the event following the mode change  $\gamma : \mathbb{F} \rightarrow \mathbb{T}$ , i.e., when System (33) gets executed to evaluate restart values for velocities. This again indicates a Dirac impulse for the tension. Note that having a profile  $\dots, 0, k, 0, \dots$  for an impulsive variable, where  $k$  is a positive integer, still indicates a Dirac behavior. The amplitude does not matter: only the profile matters. Developing a systematic analysis of this kind requires further study and is left for future investigations.

## 8.2. Nonstandard Structural Analysis

The motivation for this section was developed in Section 2.5. The following two questions will be addressed:

1. We apply structural analysis to the nonstandard semantics. Is this justified? In the standard setting, the structural analysis was motivated by the quest for low cost criteria for almost everywhere nonsingularity of systems of algebraic equations. Can we have a similar argument here?
2. We need to formalize and generalize the Comment 4 of Section 2.5.2, where we discussed the difference between standardizing functions and standardizing equations.

We begin with question 1. Algorithm 13 returns at its Line 10 systems of algebraic equations of the form  $H(Z, X, \partial) = 0$  introduced in (105), where tuple  $Z$  collects the dependent variables and tuple  $X$  collects the state variables. We justify the call of the atomic action  $\text{Eval}(\sigma, H)$  by the fact that “ $H$  is structurally nonsingular”. But the justification of structural analysis was developed in Section 5.1.1, i.e., in a standard setting. We thus need a justification for taking the liberty of using it for nonstandard semantics as well. Formally, we consider nonstandard systems of equations of the form

$${}^*H(Z, {}^*X, \partial) = 0, \text{ where } {}^*X \in {}^*\mathbb{R}^m, Z \text{ collects the dependent (nonstandard) variables, } {}^*H = [H, H, H, \dots] \text{ is the nonstandard lifting of the standard function } H, \text{ and } \partial = [\delta_n] \text{ is infinitesimal.} \quad (125)$$

**Definition 56.** *Say that  ${}^*H(Z, {}^*X, \partial) = 0$  is structurally nonsingular if and only if so is the standard equation  $H(Z, X, \delta) = 0$  in the sense of Section 5.1.1, where  $\delta > 0$  is a standard step size.*

**Lemma 57.** *Lemma 11 extends to nonstandard systems of the form (125).*

*Proof.* Select  ${}^*z = [z_n]$  and  ${}^*x = [x_n]$  satisfying  ${}^*H({}^*z, {}^*x, \partial) = 0$ . By definition, we have  $[H(z_n, x_n, \delta_n)] \sim [0, 0, \dots]$ , where relation  $\sim$  was defined in (90). Equivalently,

$$\mu(A) = 1, \text{ where } A = \{n \mid H(z_n, x_n, \delta_n) = 0\} \quad (126)$$

and  $\mu$  is the finitely additive measure over  $\mathbb{N}$  introduced in (89). Now, assume that  ${}^*H$  is structurally nonsingular according to Definition 56. Then, using notations (51) of Section 5.1.1, for every  $n \in A$ , there exist  $\mathbf{U}_n =_{\text{def}} U_{H_n}(z_n, x_n, \delta_n)$  and  $\mathbf{V}_n =_{\text{def}} V_{H_n}(z_n, x_n, \delta_n)$  satisfying property (50) with respect to the system  $H(Z_n, X_n, \delta_n) = 0$ .



Define the internal sets  $\mathbf{U} = [\mathbf{U}_n]$  and  $\mathbf{V} = [\mathbf{V}_n]$ . We have  ${}^*\lambda(\mathbf{V} \setminus \mathbf{U}) = 0$ , where  ${}^*\lambda$  is the internal set function  ${}^*\lambda = [\lambda, \lambda, \lambda, \dots]$ . For every  $n \in A$ , every  $\mathcal{J}_n \in \mathbf{V}_n$  yields an invertible matrix, implying that  $\mathcal{J} =_{\text{def}} [\mathcal{J}_n]$  also yields an invertible matrix. As a result, we have exhibited two sets  $\mathbf{U}$  and  $\mathbf{V}$  satisfying property (50) with respect to the system  ${}^*H(Z, {}^*X, \partial) = 0$ .  $\square$

Definition 56 thus formally justifies our shameless use of the structural analysis in nonstandard domains. Note that Definition 56 tells nothing about how  $\mathbf{U}$  is—it could be infinitely small. Neither does it require uniqueness of the solutions of the considered equations.

We now move to question 2 and investigate how equations should be standardized.

**Theorem 58 (standardizing equations).** *Consider a square system of the form (125), and assume that  $H$  is of class  $\mathcal{C}^1$ . The following two properties hold:*

1. *If the standard pair  $(z, x)$  satisfies  $H(z, x, 0) = 0$  and the Jacobian  $\mathbf{J}_Z H$  at the point  $(z, x, 0)$  is nonsingular, then, for every nonstandard pair  $({}^*x, \partial)$  such that  $({}^*x, \partial) \approx (x, 0)$ , there exists  ${}^*z$  such that  $\text{st}({}^*z) = z$  and  ${}^*H({}^*z, {}^*x, \partial) = 0$ .*
2. *Vice-versa, let  $({}^*z, {}^*x, \partial)$  be a finite nonstandard triple satisfying  ${}^*H({}^*z, {}^*x, \partial) = 0$ , and let  $(z, x, 0) = \text{st}({}^*z, {}^*x, \partial)$  be its standard part. Then  $H(z, x, 0) = 0$  holds.*

*Under these assumptions, the standardization of the equation  ${}^*H({}^*z, {}^*x, \partial) = 0$  is independent from the particular value for  $\partial$  we have chosen (provided that it is infinitesimal).*

**Corollary 59.** *We can standardize  ${}^*H(Z, X, \partial) = 0$  as the system  $H(Z, X, 0) = 0$ , provided that the latter is structurally nonsingular.*

The nonsingularity condition is needed for statement 1 to hold. We proceed to the proof of Theorem 58.

*Proof.* We successively prove the two statements. For Statement 1, by the Implicit Function Theorem applied to the standard function  $H$ , there exists a  $\mathcal{C}^1$ -class function  $(x', \delta) \mapsto G(x', \delta)$  such that  $H(G(x', \delta), x', \delta) = 0$  holds for  $(x', \delta)$  ranging over a neighborhood  $V$  of  $(x, 0)$ . Let  $({}^*x, \partial) = ([x_n], [\delta_n])$  be such that  $\text{st}({}^*x, \partial) = (x, 0)$ . Then, by definition of the standard part, there exists  $B \subseteq \mathbb{N}$  with  $\mu(B) = 1$  such that, for every  $n \in B$ ,  $(x_n, \delta_n) \in V$ . The sequence  $(G(x_n, \delta_n))_{n \in B}$  satisfies

$$H(G(x_n, \delta_n), x_n, \delta_n) = 0, \text{ for every } n \in B. \quad (127)$$

We extend the sequence  $(G(x_n, \delta_n))_{n \in B}$  to  $\mathbb{N}$  by assigning arbitrary values to indices  $n \in \mathbb{N} \setminus B$  and consider the nonstandard  ${}^*z =_{\text{def}} [G(x_n, \delta_n)]$ . By (127), we get  ${}^*H({}^*z, {}^*x, \partial) = 0$ .

For Statement 2, let the finite triple  $({}^*z, {}^*x, \partial)$  satisfy  ${}^*H({}^*z, {}^*x, \partial) = 0$  and let  $(z, x, 0) = \text{st}({}^*z, {}^*x, \partial)$ . Since  $(z, x, 0) \approx ({}^*z, {}^*x, \partial)$  and  $H$  is of class  $\mathcal{C}^1$ , we can apply Statement 2 of Theorem 53 to derive  $H(z, x, 0) \approx {}^*H({}^*z, {}^*x, \partial) = 0$ , whence  $H(z, x, 0) = 0$  follows since  $H(z, x, 0)$  is standard.  $\square$

### 8.3. DAE of index zero

Invoking the implicit function theorem, one can also address DAEs of index 0, i.e., of the form:

$$F(X', X) = 0 \quad (128)$$

where the Jacobian  $\nabla_V F(V, X)$  is invertible in a neighborhood of the solution of the equation  $F(V, X) = 0$ , where  $V$  is the vector of dependent variables.

**Theorem 60.** *Let  $F(X', X) = 0$  be a structurally nonsingular DAE system of index 0, and let  ${}^*F = [F, F, F, \dots]$  be the nonstandard lifting of  $F$ . Consider the following nonstandard transition system having  $X^\bullet$  as dependent variables:*

$$0 = {}^*F \left( \frac{X^\bullet - X}{\partial}, X \right), X(0) = X_0 \quad (129)$$

*Then, system (129) possesses a unique solution and we denote by  ${}^*X(t), t \in \mathbb{T}_+$  its piecewise constant interpolation. Furthermore, the solution  $X$  of DAE system  $F(X', X) = 0, X(0) = X_0$  satisfies  $X = \text{st}({}^*X)$ .*

*Proof.* Recall that the time basis of nonstandard semantics is  $\mathbb{T} = \{k \times \partial \mid k \in \mathbb{N}\}$ . Denoting by  $t_k$  the generic element of  $\mathbb{T}$ , transition system (129) rewrites:

$${}^*F \left( \frac{X(t_{k+1}) - X(t_k)}{\partial}, X(t_k) \right) = 0. \quad (130)$$

Using the Implicit Function Theorem, let  $G$  be the standard function such that  $X' = G(X)$  solves (128) for  $X'$ , and set  ${}^*G = [G, G, G, \dots]$ . Then, (130) rewrites

$$X(t_{k+1}) = X(t_k) + \partial \cdot {}^*G(X(t_k)), X(0) = X_0. \quad (131)$$

Theorem 60 follows from Theorem 52 applied to (131).  $\square$

In the next two sections we investigate the standardization of the blocks returned by Algorithm 13, as identified in Section 6.3.5. Since, in this section, we pay attention to the status standard/nonstandard, we will make this status explicit in our development whenever needed.

### 8.4. Long modes of mDAE

We use the notations and results of Section 6.3.5. More precisely, we reuse notations (80) with  $\prime$  replacing  $\bullet$ , and we invoke Lemma 40. We consider a system  $B_\prime(H, N) = (B_\prime^\dagger(H, N), B_\prime^h(H, N))$  returned by Algorithm 13 for long modes. By construction of Algorithm 13, the algebraic system  $G_\Sigma = 0$ , where  $G_\Sigma =_{\text{def}} B_\prime^\dagger(H, N)$ , is structurally nonsingular and admits  $\tilde{G}_\Sigma = 0$ , where  $\tilde{G}_\Sigma =_{\text{def}} B_\prime^h(H, N)$ , as associated consistency conditions.

At this point, we need to relate the above blocks to their standard continuous-time counterparts. The reduced index systems  $(G_\Sigma, \tilde{G}_\Sigma)$  returned by Algorithm 8 (IndexReduc) in case of success, coincide with the two systems returned by the  $\Sigma$ -method, when replacing differentiation by shifting (see the emphasized sentence when introducing this algorithm).

In our study of standardization, we will also need to consider the original outcome of the  $\Sigma$ -method (using differentiation); we denote this pair by  $(G_\Sigma^c, \bar{G}_\Sigma^c)$ , where the superscript  $c$  refers to continuous time. Continuous-time systems  $(G_\Sigma^c, \bar{G}_\Sigma^c)$  define the leading equations and consistency conditions of a (standard) DAE system that is generically regular. In the following theorem, we denote by  $X = (X(t))_{t \in \mathbb{R}_+}$  the solution of this DAE system.

Since we are concerned with standardization, we carefully distinguish between  $H$ , a standard vector function, and  ${}^*H = [H, H, \dots]$ , its nonstandard lifting.

**Theorem 61.** *The dAE system  $B_\bullet^t({}^*H, N)=0$ , with consistency conditions  $B_\bullet^h({}^*H, N)=0$ , possesses a unique solution. Its piecewise constant interpolation, denoted by  ${}^*X(t)$ ,  $t \in \mathbb{R}_+$ , satisfies  $X = \text{st}({}^*X)$ . Note that this standardization is independent from the particular infinitesimal time step  $\partial$ .*

Theorem 61 formalizes the practical informal rule that

$$\text{for long modes, the standardization of } B_\bullet^t(H, N)=0 \text{ with consistency conditions } B_\bullet^h(H, N)=0 \text{ is } B_i^t(H, N)=0 \text{ with consistency conditions } B_i^h(H, N)=0, \quad (132)$$

where

$$\begin{aligned} B_i^h(f, m) &=_{\text{def}} \begin{bmatrix} f \\ f' \\ \vdots \\ f^{(n-1)} \end{bmatrix} \text{ and} \\ B_i^t(f, m) &=_{\text{def}} f^{(m)}, \end{aligned} \quad (133)$$

whereas  $B_i^t(H, N)$  and  $B_i^h(H, N)$  are constructed from (133) as in Notations 39.

*Proof.* The theorem is proved by componentwise induction on  $N$ , with the help of the Lemma 62 to follow.  $\square$

**Lemma 62.** *Consider a structurally nonsingular nonstandard dAE system of the form*

$$\begin{aligned} {}^*H(X^\bullet, Y^\bullet, X, Y) &=_{\text{def}} \begin{bmatrix} {}^*G(X^\bullet, Y^\bullet, X, Y) \\ {}^*g(X) \\ {}^*g(X^\bullet) \end{bmatrix} \\ &= 0 \end{aligned} \quad (134)$$

with leading variables  $X^\bullet, Y^\bullet$ , and consider

$$\begin{aligned} \widehat{{}^*H}(X^\bullet, Y^\bullet, X, Y) &=_{\text{def}} \begin{bmatrix} {}^*G(X^\bullet, Y^\bullet, X, Y) \\ {}^*g(X) \\ {}^*(\mathbf{J}g)(X) \cdot X' \end{bmatrix} \\ &= 0 \end{aligned} \quad (135)$$

where  $\mathbf{J}g(X)$  is the Jacobian of the function  $g$  at  $X$ . Let  $({}^*X(t), {}^*Y(t))$ ,  $t \in \mathbb{R}_+$ , and  $(\widehat{{}^*X}(t), \widehat{{}^*Y}(t))$ ,  $t \in \mathbb{R}_+$ , be the piecewise constant interpolations of the solutions of nonstandard dAE systems (134) and (135), respectively. Assume that we know a standard DAE block

$$\widehat{G}(X', Y', X, Y) = 0$$

such that the solution  $(X(t), Y(t))_{t \in \mathbb{R}}$  of

$$\begin{aligned} \widehat{H}(X', Y', X, Y) &=_{\text{def}} \begin{bmatrix} \widehat{G}(X', Y', X, Y) \\ g(X) \\ \mathbf{J}g(X).X' \end{bmatrix} \\ &= 0 \end{aligned} \quad (136)$$

exists, is unique, and satisfies  $(X, Y) = \text{st}(\widehat{X}, \widehat{Y})$ . Then, we also have  $(X, Y) = \text{st}(*X, *Y)$ .

*Proof.* By statement 3 of Theorem 53, we have

$$*g(X^\bullet) = *g(X) + \partial.*(\mathbf{J}g)(X).(X' + \varepsilon)$$

where  $\varepsilon$  is infinitesimal. Hence  $*H$  rewrites

$$*H(X^\bullet, Y^\bullet, X, Y) = \begin{bmatrix} *G(X^\bullet, Y^\bullet, X, Y) \\ *g(X) \\ *g(X) + \partial.*(\mathbf{J}g)(X).(X' + \varepsilon) \end{bmatrix}$$

Since  $*H(X^\bullet, Y^\bullet, X, Y) = 0$  implies in particular  $*g(X) = 0$ , (134) rewrites

$$\begin{bmatrix} *G(X^\bullet, Y^\bullet, X, Y) \\ *g(X) \\ *(\mathbf{J}g)(X).(X' + \varepsilon) \end{bmatrix} = 0$$

and removing the infinitesimal  $\varepsilon$  does not change the standardization of this dAE system. We conclude by invoking the assumption of the lemma.  $\square$

### 8.5. Mode changes of mDAE

Assumption 6 (stated on page 72) is in force throughout this section. We study the standardization of mode change events. Having successfully performed the nonstandard structural analysis, the system at mode changes can be put in the generic form (104), and we use the decomposition (113) of the algebraic variables into impulsive and non-impulsive ones. For convenience, we repeat these formulas here:

$$(104) : X' = \frac{X^\bullet - X}{\partial} \text{ in } 0 = \mathbf{H}(X', X^\bullet, V, X) \quad (113) : V = W^i \cup W$$

Note that, in (104), the infinitesimal parameter  $\partial$  is involved in both time (by being the time step of  $\mathbb{T}$ ) and space (since derivatives  $x'$  in the original standard model are expanded as  $\frac{x^\bullet - x}{\partial}$ ). Concrete instances of this are found in the clutch example, e.g., in System (15).

Compared with the previous section that dealt with long modes, the novelty, here, is that we take a different target domain for the standardization, namely the domain of *discrete-time* transition systems—standard transition systems indexed by  $\mathbb{N}$ . Clearly, we cannot expect mapping to this domain the entire trajectory of System (104). We only aim at standardizing the finite prefix of it needed to specify the restart conditions for the new mode. The “ $\bullet$ ” shift operator defines the discrete time of this dynamics, hence we

do not need to eliminate  $\partial$  from the time domain. In contrast, we will have to eliminate the occurrences of the infinitesimal parameter  $\partial$  acting in space.

Consider again the clutch example and in particular System (24), used to determine the restart values  $\omega_1^\bullet, \omega_2^\bullet$  and the current torques  $\tau_1, \tau_2$  from the current values for the states  $\omega_1, \omega_2$ . The occurrence of  $\partial$  in the right-hand sides of equations  $(e_1^\partial), (e_2^\partial)$  is the difficulty, since  $\partial$  is nonstandard—see Comment 4 regarding this.

The standardization of (104) relies on Theorem 58. Unfortunately, it is not true in general that substituting 0 for  $\partial$  in (104) yields a structurally nonsingular system with respect to the dependent variables  $Z =_{\text{def}} (X^\bullet, V)$ . We thus need further work to be able to invoke Theorem 58. This will consist in eliminating, from System (104), the impulsive variables. For the example of the clutch, the corresponding analysis was developed in Section 2.5.2. In the following, we propose a systematic approach.

To this end, we also use the rewriting of (104), i.e., we use both:

$$(104) \quad : \quad 0 = \mathbf{H}(X', X^\bullet, V, X)$$

$$(105) \quad : \quad 0 = H(Z, X, \partial)$$

where we recall that the form (105) is obtained by expanding  $X'$  as  $\frac{X^\bullet - X}{\partial}$  in (104)—the same notational convention is used in the assumption below. Using decomposition (113) of algebraic variables into impulsive and non-impulsive ones:  $V = W^i \cup W$ , and writing

$$Y =_{\text{def}} (X^\bullet, W),$$

we assume the following about (104):

**Assumption 8.** *There exists a reduced system  $\mathbf{K}(X', X^\bullet, W, X) = 0$  with dependent variables  $Y$ , satisfying the following two conditions:*

1. *System  $K(Y, X, 0) = 0$  is structurally nonsingular, where  $K(Y, X, \partial)$  is derived from  $\mathbf{K}(X', X^\bullet, W, X)$  as explained above;*
2. *The considered model is a reduction of the original one, in that, equivalently:*

$$\begin{aligned} \mathbf{K}(X', X^\bullet, W, X) = 0 &\iff \exists W^i. \mathbf{H}(X', X^\bullet, V, X) = 0 \\ K(Y, X, \partial) = 0 &\iff \exists W^i. H(Z, X, \partial) = 0 \end{aligned}$$

In words, we can eliminate, from the original model, the impulsive variables (belonging to  $W^i$ ), and the result is structurally nonsingular for  $\partial = 0$ . We constructed such a reduced system in Section 2.5.2 by eliminating the torques for the released-to-engaged mode change in the clutch example.

**Theorem 63.** *Assumptions 6 (see page 72) and 8 (see page 84) are in force. Then, standardizing System (104) is performed by solving, for  $Y$ , the structurally nonsingular system  $K(Y, X, 0) = 0$ . Note that this standardization is independent from the particular infinitesimal time step  $\partial$ .*

Note that the latter system is standard.

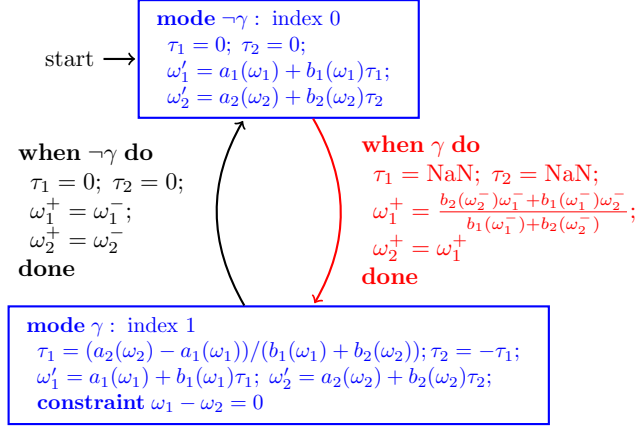


Figure 14: Standardization of ExecRun for the ideal clutch (Figure 11). Blocks have been standardized and then symbolically pivoted.  $x^-$  is the previous value of state variable  $x$ , i.e., the left limit of  $x$  when exiting a mode;  $x^+$  is the value for restart. Long modes are colored blue; non-impulsive (resp. impulsive) state-jumps are colored black (resp. red). The dynamics in mode  $\neg\gamma$  is defined by an ODE system, while in mode  $\gamma$ , it is defined by an over-determined index-1 DAE system consisting of an ODE system coupled to an algebraic constraint. In the transition from mode  $\neg\gamma$  to mode  $\gamma$ , variables  $\tau_1$  and  $\tau_2$  are impulsive, and their standardization is undefined. This explains why they are set to NaN (Not a Number).

*Proof.* Theorem 63 is a consequence of Theorem 58.  $\square$

Theorem 63 allows us to compute the restart values for the state variables that are not impulsive, which is sufficient to restart the dynamics in the coming long mode (impulsive variables at mode change cannot be initial conditions for state variables in the new mode).

We have shown in Section 2.5.2 that Assumption 8 holds for the clutch example. In general, the method would consist in, first, computing the decomposition (113), and, then, eliminating impulsive variables. Figure 14 shows the effective code resulting from the standardization of the labeled automaton of Figure 11.

To get practical, Assumption 8 (see page 84) relies on computer algebraic techniques for eliminating variables from systems of equations. Algebraic elimination is easy for systems in which impulsive variables enter linearly. If they enter polynomially, Gröbner bases can be used, but at a high computational cost. Effective elimination techniques may not exist in general, or may be too costly. Hence, the alternative based on the numerical scheme developed in Section 9.3 is of great practical interest.

## 9. Main results

This section collects results that are essential in grounding our approach. First, we prove that our approach is intrinsic, in that the actual code we generate does not depend on the particular nonstandard scheme used, instead of the explicit first-order Euler scheme, to expand the derivative. Then, we compare our approach with known schemes, for subclasses for which they are known to exist. Finally, we propose a numerical scheme to compute restart values at mode changes, in full generality, for programs accepted by

our structural analysis. This numerical scheme does not require eliminating impulsive variables using computer algebra.

### 9.1. What if we change the interpretation of the derivative?

So far, an essential step in our reasoning was the interpretation (70) for the derivative  $\omega'$ . However, several nonstandard interpretations of the derivative can be equally considered. What happens if we take an interpretation different from (70)? We first discuss this on the clutch example, then provide a general result.

#### 9.1.1. The clutch example

If the system is within a long continuous mode, the following alternative expansion exists for the derivative of  $\omega$ : for  $n \in \mathbb{N}$  finite but otherwise arbitrary,  $\omega' \approx \frac{(\omega^\bullet - \omega)^{\bullet n}}{\partial}$ . More generally, for  $\mathbf{a}=(\alpha_P, \dots, \alpha_N)$  such that  $P \leq N$ ,  $N \geq 0$ ,  $\alpha_N \neq 0$ , and  $\sum_{n=P}^N \alpha_n = 1$ , one has:

$$\omega' \approx \frac{1}{\partial} \sum_{n=P}^N \alpha_n (\omega^\bullet - \omega)^{\bullet n} . \quad (137)$$

Let us investigate how the analysis of the clutch is modified when using the nonstandard interpretation

$$\omega'_{\mathbf{a}} \stackrel{\text{def}}{=} \frac{1}{\partial} \sum_{n=P}^N \alpha_n (\omega^\bullet - \omega)^{\bullet n} \quad (138)$$

for the derivative. The nonstandard semantics of the clutch model is:

$$\left\{ \begin{array}{ll} \omega'_{\mathbf{a}1} = f_1(\omega_1, \tau_1) & (e_1^\partial) \\ \omega'_{\mathbf{a}2} = f_2(\omega_2, \tau_2) & (e_2^\partial) \\ \text{if } \gamma \text{ then } \omega_1 - \omega_2 = 0 & (e_3) \\ \text{and } \tau_1 + \tau_2 = 0 & (e_4) \\ \text{if not } \gamma \text{ then } \tau_1 = 0 & (e_5) \\ \text{and } \tau_2 = 0 & (e_6) \end{array} \right. \quad (139)$$

Algorithm 13 tells us that we must shift forward  $(e_3)$  by  $N+1$  and keep other equations unchanged; other latent equations  $(e_3^{\bullet n})$  for  $n \leq N$  are erased when solving the conflict with the previous mode:

$$\left\{ \begin{array}{ll} \omega'_{\mathbf{a}1} = f_1(\omega_1, \tau_1) & (e_1^\partial) \\ \omega'_{\mathbf{a}2} = f_2(\omega_2, \tau_2) & (e_2^\partial) \\ \text{if } \gamma \text{ then } \omega_1^{\bullet(N+1)} - \omega_2^{\bullet(N+1)} = 0 & (e_3^{\bullet(N+1)}) \\ \text{and } \tau_1 + \tau_2 = 0 & (e_4) \\ \text{if not } \gamma \text{ then } \tau_1 = 0 & (e_5) \\ \text{and } \tau_2 = 0 & (e_6) \end{array} \right. \quad (140)$$

Using  $\tau_1 = -\tau_2 \stackrel{\text{def}}{=} \tau$ , the resetting equations when entering the mode  $\gamma = \text{T}$  are

$$\left\{ \begin{array}{l} \omega'_{\mathbf{a}1} = f_1(\omega_1, \tau) \\ \omega'_{\mathbf{a}2} = f_2(\omega_2, -\tau) \\ \omega_1^{\bullet(N+1)} - \omega_2^{\bullet(N+1)} = 0 \end{array} \right. \quad (141)$$

Assuming again that the  $f_i$ 's are linear in  $\tau$ , i.e., of the form  $f_i(\omega_i, \tau) = a_i(\omega_i) + b_i(\omega_i)\tau$ , (141) rewrites:

$$\begin{cases} \sum_{n=P}^N \alpha_n (\omega_1^\bullet - \omega_1)^{\bullet n} = \partial.(a_1(\omega_1) + b_1(\omega_1)\tau) \\ \sum_{n=P}^N \alpha_n (\omega_2^\bullet - \omega_2)^{\bullet n} = \partial.(a_2(\omega_2) - b_2(\omega_2)\tau) \\ \omega_1^{\bullet(N+1)} - \omega_2^{\bullet(N+1)} = 0 \end{cases} .$$

We can eliminate  $\tau$ :

$$\begin{aligned} 0 &= \sum_{n=P}^N \alpha_n (b_2(\omega_2)(\omega_1^\bullet - \omega_1)^{\bullet n} + b_1(\omega_1)(\omega_2^\bullet - \omega_2)^{\bullet n}) \\ &\quad - \partial.(a_1(\omega_1)b_2(\omega_2) + a_2(\omega_2)b_1(\omega_1)) \\ 0 &= \omega_1^{\bullet(N+1)} - \omega_2^{\bullet(N+1)} \end{aligned} \quad (142)$$

We can zero the term in  $\partial$  in the first equation: doing so keeps the system structurally nonsingular, hence it is safe. It then remains to standardize

$$\begin{aligned} 0 &= \sum_{n=P}^N \alpha_n (b_2(\omega_2)(\omega_1^\bullet - \omega_1)^{\bullet n} + b_1(\omega_1)(\omega_2^\bullet - \omega_2)^{\bullet n}) \\ 0 &= \omega_1^{\bullet(N+1)} - \omega_2^{\bullet(N+1)} \end{aligned} \quad (143)$$

and, finally, solve it for the dependent variable  $\omega^{\bullet(N+1)} =_{\text{def}} \omega_1^{\bullet(N+1)} = \omega_2^{\bullet(N+1)}$ . For  $n = P, \dots, N$ ,  $\omega_i^{\bullet n}$  standardizes as the left-limit  $\omega_i^-$  for  $i = 1, 2$ , whereas  $\omega^{\bullet(N+1)}$  standardizes as the right-limit  $\omega^+$ . Hence, (143) standardizes as

$$0 = \alpha_N (b_2(\omega_2^-)(\omega^+ - \omega_1^-) + b_1(\omega_1^-)(\omega^+ - \omega_2^-))$$

which does not depend on  $\alpha_N$  since  $\alpha_N \neq 0$ . Finally, we get

$$\omega^+ = \frac{b_2(\omega_2^-)\omega_1^- + b_1(\omega_1^-)\omega_2^-}{b_1(\omega_1^-) + b_2(\omega_2^-)},$$

which coincides with (26).

### 9.1.2. A general result

In this section, we consider again the generic nonstandard interpretation (138) for the derivative. We investigate under which assumptions the result of our whole compilation suite is independent from the particular way derivatives are expanded when mapping the considered mDAE system to its nonstandard semantics. We consider the following assumption on  $S$ , thus characterizing a restricted class of  $\partial\text{AE}$ :

**Assumption 9.** *mDAE-system  $S$  involves only long modes, no transient modes.*

Let us compare, under Assumption 9 and Comment 54, the results of Algorithm 10 on  $S$  for the following two nonstandard semantics for the derivative:

$$\begin{aligned} \partial.x' &= x^\bullet - x \\ x^\bullet &= \partial.x' + x \end{aligned} \quad (144)$$

and



$$\begin{aligned}\partial.x'_{\mathbf{a}} &= \sum_{n=P}^N \alpha_n (x^\bullet - x)^{\bullet n} \\ \alpha_N .x^{\bullet(N+1)} &= \partial.x'_{\mathbf{a}} + \sum_{n=P}^N (\alpha_n - \alpha_{n-1}) x^{\bullet n}\end{aligned}\tag{145}$$

where  $\mathbf{a} = (\alpha_P, \dots, \alpha_N)$  satisfies  $P \leq N$ ,  $N \geq 0$ ,  $\alpha_N \neq 0$ , and  $\sum_{n=P}^N \alpha_n = 1$ , whereas  $\alpha_{P-1} = 0$  by convention. Note that nothing prevents  $P$  from being negative. For the two formulas, we have given both the direct and inverse forms. Equations (144) are obtained by taking  $P=N=0$  and  $\alpha_N=1$  in (145).

### Comparing the structural analyses

Let  $S$  be an mDAE and let  $H_S$  be the system  $H$  returned by Algorithm 10 at Line 6, i.e., the outcome of index reduction, for both a long mode and a mode change event.

We will consider the following problem: how should  $H_S$  be modified to become the outcome of the index reduction of  $S$ , when derivatives are expanded using (145) instead of the first-order Euler expansion (144)? We denote by  $^*S_{\mathbf{a}}$  and  $^*S$  the corresponding two nonstandard semantics for  $S$ . Regarding this question, two lines of Algorithm 10 are important: Line 3 and Line 6. We investigate them successively.

*Line 3 of Algorithm 10.* This consists in applying the  $\Sigma$ -method to the original mDAE system  $S$ , within a given mode; let  $G=0$  be the dynamics of  $S$  in this mode. The outcome of this algorithm only depends on the weighted bipartite graph  $\mathcal{G}$  associated to  $G$ . In the original  $\Sigma$ -method, the weight of an edge  $(f, x)$  of graph  $\mathcal{G}$  is equal to the highest differentiation degree of  $x$  in  $f$ . Similarly, if  $G$  is a dAE system (in discrete time), the weight of an edge  $(f, x)$  of  $\mathcal{G}$  is equal to the highest shifting degree of  $x$  in  $f$ . Having recalled this, we can now compare the two expansions:

- *Using expansion (144):*
  1. We begin with DAE system  $G$  and its bipartite graph  $\mathcal{G}$ ;
  2. Replacing DAE system  $G$  by dAE system  $^*G$  using expansion (144) yields  $^*\mathcal{G}$ , that is identical to  $\mathcal{G}$ .
- *Using expansion (145):*
  1. We begin with DAE system  $G$  and its bipartite graph  $\mathcal{G}$ ;
  2. Replacing DAE system  $G$  by dAE system  $^*G$  using expansion (145) yields  $^*\mathcal{G}_{\mathbf{a}}$ , in which the vertices and edges are exactly those in  $\mathcal{G}$  and all weights have been multiplied by

$$M \stackrel{\text{def}}{=} N+1 .$$

The case  $N = 0$  is the base case of the forward first-order Euler scheme, so that we do not need to consider it. Therefore, in the sequel, we assume  $N > 0$ , i.e.,  $M > 1$ .

Using the notations of (61), let  $(c_f)_{f \in F}$  and  $(c_f^{\mathbf{a}})_{f \in F}$  be the offsets associated to  $^*\mathcal{G}$  and  $^*\mathcal{G}_{\mathbf{a}}$ . Then, by Lemma 23, one gets:

**Lemma 64.** *The  $\Sigma$ -method succeeds on  $^*\mathcal{G}$  if and only if it succeeds on  $^*\mathcal{G}_{\mathbf{a}}$ , and, in this case,  $c_f^{\mathbf{a}} = M \times c_f$  for every  $f \in F$  and  $d_x^{\mathbf{a}} = M \times d_x$  for every  $x \in X(F)$ .*

In particular, each guarded equation in  $\ast S_{\mathbf{a}}$  corresponds to an equation in  $\ast S$  shifted  $M$  times. We denote by

$$(G_{\Sigma}, \overline{G}_{\Sigma}) \text{ and } (G_{\Sigma}^{\mathbf{a}}, \overline{G}_{\Sigma}^{\mathbf{a}}) \quad (146)$$

the corresponding two returns of the  $\Sigma$ -method.

*Line 6 of Algorithm 10.* If the current instant sits within a long mode, then this line just erases the consistency conditions  $\overline{G}_{\Sigma}$  and  $\overline{G}_{\Sigma}^{\mathbf{a}}$ . So, we are left with  $G_{\Sigma}$  if we use expansion (144), and  $G_{\Sigma}^{\mathbf{a}}$  if we use expansion (145). Now, dividing the infinitesimal step  $\partial$  by  $M$  in  $G_{\Sigma}$  yields  $G_{\Sigma}^{\mathbf{a}}$ . Theorem 61 tells us that the standardizations of the two systems  $G_{\Sigma}$  and  $G_{\Sigma}^{\mathbf{a}}$  yield the same DAE system.

As such, the difficulty is the handling of mode changes by Line 6, when the two different expansions are used. Is the handling of conflicts preserved, when the two different expansions are used? More precisely, for  $t_o \in \mathbb{T}$  an instant of mode change:

- Let  $G=0$  and  $\bullet G=0$  be the DAE system at the current ( $t \geq t_o$ ) and previous ( $t < t_o$ ) modes, respectively; we denote by  $(c_f, d_x)$  and  $(c_g^-, d_y^-)$  the offsets returned by the  $\Sigma$ -method for the current and previous mode, respectively;
- Let  $(G_{\Sigma}, \overline{G}_{\Sigma})$  and  $(G_{\Sigma}^{\mathbf{a}}, \overline{G}_{\Sigma}^{\mathbf{a}})$  be the returns of the  $\Sigma$ -method when applied to  $G=0$  and  $G^{\mathbf{a}}=0$ ;
- Let  $(\bullet G_{\Sigma}, \bullet \overline{G}_{\Sigma})$  and  $(\bullet G_{\Sigma}^{\mathbf{a}}, \bullet \overline{G}_{\Sigma}^{\mathbf{a}})$  be the returns of the  $\Sigma$ -method when applied to  $\bullet G=0$  and  $\bullet G^{\mathbf{a}}=0$ ;
- We denote by  $\Delta$  and  $\Delta^{\mathbf{a}}$ , and by  $\Phi$  and  $\Phi^{\mathbf{a}}$ , the contexts and sets of facts constructed by ExecRun if expansions (144) and (145) are used.

Line 6 takes as inputs the systems  $(G_{\Sigma} \cup \overline{G}_{\Sigma}) \setminus \Phi$  and  $(G_{\Sigma}^{\mathbf{a}} \cup \overline{G}_{\Sigma}^{\mathbf{a}}) \setminus \Phi^{\mathbf{a}}$ . When applying SolveConflict, the dependent variables for  $(G_{\Sigma} \cup \overline{G}_{\Sigma}) \setminus \Phi$  are all the variables of this system, both leading and non-leading, whose values were not already set by executing the previous nonstandard instant. The same holds for the dependent variables of  $(G_{\Sigma}^{\mathbf{a}} \cup \overline{G}_{\Sigma}^{\mathbf{a}}) \setminus \Phi^{\mathbf{a}}$ . The following result holds:

**Theorem 65.** *Set  $F =_{\text{def}} (G_{\Sigma} \cup \overline{G}_{\Sigma}) \setminus \Phi$  with dependent variables defined by context  $\Delta$ , and  $F^{\mathbf{a}} =_{\text{def}} (G_{\Sigma}^{\mathbf{a}} \cup \overline{G}_{\Sigma}^{\mathbf{a}}) \setminus \Phi^{\mathbf{a}}$  with dependent variables defined by context  $\Delta^{\mathbf{a}}$ . Then, the mappings*

$$\begin{aligned} \psi_f : F^{\mathbf{a}} &\mapsto F & \text{defined by } \psi_f (f^{\bullet M(c_f - m)}) &= f^{\bullet(c_f - m)} \\ \psi_x : X(F^{\mathbf{a}}) &\mapsto X(F) & \text{defined by } \psi_x (x^{\bullet M(d_x - m)}) &= x^{\bullet(d_x - m)} \end{aligned}$$

*preserve the Dulmage-Mendelsohn decompositions of  $\mathcal{G}_F$  and  $\mathcal{G}_{F^{\mathbf{a}}}$ , i.e.,  $g \in \beta_{\mathcal{D}/\mathcal{E}/\mathcal{U}}^{\mathbf{a}}$  if and only if  $\psi_f(g) \in \beta_{\mathcal{D}/\mathcal{E}/\mathcal{U}}$ , and  $y \in \beta_{\mathcal{D}/\mathcal{E}/\mathcal{U}}^{\mathbf{a}}$  if and only if  $\psi_x(y) \in \beta_{\mathcal{D}/\mathcal{E}/\mathcal{U}}$ .*

The proof of this theorem decomposes into a series of lemmas. Let  $f=0$  range over the set of equations of system  $G=0$ , and let  $n = c_f$  be the equation offset computed for  $f$  when applying the  $\Sigma$ -method to  $G=0$ .

**Lemma 66.**

1. The following formulas hold:

$$\begin{aligned}
G_\Sigma \cup \overline{G}_\Sigma &= \bigcup_{f \in G} B(f, c_f) \quad \text{and} \quad G_\Sigma^{\mathbf{a}} \cup \overline{G}_\Sigma^{\mathbf{a}} = \bigcup_{f \in G} B(f, M c_f) & (a) \\
\Delta &= \bigcup_{f \in \bullet G} B(f, c_f^- - 1) \quad \text{and} \quad \Delta^{\mathbf{a}} = \bigcup_{f \in \bullet G} B(f, M c_f^- - 1) & (b) \\
\Phi &= \bigcup_{f \in \bullet G \cap G} B(f, k_f) \quad \text{and} \quad \Phi^{\mathbf{a}} = \bigcup_{f \in \bullet G \cap G} B(f, K_f) & (c) \\
\text{where } k_f &= \min(c_f^- - 1, c_f) \quad \text{and} \quad K_f = \min(M c_f^- - 1, M c_f) & (d)
\end{aligned} \tag{147}$$

2. If  $c_f^- \leq c_f$ , then  $k_f = c_f^- - 1$  and  $K_f = M c_f^- - 1$ ;  
If  $c_f^- > c_f$ , then  $k_f = c_f$  and  $K_f = M c_f$ .

*Proof.* For the first statement, (147-a) is a direct consequence of Lemma 64 and Lemma 40 from Section 6.3.5. (147-b) and (147-c) follow from Definition 38 and Lemma 64. For the second statement, we note that, since  $M > 1$ ,  $c_f^- \leq c_f$  is equivalent to  $c_f^- - \frac{1}{M} \leq c_f$ , or  $M c_f^- - 1 \leq M c_f$ . The case  $c_f^- \leq c_f$  follows. The other case is proved similarly.  $\square$

In the following lemma, we use Notations 39 and recall that  $X(G)$  denotes the set of variables involved in the set  $G$  of functions. We also introduce the following notations.

**Notations 67.** For  $G=0$  a DAE system, we extend by convention its offsets  $(c_f)_{f \in G}$  and  $(d_x)_{x \in X(G)}$  to arbitrary  $f$  and  $x$  by setting  $c_f = d_x = -\infty$  if  $f \notin G$  and  $x \notin X(G)$ . For  $-\infty \leq m$  and  $0 \leq n$ , and  $\zeta$  a variable or a function, we set

$$\mathbf{B}(\zeta, m, n) =_{\text{def}} \begin{bmatrix} \zeta^{\bullet \widehat{m}} \\ \zeta^{\bullet(\widehat{m}+1)} \\ \vdots \\ \zeta^{\bullet n} \end{bmatrix}, \quad \text{where } \widehat{m} = \max(m, 0)$$

with the convention that  $\mathbf{B}(\zeta, m, n)$  is empty if  $m > n$ .

**Lemma 68.** Set  $F =_{\text{def}} (G_\Sigma \cup \overline{G}_\Sigma) \setminus \Phi$  and  $F^{\mathbf{a}} =_{\text{def}} (G_\Sigma^{\mathbf{a}} \cup \overline{G}_\Sigma^{\mathbf{a}}) \setminus \Phi^{\mathbf{a}}$ . Then:

1.  $F$  and  $F^{\mathbf{a}}$  rewrite as follows:

$$F = \bigcup_{f \in G} \mathbf{B}(f, k_f + 1, c_f) \quad \text{and} \quad F^{\mathbf{a}} = \bigcup_{f \in G} \mathbf{B}(f, K_f + 1, M c_f) . \tag{148}$$

2. Let  $Z$  and  $Z^{\mathbf{a}}$  be the sets of dependent variables of  $F=0$  and  $F^{\mathbf{a}}=0$ . Then:

$$Z = \bigcup_{x \in X(G)} \mathbf{B}(x, d_x^-, d_x) \quad \text{and} \quad Z^{\mathbf{a}} = \bigcup_{x \in X(G)} \mathbf{B}(x, M d_x^-, M d_x) . \tag{149}$$

*Proof.* Obvious. □

If  $c_f^- \leq c_f$ , then  $k_f+1 = c_f^-$  and  $K_f+1 = Mc_f^-$ ; thus,  $\mathbf{B}(f, k_f+1, c_f)$  and  $\mathbf{B}(f, K_f+1, c_f)$  are both nonempty. If  $c_f^- > c_f$ , then  $k_f+1 = c_f+1$  and  $K_f+1 = Mc_f+1$ ; thus,  $\mathbf{B}(f, k_f+1, c_f)$  and  $\mathbf{B}(f, K_f+1, c_f)$  are both empty. Hence, (148) rewrites:

$$F = \bigcup_{f \in G: c_f^- \leq c_f} \mathbf{B}(f, c_f^-, c_f) \text{ and } F^{\mathbf{a}} = \bigcup_{f \in G: c_f^- \leq c_f} \mathbf{B}(f, Mc_f^-, Mc_f) . \quad (150)$$

Similarly, (149) rewrites:

$$Z = \bigcup_{x \in X(G): d_x^- \leq d_x} \mathbf{B}(x, d_x^-, d_x) \text{ and } Z^{\mathbf{a}} = \bigcup_{x \in X(G): d_x^- \leq d_x} \mathbf{B}(x, Md_x^-, Md_x) . \quad (151)$$

In words,  $F^{\mathbf{a}}$  and its set of dependent variables are obtained, from  $F$  and its set of dependent variables, via an  $M$ -stretching. From formulas (150) and (151), we deduce the following result:

**Lemma 69.** *Let  $f \in G$  and  $x \in X(G)$ . Then,  $(f^{\bullet k}, x^{\bullet l})$  is an edge of  $\mathcal{G}_F$  if and only if  $(f^{\bullet Mk}, x^{\bullet Ml})$  is an edge of  $\mathcal{G}_{F^{\mathbf{a}}}$ .*

The situation is illustrated in Figures 15, 16, and 17, for the clutch and two variants of it, for  $N = 1$  and  $N = 2$ .

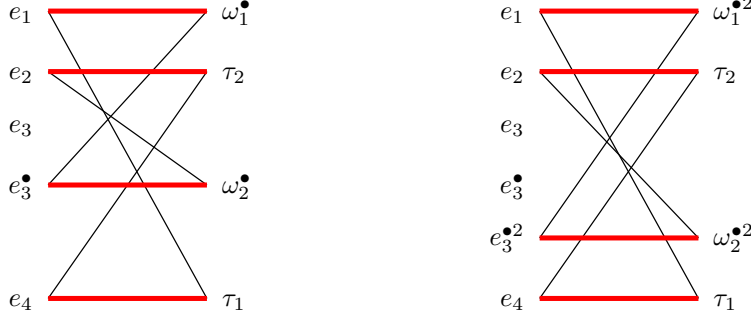


Figure 15: The clutch is an example in which  $c_f^- = c_f$  holds. Maximal cardinality matchings for  $N = 1$  (left) and  $N = 2$  (right) are given. Subsequently applying the Dulmage-Mendelsohn decomposition returns  $\beta_{\mathcal{D}} = \{e_3\}$  on the left, and  $\beta_{\mathcal{D}} = \{e_3, e_3^{\bullet}\}$  on the right.

#### *Procedure for constructing matchings of maximal cardinality for $F$ and $F^{\mathbf{a}}$*

We want to design a procedure for determining such matchings so that the two constructions parallel each other. This will allow us to compare the results.

*Step 1: starting from  $G_{\Sigma}$  and  $G_{\Sigma}^{\mathbf{a}}$ .* By construction,  $G_{\Sigma}$  and  $G_{\Sigma}^{\mathbf{a}}$  are both structurally nonsingular and determine the leading variables of the new mode, given a consistent valuation of other variables. They both admit a complete matching.

$$\text{Let } \mathcal{N}_1 \text{ be a complete matching for } G_{\Sigma}. \text{ Its edges have the form } (f^{\bullet c_f}, x^{\bullet d_x}) \text{ for appropriate pairs } (f, x) \in G \times X(G). \text{ Then, } \mathcal{N}_1^{\mathbf{a}} =_{\text{def}} \{ (f^{\bullet Mc_f}, x^{\bullet Md_x}) \mid (f^{\bullet c_f}, x^{\bullet d_x}) \in \mathcal{N}_1 \} \text{ is a complete matching for } G_{\Sigma}^{\mathbf{a}}. \quad (152)$$

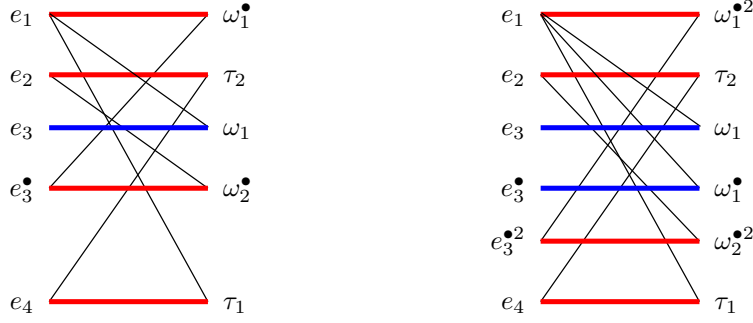


Figure 16: An example with  $c_f^- < c_f$ : a variation on the clutch. The previous mode does not involve  $\omega_1^\bullet$ , hence  $\omega_1$  remains a dependent variable when entering the new mode. Maximal cardinality matchings for  $N = 1$  (left) and  $N = 2$  (right) are shown. With reference with the example of Figure 15, the added edges correspond to consistency equations and are colored in blue. Subsequently applying the Dulmage-Mendelsohn decomposition returns  $\beta_\Sigma = \emptyset$  on both. Actually, the two matchings are complete. As a result, the mode change is fully determined, despite the new mode for the clutch having underdetermined consistency conditions (because the common rotation velocity is free).

This parallel construction is illustrated in Figure 15.

*Step 2: Removing, from the set of dependent variables of  $G_\Sigma$  and  $G_\Sigma^{\mathbf{a}}$ , the variables involved in the contexts  $\Delta$  and  $\Delta^{\mathbf{a}}$ .* Let

$$Y =_{\text{def}} \{x \in X(G) \mid x^{\bullet d_x} \in X(\Delta)\} \text{ and } Y^{\mathbf{a}} =_{\text{def}} \{x \in X(G) \mid x^{\bullet M d_x} \in X(\Delta^{\mathbf{a}})\} .$$

Then,  $Y = Y^{\mathbf{a}}$  holds. Denote by  $\widehat{G}_2$  (resp.  $\widehat{G}_2^{\mathbf{a}}$ ) the system  $G_\Sigma$  (resp.  $G_\Sigma^{\mathbf{a}}$ ) where the variables involved in the context  $\Delta$  (resp.  $\Delta^{\mathbf{a}}$ ) have been removed. The bipartite graphs of  $\widehat{G}_2$  and  $\widehat{G}_2^{\mathbf{a}}$ , denoted by  $\widehat{\mathcal{G}}_2$  and  $\widehat{\mathcal{G}}_2^{\mathbf{a}}$ , are isomorphic. Hence:

$$\text{If } \widehat{\mathcal{N}}_2 \text{ denotes a matching of maximal cardinality of } \widehat{\mathcal{G}}_2, \text{ then, } \widehat{\mathcal{N}}_2^{\mathbf{a}} =_{\text{def}} \{(f^{\bullet M c_f}, x^{\bullet M d_x}) \mid (f^{\bullet c_f}, x^{\bullet d_x}) \in \widehat{\mathcal{N}}_2\} \text{ is a matching of maximal cardinality for } \widehat{\mathcal{G}}_2^{\mathbf{a}}. \quad (153)$$

In particular, their Dulmage-Mendelsohn decompositions are isomorphic too. An example of this is illustrated in Figure 17.

We start from  $\widehat{\mathcal{N}}_2$  and  $\widehat{\mathcal{N}}_2^{\mathbf{a}}$ , perform their respective Dulmage-Mendelsohn decompositions, and remove, from both bipartite graphs, their respective overdetermined blocks  $\beta_\Sigma$  and  $\beta_\Sigma^{\mathbf{a}}$ . The resulting bipartite graphs are denoted by  $\mathcal{G}_2$  and  $\mathcal{G}_2^{\mathbf{a}}$ : they are also isomorphic. Hence:

$$\text{Let } \mathcal{N}_2 \text{ be a matching of maximal cardinality for } \mathcal{G}_2. \text{ The formula } \mathcal{N}_2^{\mathbf{a}} =_{\text{def}} \{(f^{\bullet M c_f}, x^{\bullet M d_x}) \mid (f^{\bullet c_f}, x^{\bullet d_x}) \in \mathcal{N}_2\} \text{ defines a matching of maximal cardinality for } \mathcal{G}_2^{\mathbf{a}}. \quad (154)$$

$\mathcal{N}_2$  involves all equations of  $\mathcal{G}_2$  since, by Lemma 15, both  $\mathcal{G}_2$  and  $\mathcal{G}_2^{\mathbf{a}}$  have empty overdetermined blocks in their respective Dulmage-Mendelsohn decompositions.



Figure 17: An example with  $c_f^- > c_f$ : another variation on the clutch. The previous mode determines  $\omega_2^{\bullet 2}$ , hence  $\omega_2^{\bullet}$  is no longer a dependent variable when entering the new mode. Maximal cardinality matchings for  $N = 1$  (left) and  $N = 2$  (right) are shown. Subsequently applying the Dulmage-Mendelsohn decomposition on the left-hand graph returns  $\beta_\Sigma = \mathcal{G}_F$ , hence Line 6 of Algorithm 10 returns  $H = \emptyset$ . No equation of the new mode is satisfied at the first instant and we move to the next instant. The same holds for the right-hand graph.

*Step 3: Adding to  $G_\Sigma$  the consistency subsystem  $\overline{G}_\Sigma$  and subtracting  $\Phi$ ; performing the corresponding changes to  $G_\Sigma^a$ .* We start from  $\mathcal{G}_2, \mathcal{N}_2$  and  $\mathcal{G}_2^a, \mathcal{N}_2^a$ . We do not need to remove equations belonging to  $\Phi$  and  $\Phi^a$ , since such equations, if any, have already been removed at Step 2, when moving from  $\widehat{\mathcal{N}}_2$  to  $\mathcal{N}_2$ . Therefore, the only final modification we need to perform is adding the consistency equations. The following result holds:

- Lemma 70.**
1. Let  $f \in G$  such that  $c_f^- \leq c_f$  and let  $(f^{\bullet c_f}, x^{\bullet d_x}) \in \mathcal{N}_2$ . For every  $m$  such that  $0 < m \leq c_f - c_f^-$ , the edge  $(f^{\bullet(c_f-m)}, x^{\bullet(d_x-m)})$  is an augmenting path of  $\mathcal{N}_2$  with respect to system  $F$ .
  2. Let  $f \in G$  such that  $c_f^- \leq c_f$  and let  $(f^{\bullet M c_f}, x^{\bullet M d_x}) \in \mathcal{N}_2^a$ . For every  $m$  such that  $0 < m \leq M(c_f - c_f^-)$ , the edge  $(f^{\bullet(M c_f - m)}, x^{\bullet(M d_x - m)})$  is an augmenting path of  $\mathcal{N}_2^a$  with respect to system  $F^a$ .

*Proof.* Consider Statement 1. If  $c_f \leq c_f^-$  for every  $f \in G$ , then matching  $\mathcal{N}_2$  is of maximal cardinality for  $F$ , so that no augmenting path exists. From now on, we assume  $c_f > c_f^-$  for some  $f \in G$ . Select any such  $f$  and let  $(f^{\bullet c_f}, x^{\bullet d_x})$  be the edge of  $\mathcal{N}_2$  involving  $f$ . Then,  $f^{\bullet(c_f-1)}$  and  $x^{\bullet(d_x-1)}$  are involved in  $F$ , but  $x^{\bullet(d_x-1)}$  is not a vertex of  $\mathcal{N}_2$ . Thus, for every  $f \in G$  such that  $c_f > c_f^-$ , then  $(f^{\bullet(c_f-1)}, x^{\bullet(d_x-1)})$  is an augmenting path of  $\mathcal{N}_2$  in  $F$ . The same reasoning applies for  $(f^{\bullet(c_f-m)}, x^{\bullet(d_x-m)})$ , provided that  $0 < m \leq c_f - c_f^-$ . This proves Statement 1. The same reasoning applies for proving Statement 2.  $\square$

- Lemma 71.**
1. The following formulas define matchings of maximal cardinality for  $F$  and  $F^a$ , respectively:

$$\begin{aligned} \mathcal{M} &= \left\{ (f^{\bullet(c_f-m)}, x^{\bullet(d_x-m)}) \mid (f^{\bullet c_f}, x^{\bullet d_x}) \in \mathcal{N}_2, 0 < m \leq c_f - c_f^- \right\} \\ \mathcal{M}^a &= \left\{ (f^{\bullet M(c_f-m)}, x^{\bullet M(d_x-m)}) \mid (f^{\bullet c_f}, x^{\bullet d_x}) \in \mathcal{N}_2, 0 < m \leq c_f - c_f^- \right\} \end{aligned}$$

2. Consider the surjective mapping  $\psi : \mathcal{M}^{\mathbf{a}} \mapsto \mathcal{M}$ , defined by

$$\psi(f^{\bullet M(c_f-m)}, x^{\bullet M(d_x-m)}) = (f^{\bullet(c_f-m)}, x^{\bullet(d_x-m)})$$

and define

$$\begin{aligned} \psi_f : F^{\mathbf{a}} &\mapsto F & \text{by } \psi_f(f^{\bullet M(c_f-m)}) &= f^{\bullet(c_f-m)} , \\ \psi_x : X(F^{\mathbf{a}}) &\mapsto X(F) & \text{by } \psi_x(x^{\bullet M(d_x-m)}) &= x^{\bullet(d_x-m)} . \end{aligned}$$

Then,  $\psi$  preserves the Dulmage-Mendelsohn decompositions of  $\mathcal{G}_F$  and  $\mathcal{G}_{F^{\mathbf{a}}}$ , meaning that  $g \in \beta_{\mathfrak{D}/\mathfrak{E}/\mathfrak{U}}^{\mathbf{a}}$  if and only if  $\psi_f(g) \in \beta_{\mathfrak{D}/\mathfrak{E}/\mathfrak{U}}$ , and  $y \in \beta_{\mathfrak{D}/\mathfrak{E}/\mathfrak{U}}^{\mathbf{a}}$  if and only if  $\psi_x(y) \in \beta_{\mathfrak{D}/\mathfrak{E}/\mathfrak{U}}$ .

*Proof.* Statement 1 is a direct consequence of Lemma 70 and (154). We thus focus on Statement 2. By Statement 1, we can see that  $(f^{\bullet(c_f-m)}, x^{\bullet(d_x-m)}) \in \mathcal{M}$  if and only if  $(f^{\bullet M(c_f-m)}, x^{\bullet M(d_x-m)}) \in \mathcal{M}^{\mathbf{a}}$ . Let  $g = f^{\bullet(c_f-m)}$  be a vertex of  $\mathcal{G}_F$  associated with a function. Then,  $g$  is unmatched in  $\mathcal{M}$  if and only if  $\psi_f(g)$  is unmatched in  $\mathcal{M}^{\mathbf{a}}$ . We conclude by invoking Lemma 69: the mappings  $\psi_x$  and  $\psi_f$  preserve the overdetermined blocks. A similar reasoning holds regarding the preservation of underdetermined blocks, by starting from unmatched vertices of variable type.  $\square$

This finishes the proof of Theorem 65.

#### Comparing the standardizations

The independence of the standardization with respect to the particular semantics of the derivative is immediate within long (continuous) modes. We thus focus on mode changes, where we must target a standard discrete-time dynamical system indexed by  $\mathbb{N}$ . Throughout this section, Assumptions 6 (see page 72), 8 (see page 84) and 9 (see page 87) are in force. We want to compare the respective standardizations of the returns of Algorithm 10 (ExecRun), when using the two expansions (144) and (145) for the derivatives.

We first adapt the generic system (104), computing the restart at a mode change:

$$\text{expand } X' \text{ as } \frac{1}{\partial} \sum_M^N \alpha_n (X^{\bullet} - X)^{\bullet n} \text{ in } 0 = \mathbf{H}(X', X^{\bullet}, V, X) , \quad (155)$$

$$\text{also written } 0 = H(Z, X, \partial) . \quad (156)$$

In (155),  $X$  collects the state variables, while  $V$  collects the algebraic dependent variables, decomposed into impulsive and non-impulsive variables:  $V = W^i \cup W$ . Since  $N \geq 0$ , the dependent variables of system (155) are  $Z =_{\text{def}} (X^{\bullet(N+1)}, V)$ , and we set  $Y =_{\text{def}} (X^{\bullet(N+1)}, W)$ ; in words,  $Y$  collects the non-impulsive dependent variables.

In (155), the functions  $F$  and  $G$  are standard; so is the forward shift  $x \mapsto x^{\bullet}$ , since it will be mapped to a corresponding forward shift in the discrete time index  $\mathbb{N}$ . Hence, the infinitesimal time step  $\partial$ , arising in the expansion of the derivative in  $F$ , is the only reason for System (155) to be nonstandard.

Assumption 8 states that impulsive variables can be eliminated from system (155), resulting in the following reduced system:

$$\mathbf{K}(X', X^{\bullet}, W, X) = 0 \iff \exists W^i. \mathbf{H}(X', X^{\bullet}, V, X) = 0 . \quad (157)$$

Recall that  $\partial$  enters System (157)-left, i.e., the left-hand part of the equivalence given by Equation (157), via the expansion used for  $X'$  in (155). As part of Assumption 8:

$$\text{System (157)-left remains structurally nonsingular when } \partial = 0. \quad (158)$$

In (157),  $(X^{\bullet(N+1)}, W)$  standardizes as the right-limit  $(X^+, W^+)$  whereas all the  $X^{\bullet n}$ ,  $n = P, \dots, 0$  standardize as the left-limit  $X^-$ . Therefore, in the expansion used for  $X'$  in (155), all terms cancel out but the last one. Hence, System (157) standardizes as:

- Case  $N > 0$ : substitute  $\partial \leftarrow 0$  in  $\mathbf{K} \left( \frac{\alpha_N}{\partial} (X^+ - X^-), \mathbf{X}^-, W, X \right) = 0$ ;
- Case  $N = 0$ : substitute  $\partial \leftarrow 0$  in  $\mathbf{K} \left( \frac{\alpha_N}{\partial} (X^+ - X^-), \mathbf{X}^+, W, X \right) = 0$ .

The difference lies in the  $X^-$  vs.  $X^+$ , highlighted in red. These substitutions are both legitimate, thanks to (158). In both cases, the result does not depend on  $\alpha_N$ , since  $\alpha_N \neq 0$ . Hence, we get the following result:

**Theorem 72.** *Under Assumptions 6, 8, and 9 (see pages 72, 84, and 87), the standardization of the restart conditions does not depend on the expansion (145) for the derivative.*

**Comment 73.** So far, Theorem 72 applies to mDAE systems possessing only long modes. We could extend this theorem in several directions. First, we could include in our modeling language left- or right-limits:  $x^-$  or  $x^+$ . Then, our reasoning would apply with no change provided that  $x^-$  and  $x^+$  are encoded as  $\bullet x$  or  $x \bullet$  if expansion (144) is used, and  $\bullet^{(M+1)}x$  or  $x \bullet^{(M+1)}$  if expansion (145) is used. Removing Assumption 9 could be handled similarly, by redefining transient modes as modes of nonstandard duration  $M\partial$  if expansion (145) is used, instead of  $\partial$  if expansion (144) is used.

## 9.2. A correctness result for the subclass of semilinear systems

Our approach, consisting of the structural analysis in the nonstandard domain (Algorithm 10 and its refined version Algorithm 13) followed by standardization, either fails, in which case it returns proper diagnosis, or succeeds and produces executable simulation code. The natural question is:

**Question 74.** Does this approach compute the solution of the given mDAE system?

Once again, answering Question 74 in its full generality is out of reach as, to the best of our knowledge, there is no general mathematical definition of the solution of an mDAE system. Pathologies may exist in the complement of the long ‘continuous’ modes: possibly unbounded cascades of events, chattering or sliding modes, Zeno behaviors for events of mode changes or transient modes, and more. Even when the set of events of mode changes is ‘gentle’ (a finite or diverging sequence of instants), the characterization of the impulses may not be available. Still, as mentioned in the Introduction, notions of solutions of mDAE systems exist for some restricted classes.

In this section, we provide an answer to Question 74 for a restricted subclass of mDAE systems, by reproducing the results of [8] and giving their detailed proofs. For the description of the class of systems considered, we use the notion of function of bounded



variation: a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  has *bounded variation* if it is the primitive of a Lebesgue integrable function [26]. As a consequence, if  $f$  has bounded variation, then

$$\lim_{h \searrow 0} \int_t^{t+h} f'(s) ds = 0 . \quad (159)$$

**Definition 75 (semilinear systems).** Call semilinear an mDAE system  $S$  such that, for each mode  $\mu$ , the active DAE system takes the restricted form

$$\begin{cases} 0 = A(X_s)X' + B_\mu(X) \\ 0 = C_\mu(X) \end{cases} \quad (160)$$

where:

1.  $X_s$  collects the smooth elements of  $X$ , i.e., those being continuous and of bounded variation around each instant of mode change. Other elements of  $X$  might be discontinuous.
2. Both matrix  $A(\cdot)$  and the mode-dependent vectors  $B_\mu(\cdot)$  and  $C_\mu(\cdot)$  are continuous functions of their arguments, and remain bounded in a neighborhood of each instant of mode change.
3. The matrix

$$\mathbf{J} = \begin{bmatrix} A(X_s) \\ \mathbf{J}_X C_\mu(X) \end{bmatrix}$$

is regular, where  $\mathbf{J}_X C_\mu(X)$  denotes the Jacobian matrix of  $C_\mu$  with respect to  $X$ , evaluated at  $X$ .

**Comment 76.** The special form (160) and Condition 3 for the Jacobian matrix express that mode dynamics have been processed using the  $\Sigma$ -method: differentiating the algebraic constraint once yields a regular linear system determining  $X'$ . Then, Conditions 1 and 2 together make the impulse analysis easier.

**Comment 77.** In [8], Section 4.3, it is shown that multi-body systems with contacts yield semilinear mDAE systems.

For semilinear mDAE systems, the solution at an instant  $t_*$  of mode change can be determined as follows. After proper reordering of the entries of  $X$ , we decompose vector  $X$  into its smooth and nonsmooth parts and decompose matrix  $A$  accordingly:

$$X = \begin{bmatrix} X_s \\ X_n \end{bmatrix} \quad \text{and} \quad A = [A_s \quad A_n] .$$

Since  $X_s$  is smooth at any instant, by Condition 1, matrix  $A(X_s(t_*))$  is known prior to computing the mode change. In contrast,  $X'_n(u)du$  is a Dirac measure at mode change.

To evaluate the restart values for all elements of  $X$ , we integrate the system dynamics over the interval  $[t_* - \varepsilon, t_* + \varepsilon]$  for  $\varepsilon > 0$  sufficiently small, taking into account the dynamics before and after the mode change. In the following calculations,  $\mu^-$  and  $\mu^+$

denote the mode prior and after the mode change, and  $B_+$  and  $C_+$  stand for  $B_{\mu^+}$  and  $C_{\mu^+}$ , respectively.

$$\begin{aligned}
& \int_{[t_*-\varepsilon, t_*]} (A(X_s(u))X'(u)+B_-(X(u)))du + \int_{[t_*, t_*+\varepsilon]} (A(X_s(u))X'(u)+B_+(X(u)))du \\
= & \int_{[t_*-\varepsilon, t_*]} (A_s(X_s(u))X'_s(u)+B_-(X(u)))du + \int_{[t_*, t_*+\varepsilon]} (A_s(X_s(u))X'_s(u)+B_+(X(u)))du \\
& + \int_{[t_*-\varepsilon, t_*+\varepsilon]} A_n(X_s(u))X'_n(u)du \\
\approx & 0 + 0 + A_n(X_s(t_*))(X_n(t_*+\varepsilon) - X_n(t_*-\varepsilon)) \\
\approx & A(X_s(t_*))(X(t_*+\varepsilon) - X(t_*-\varepsilon))
\end{aligned}$$

The evaluation to zero of the first two integrals follows from (159). For the third integral, we use the fact that  $X'_n(u)du$  is approximately a Dirac measure at  $t_*$  in the indicated interval. This reasoning leads to the following result determining the restart conditions at  $t_*$ :

**Theorem 78 ([8]).** *The restart conditions  $X^+$  are determined from the left-limits  $X^-$  by using the following system of equations:*

$$\begin{cases} 0 = A(X_s^-)(X^+ - X^-) \\ 0 = C_+(X^+) \end{cases} \quad (161)$$

which is locally regular thanks to Condition 3 of Definition 75.

Whether we can solve system (161) depends on the particular system at hand. It remains to compare scheme (161) with our approach. The nonstandard semantics of (160) is:

$$\begin{cases} 0 = A(X_s)(X^\bullet - X) + \partial \times B_\mu(X) \\ 0 = C_\mu(X) \end{cases}$$

and applying the  $\Sigma$ -method to it yields the augmented system

$$\begin{cases} 0 = A(X_s)(X^\bullet - X) + \partial \times B_\mu(X) \\ 0 = C_\mu(X) \\ 0 = C_\mu(X^\bullet) \end{cases} . \quad (162)$$

The reasoning performed at Section 6.2 applies to semilinear mDAE systems. Following (73), we unfold (162) at the two successive instants  ${}^\bullet t_*$  and  $t_*$ :

$$\begin{aligned}
{}^\bullet t_* : & \begin{cases} 0 = A({}^\bullet X_s)(X - {}^\bullet X) + \partial \times B_-({}^\bullet X) \\ 0 = C_-({}^\bullet X) \\ 0 = C_-(X) \end{cases} \\
t_* : & \begin{cases} 0 = A(X_s)(X^\bullet - X) + \partial \times B_+(X) \\ 0 = C_+(X) \\ 0 = C_+(X^\bullet) \end{cases} \end{aligned} \quad (163)$$

Eliminating  ${}^\bullet X$  in (163) yields

$$\begin{cases} 0 = C_-(X) & (e^-) \\ 0 = A(X_s)(X^\bullet - X) + \partial \times B_+(X) & (e_1^+) \\ 0 = C_+(X) & (e_2^+) \\ 0 = C_+(X^\bullet) & (e_3^+) \end{cases}$$

Equations  $(e^-)$  and  $(e_2^+)$  are in conflict: we discard the latter. This finally yields

$$\begin{cases} 0 = C_-(X) & (e^-) \\ 0 = A(X_s)(X^\bullet - X) + \partial \times B_+(X) & (e_1^+) \\ 0 = C_+(X^\bullet) & (e_3^+) \end{cases} \quad (164)$$

In System (164),  $(e^-)$  is a context, whereas  $(e_1^+)$  and  $(e_3^+)$  determine  $X^\bullet$  as a function of  $X$ . Setting  $\partial \leftarrow 0$  in  $(e_1^+, e_3^+)$  yields a structurally nonsingular system by Condition 3 of Definition 75. Hence, by Theorem 58, setting  $\partial \leftarrow 0$  in  $(e_1^+, e_3^+)$  performs the standardization: the scheme of Theorem 78 is recovered.

**Theorem 79 (correctness result).** *For semilinear mDAE systems, our approach computes the correct solution.*

**Comment 80.** Our approach computes solutions for any mDAE system whose structural analysis (Algorithm 10 or Algorithm 13) succeeds. This is way beyond the class of semilinear mDAE systems. It would be desirable to compare the so obtained scheme with schemes known for more dedicated physics.  $\square$

**Comment 81.** Hilding Elmqvist and Martin Otter have developed the [ModiaMath](#) tool for semilinear mDAE systems by effectively implementing scheme (161).  $\square$

### 9.3. Numerical scheme for restarts

In this section, we focus on the system of equations defining restarts, in its form (105), which we recall for convenience:

$$H(Z, X, \partial) = 0 \quad . \quad (165)$$

Our aim is to justify our claim that, if we remain ‘blind’ with respect to impulsive variables, the solution of this system is well approximated by the standard system

$$H(Z, X, \delta) = 0 \quad , \quad (166)$$

where the infinitesimal  $\partial$  has been replaced by a small standard parameter  $\delta > 0$ . This is essentially the proposed numerical scheme. Throughout this section, Assumptions 6 and 8 of Section 8.5 are in force.

We briefly recall the corresponding notations for convenience: in (165),  $X$  collects the state variables and  $Z =_{\text{def}} (X^\bullet, V)$  collects the dependent variables. The algebraic variables  $V$  can be split into impulsive and non-impulsive variables:  $V = W^i \cup W$ , and we set  $Y =_{\text{def}} (X^\bullet, W)$ . Assumption 8 states that impulsive variables can be eliminated from System (166), resulting in the reduced system

$$K(Y, X, \delta) = 0 \quad , \quad (167)$$

that remains structurally nonsingular when  $\delta = 0$ .

**Theorem 82.** *The following statements hold:*

1. Let  $Y_*$  be a solution of (167). There exists a selection map  $\delta \mapsto (Y(\delta), W^i(\delta))$ , from  $\mathbb{R}_{>0}$  to the set of solutions of the system (166), such that  $\lim_{\delta \searrow 0} Y(\delta) = Y_*$ .
2. Conversely, let  $(Y(\delta), W^i(\delta))$  be a solution of (166) such that  $Y_* =_{\text{def}} \lim_{\delta \searrow 0} Y(\delta)$  exists; then,  $Y_*$  is a solution of (167).

*Proof.* We successively prove the two statements.

*Statement 1.* By the Implicit Function Theorem applied to the standard function  $K$ , there exists a function  $(\hat{X}, \delta) \mapsto Y(\hat{X}, \delta)$  of class  $C^1$  such that  $Y(X, 0) = Y_*$  and  $K(Y(\hat{X}, \delta), \hat{X}, \delta) = 0$  holds for  $(\hat{X}, \delta)$  ranging over a neighborhood of  $(X, 0)$ . Then, by Assumption 8,  $\exists W^i. H(Z, X, \delta) = 0 \iff K(Y, X, \delta) = 0$  holds, which concludes the proof of this statement.

*Statement 2.* This is trivial and does not use Assumption 8. Since  $K$  is smooth, the two conditions (1)  $K(Y(\delta), X, \delta) = 0$  for every  $\delta > 0$ , and (2)  $Y(\delta)$  converges to  $Y_*$ , imply  $K(Y_*, X, 0) = 0$ .  $\square$

**Corollary 83.** *Let Assumption 8 be in force and let the systems (166) and (167) possess a unique solution for  $Z$  and  $Y$ , respectively. Then,  $H(Z(\delta), X, \delta) = 0$  implies  $K(\lim_{\delta \searrow 0} Y(\delta), X, 0) = 0$ .*

For its practical application, this scheme requires the prior knowledge of the decomposition  $V = W^i \cup W$  of algebraic variables into impulsive and non-impulsive ones: we solve System (166) for  $Z$  and simply discard the values of  $W^i$  from this solution. The drawback of this scheme is that system (166) will be likely ill-conditioned.

Thus, in practice, we would recommend a rescaling of the impulsive variables based on the impulsion order: for every impulsive variable  $w$  or order  $\mu =_{\text{def}} \llbracket w \rrbracket$ , we perform the substitution  $w \leftarrow \nu \times \delta^{-\mu}$ , where  $\nu$  is the new variable replacing  $w$ —such a rescaling works only for variables of finite impulsion order, see Section 8.1.

#### 9.4. Experimental results: mode changes of a non semilinear system

To exercise a numerical method computing state-jumps at mode changes, when some variables are impulsive, let us consider the following multimode DAE system:

$$\left\{ \begin{array}{ll} & \omega'_1 = a_1 \omega_1 + b_1 \tau_1^3 \quad (e_1) \\ & \omega'_2 = a_2 \omega_2 + b_2 \tau_2 \quad (e_2) \\ \text{if } \gamma \text{ then} & \omega_1 - \omega_2 = 0 \quad (e_3) \\ & \text{and } \tau_1 + \tau_2 = 0 \quad (e_4) \\ \text{if not } \gamma \text{ then} & \tau_1 = 0 \quad (e_5) \\ & \text{and } \tau_2 = 0 \quad (e_6) \end{array} \right. \quad (168)$$

This system is identical to the clutch model (3), with the only difference that torque  $\tau_1$  appears non-linearly in equation  $(e_1)$ . Therefore, this system is not semilinear, and the direct approach based on linear-algebraic techniques used in Section 9.2 does not apply. Nevertheless, we shall demonstrate on this example that our general approach applies. We focus on the mode change  $\gamma : F \rightarrow T$ . First, we perform impulse analysis following Section 8.1. Based on this impulse analysis, we apply a change of variable that bypasses the need for eliminating impulsive variables. We then standardize the resulting algebraic system of equations and we prove that it possesses a unique solution.

Using the structural analysis method detailed in Section 2.4, the state-jump is solution

of the following system of algebraic equations, where  $\omega_1^\bullet, \omega_2^\bullet, \tau_1$  and  $\tau_2$  are the unknowns:

$$\begin{cases} 0 = \frac{\omega_1^\bullet - \omega_1}{\partial} - a_1\omega_1 - b_1\tau_1^3 & (e_1^\partial) \\ 0 = \frac{\omega_2^\bullet - \omega_2}{\partial} - a_2\omega_2 - b_2\tau_2 & (e_2^\partial) \\ 0 = \omega_1^\bullet - \omega_2^\bullet & (e_3^\bullet) \\ 0 = \tau_1 + \tau_2 & (e_4) \end{cases} \quad (169)$$

Applying the method of Section 8.1, the impulse orders of the dependent variables  $\omega_1^\bullet - \omega_1, \omega_2^\bullet - \omega_2, \tau_1, \tau_2$  are solution of the following system of equalities and inequalities:

$$\begin{cases} 1 + \llbracket \omega_1^\bullet - \omega_1 \rrbracket \leq \max\{\llbracket a_1\omega_1 \rrbracket, \llbracket b_1 \rrbracket + 3\llbracket \tau_1 \rrbracket\} \\ \llbracket a_1\omega_1 \rrbracket \leq \max\{1 + \llbracket \omega_1^\bullet - \omega_1 \rrbracket, \llbracket b_1 \rrbracket + 3\llbracket \tau_1 \rrbracket\} \\ \llbracket b_1 \rrbracket + 3\llbracket \tau_1 \rrbracket \leq \max\{1 + \llbracket \omega_1^\bullet - \omega_1 \rrbracket, \llbracket a_1\omega_1 \rrbracket\} \\ 1 + \llbracket \omega_2^\bullet - \omega_2 \rrbracket \leq \max\{\llbracket a_2\omega_2 \rrbracket, \llbracket b_2 \rrbracket + \llbracket \tau_2 \rrbracket\} \\ \llbracket a_2\omega_2 \rrbracket \leq \max\{1 + \llbracket \omega_2^\bullet - \omega_2 \rrbracket, \llbracket b_2 \rrbracket + \llbracket \tau_2 \rrbracket\} \\ \llbracket b_2 \rrbracket + \llbracket \tau_2 \rrbracket \leq \max\{1 + \llbracket \omega_2^\bullet - \omega_2 \rrbracket, \llbracket a_2\omega_2 \rrbracket\} \\ \llbracket \omega_1^\bullet - \omega_1 \rrbracket = \llbracket \omega_2^\bullet - \omega_2 \rrbracket = 0 \\ \llbracket \tau_1 \rrbracket = \llbracket \tau_2 \rrbracket \end{cases} \quad (170)$$

System (170) admits a unique solution:  $\llbracket \omega_1^\bullet - \omega_1 \rrbracket = \llbracket \omega_2^\bullet - \omega_2 \rrbracket = 0$  and  $\llbracket \tau_1 \rrbracket = \llbracket \tau_2 \rrbracket = \frac{1}{3}$ . This impulse analysis tells us that if a solution to System (169) exists, then unknowns  $\tau_1$  and  $\tau_2$  are of order  $O(\partial^{-\frac{1}{3}})$ . The following substitution helps solving System (169):  $\tau_i = \partial^{-\frac{1}{3}} \hat{\tau}_i$ . With this substitution, System (169) becomes:

$$\begin{cases} 0 = \omega_1^\bullet - \omega_1 - (\partial a_1\omega_1 + b_1\hat{\tau}_1^3) & (e_1^\partial) \\ 0 = \omega_2^\bullet - \omega_2 - (\partial a_2\omega_2 + \partial^{\frac{2}{3}}b_2\hat{\tau}_2) & (e_2^\partial) \\ 0 = \omega_1^\bullet - \omega_2^\bullet & (e_3^\bullet) \\ 0 = \hat{\tau}_1 + \hat{\tau}_2 & (e_4) \end{cases} \quad (171)$$

Let us replace each unknown by its unique decomposition into a standard part and an infinitesimal part. Regarding velocities, we have:  $\omega_i^\bullet = \bar{\omega}_i^\bullet + \phi_i$ , where  $\phi_i \approx 0$ , for  $i = 1, 2$ . Next, since the rescaled variables  $\hat{\tau}_i$  are non-impulsive, the same substitution is legitimate by Lemma 50:  $\hat{\tau}_i = \bar{\tau}_i + \psi_i$ , with  $\psi_i \approx 0$ , for  $i = 1, 2$ . System (171) then becomes:

$$\begin{cases} \bar{\omega}_1^\bullet + \phi_1 - \omega_1 = \partial a_1\omega_1 + b_1(\bar{\tau}_1^3 + 3\bar{\tau}_1^2\psi_1 + 3\bar{\tau}_1\psi_1^2 + \psi_1^3) & (e_1^\partial) \\ \bar{\omega}_2^\bullet + \phi_2 - \omega_2 = \partial a_2\omega_2 + \partial^{\frac{2}{3}}b_2(\bar{\tau}_2 + \psi_2) & (e_2^\partial) \\ \bar{\omega}_1^\bullet + \phi_1 - \bar{\omega}_2^\bullet - \phi_2 = 0 & (e_3^\bullet) \\ \bar{\tau}_1 + \psi_1 + \bar{\tau}_2 + \psi_2 = 0 & (e_4) \end{cases} \quad (172)$$

Zeroing the infinitesimals in these equations yields the following nonsingular system of equations, which, by Theorem 58, is then the standardization of System (172):

$$\begin{cases} \bar{\omega}_1^\bullet - \omega_1 = b_1\bar{\tau}_1^3 & (\bar{e}_1^\partial) \\ \bar{\omega}_2^\bullet - \omega_2 = 0 & (\bar{e}_2^\partial) \\ \bar{\omega}_1^\bullet - \bar{\omega}_2^\bullet = 0 & (\bar{e}_3^\bullet) \\ \bar{\tau}_1 + \bar{\tau}_2 = 0 & (\bar{e}_4) \end{cases} \quad (173)$$

This standard system of equations admits a unique solution:

$$\bar{\omega}_1^\bullet = \bar{\omega}_2^\bullet = \omega_2 \quad \text{and} \quad \bar{\tau}_1 = -\bar{\tau}_2 = \sqrt[3]{\frac{\omega_2 - \omega_1}{b_1}}. \quad (174)$$

We now show that solving System (169) with a small standard positive parameter instead of the nonstandard  $\partial$  yields a numerical scheme converging to the solution (174), if we disregard impulsive variables. Algorithm 14 describes our numerical method, whereas Figure 18 displays the results.

---

**Algorithm 14** Numerical method for approximating the state-jump of non-semilinear system (168) for mode change  $\gamma : F \rightarrow T$ ;  $F_h(u)$  denotes the value of the right-hand side of System (169) for a 4-tuple  $u$  of values for the unknowns and  $\partial$  replaced by  $h$ .

---

**Require:** Difference equations  $F_h(u) : \mathbb{R} \times \mathbb{R}^N \rightarrow \mathbb{R}^N$ ,  $u_0 \in \mathbb{R}^N$  initial guess,  $\pi : \mathbb{R}^N \rightarrow \mathbb{R}^M$  projection on the set of non-impulsive variables,  $h_0 > 0$  initial time step,  $\epsilon > 0$  tolerance bound,  $0 < \theta < 1$  scaling factor ; **return**  $u$  estimated state-jump

- 1:  $u \leftarrow u_0$ ;  $x \leftarrow \pi(u_0)$ ;  $h \leftarrow h_0$
  - 2: **repeat**
  - 3:      $y \leftarrow x$
  - 4:      $u \leftarrow \text{NSolve}(\text{system} = F_h, \text{guess} = u, \text{tolerance} = h)$
  - 5:      $x \leftarrow \pi(u)$ ;  $h \leftarrow \theta h$
  - 6: **until**  $|x - y|_\infty \leq \epsilon$
- 

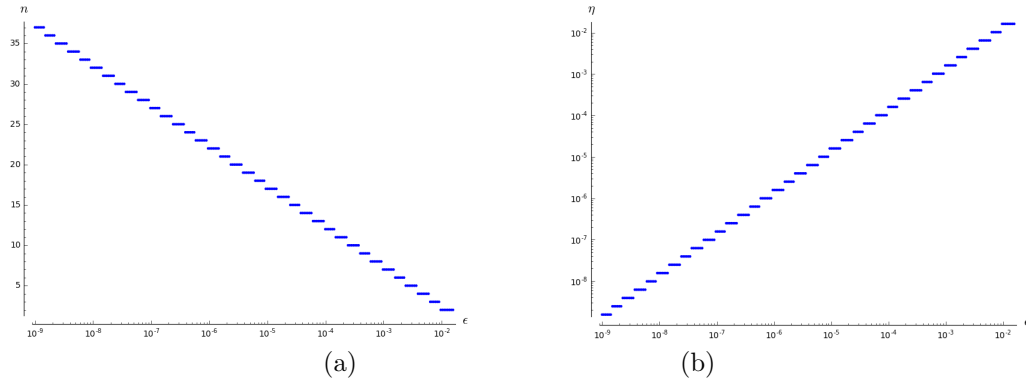


Figure 18: Numerical method for computing the state-jump of non-semilinear system (168): (a) number of iterations as a function of  $\epsilon$ ; (b) error  $\eta$  as a function of  $\epsilon$ .

Experimental results on this method have been obtained using Sage. The nonlinear solver used is the least squares solver of the Scipy's `optimize` library. Figure 18 gives, on the left, the number of iterations required, and on the right, the absolute error  $\eta = |x - \bar{x}|_\infty$  between the computed solution and the exact solution  $\bar{x} = (\omega_2, \omega_2)$ . The scaling factor is  $\theta = 0.5$ , the initial time step  $h_0 = 10^{-2}$  and  $u_0 = (\omega_1, \omega_2, 0, 0)$ . It can be observed that the method converges in not more than 37 iterations, down to  $\epsilon = 10^{-9}$ . It should also be remarked that, experimentally, the absolute error is bounded :  $\eta < 2\epsilon$ .

## 10. Toward a tool supporting our approach

The general approach we presented sets formidable challenges for a tool development. We identify three important requirements for such a tool.

- First, it should implement the theory: structural analysis must be adjusted to every different mode and mode change. Both aspects are included in what we call a ‘multimode structural analysis’.
- Second, the described approach aims at the compilation of mDAE systems. As such, ‘on-the-fly’ techniques for the structural analysis of such systems, as advocated in [17] and [27], cannot be taken into consideration. As soon as a model is structurally unsound in any given mode or at any given mode change, it should be rejected at compile time. For accepted models, the output of the tool should enable the generation of efficient simulation code.
- Third, enumerating all modes and mode changes to perform multimode structural analysis is bound to collapse. An mDAE compiler should be able to handle models in a clever way to avoid such enumeration.

We believe that we achieved an important first step towards alleviating these challenges with our tool IsamDAE under development. IsamDAE implements the structural analysis of all modes of an mDAE without any enumeration of the modes, instead relying on Binary Decision Diagrams (BDD, see for instance [18]) for implicit representations of the mode-dependent structure of an mDAE. This is described more thoroughly in [19, 20].

In order to illustrate how a multimode model is currently handled by IsamDAE, this Section focuses on the so-called RLDC2 model, a simple electronic circuit shown in Figure 19. This example, provided to us by Sven-Erik Mattsson, is actually a difficult one to handle for the existing Modelica tools because of its mode-dependent index and structure, even though it exhibits no impulsive behavior. Indeed, we show once more how this model is mishandled by two leading Modelica tools, namely, Dymola and OpenModelica, and hint at reasons for this fact, thereby justifying our work once again. We then detail its handling by the IsamDAE tool, in order to provide an insight about why this tool could be expanded into a complete multimode structural analysis suite.

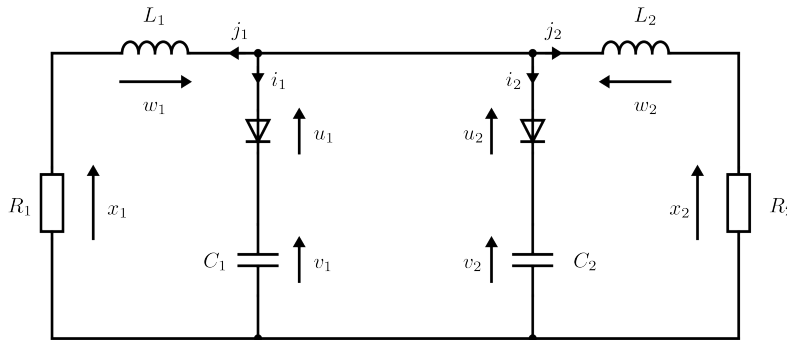


Figure 19: Schematics of the RLDC2 circuit.

### 10.1. The RLDC2 model

The circuit consists in two RLC circuits interconnected in parallel and in which two diodes have been inserted. The diodes are considered to be ideal, meaning that

they are not ruled by the customary Shockley law  $i = I(e^{u/U} - 1)$ , but rather by the complementarity condition  $0 \leq i \perp -u \geq 0$ , meaning that, at any moment, both  $i$  and  $-u$  are nonnegative and  $iu = 0$  (Section 1.2.2 already mentions complementarity conditions). Figure 20 illustrates these two diode laws.

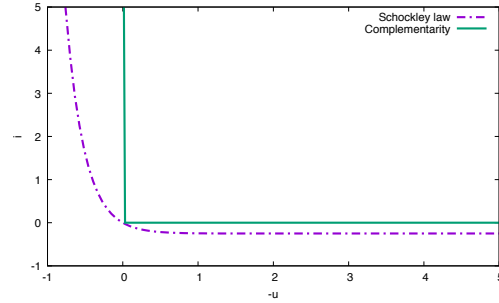


Figure 20: Shockley law (dashed, magenta) vs. ideal complementarity condition (solid, green).

As in Section 1.2.2, we redefine the graph of this complementarity condition as a parametric curve, represented by the following three equations:

$$\begin{aligned} s &= \text{if } \gamma \text{ then } i \text{ else } -u & (S) \\ 0 &= \text{if } \gamma \text{ then } u \text{ else } i & (Z) \\ \gamma &= (s \geq 0) & (G) \end{aligned}$$

To avoid a logico-numerical fixpoint, we consider instead the system obtained by replacing  $s$  with its left-limit  $s^-$  in equation (G):

$$\begin{aligned} s &= \text{if } \gamma \text{ then } i \text{ else } -u & (S) \\ 0 &= \text{if } \gamma \text{ then } u \text{ else } i & (Z) \\ \gamma &= (s^- \geq 0) & (G^-) \end{aligned}$$

Under the assumption that  $u$  and  $i$  are continuous functions of time (which turns out to be a valid assumption for the RLDC2 circuit),  $s(t) = s^-(t)$  holds at every instant  $t$ . This implies that systems (S), (Z), (G) and (S), (Z), (G<sup>-</sup>) are equivalent. Using this encoding of the complementarity condition, the RLDC2 circuit is modeled as the mDAE shown in Figure 21, with two Boolean variables  $\gamma_1$  and  $\gamma_2$  and four multimode equations (S<sub>1</sub>), (Z<sub>1</sub>), (S<sub>2</sub>) and (Z<sub>2</sub>), incident to a set of variables that depends on the mode of the system (i.e., the values of Boolean variables  $\gamma_1$  and  $\gamma_2$ ). Note that this model belongs to the class introduced in Definition 31.

### 10.2. Handling the RLDC2 model with Modelica tools

The Modelica model of the RLDC2 circuit is given in Figure 22; it is a direct translation of the model above. The authors had the opportunity to test it on two implementations of the Modelica language: OpenModelica v1.12<sup>25</sup> and Dymola 2019.<sup>26</sup>

<sup>25</sup><https://openmodelica.org>

<sup>26</sup><https://www.3ds.com/products-services/catia/products/dymola/>



$$\begin{aligned}
0 &= i_1 + i_2 + j_1 + j_2 & (K_1) \\
x_1 + w_1 &= u_1 + v_1 & (K_2) \\
u_1 + v_1 &= u_2 + v_2 & (K_3) \\
u_2 + v_2 &= x_2 + w_2 & (K_4) \\
w_1 &= L_1 \cdot j_1' & (L_1) \\
w_2 &= L_2 \cdot j_2' & (L_2) \\
i_1 &= C_1 \cdot v_1' & (C_1) \\
i_2 &= C_2 \cdot v_2' & (C_2) \\
x_1 &= R_1 \cdot j_1 & (R_1) \\
x_2 &= R_2 \cdot j_2 & (R_2) \\
s_1 &= \text{if } \gamma_1 \text{ then } i_1 \text{ else } -u_1 & (S_1) \\
s_2 &= \text{if } \gamma_2 \text{ then } i_2 \text{ else } -u_2 & (S_2) \\
0 &= \text{if } \gamma_1 \text{ then } u_1 \text{ else } i_1 & (Z_1) \\
0 &= \text{if } \gamma_2 \text{ then } u_2 \text{ else } i_2 & (Z_2) \\
\gamma_1 &= (s_1^- \geq 0) & (\gamma_1^-) \\
\gamma_2 &= (s_2^- \geq 0) & (\gamma_2^-)
\end{aligned} \tag{175}$$

Figure 21: The RLDC2 model (*n.b.*: variable  $y'$  denotes the time derivative of  $y$ ).

Identical results are obtained with both tools: the model is deemed nonsingular at compile time, yet a runtime error immediately occurs (at time  $t = 0$ ) during simulation, as shown in Figure 23.

The issue is that the structural analysis implemented in these tools treats the model as a single-mode DAE, disregarding the mode-dependent variability of the incidence relations. The consequence is that, despite the model being deemed structurally nonsingular by these tools, blocks of equations that are structurally singular in some modes are produced. This immediately leads to runtime errors: both tools attempt a pivoting of the Jacobian matrix by an element that is equal to zero in some modes.

### 10.3. Structural analysis of the RLDC2 model

The mDAE system of Figure 21 has two modes of nonzero index, when diodes are either both passing or both open. Having performed index reduction on each of them, we show the corresponding block dependency graphs in Figures 24 and 25, respectively—see Lemma 12 and Figure 10. These two graphs significantly differ, highlighting the fact that a mode-dependent structural analysis is necessary.

Figure 26 shows a graph generated by the structural analysis performed by our software IsamDAE [19, 20], for the mDAE system of Figure 21. This graph does not consist of one graph per mode, but, rather, of a set of edges labeled by a propositional formula characterizing the set of modes in which the considered branch is involved. This data structure prevents the combinatorial explosion arising in mDAE systems composed of a large number of mDAE subsystems, when modes are enumerated.

The block triangular form of the system can then be found at runtime for any given mode, for instance by evaluating the block dependency graph in this mode and performing a topological sort on the resulting graph. For the RLDC2 example, when evaluated in the mode in which both diodes are passing (resp. blocking), the graph of Figure 26 yields

```

model RLDC2_CC "RLDC2 example with complementarity conditions"
  parameter Real R1 = 10.0;
  parameter Real R2 = 15.0;
  parameter Real L1 = 1.0;
  parameter Real L2 = 1.5;
  parameter Real C1 = 0.10;
  parameter Real C2 = 0.15;
  Real i1;
  Real i2;
  Real j1(start=2.0, fixed=true);
  Real j2(start=1.0, fixed=true);
  Real u1;
  Real u2;
  Real v1(start=0.5, fixed=true);
  Real v2(start=1.0, fixed=true);
  Real w1;
  Real w2;
  Real x1;
  Real x2;
  Real s1;
  Real s2;
  Boolean g1(start=false);
  Boolean g2(start=false);
equation
  0 = j1+i1+i2+j2; // (K1)
  x1+w1 = u1+v1; // (K2)
  u1+v1 = u2+v2; // (K3)
  u2+v2 = x2+w2; // (K4)
  x1 = R1*j1; // (R1)
  x2 = R2*j2; // (R2)
  w1 = L1*der(j1); // (L1)
  w2 = L2*der(j2); // (L2)
  i1 = C1*der(v1); // (C1)
  i2 = C2*der(v2); // (C2)
  s1 = if g1 then i1 else -u1; // (S1)
  s2 = if g2 then i2 else -u2; // (S2)
  0 = if g1 then u1 else i1; // (Z1)
  0 = if g2 then u2 else i2; // (Z2)
  g1 = (s1 >= 0); // (G1)
  g2 = (s2 >= 0); // (G2)
  annotation (...);
end RLDC2_CC;

```

Figure 22: Modelica model of the RLDC2 circuit.

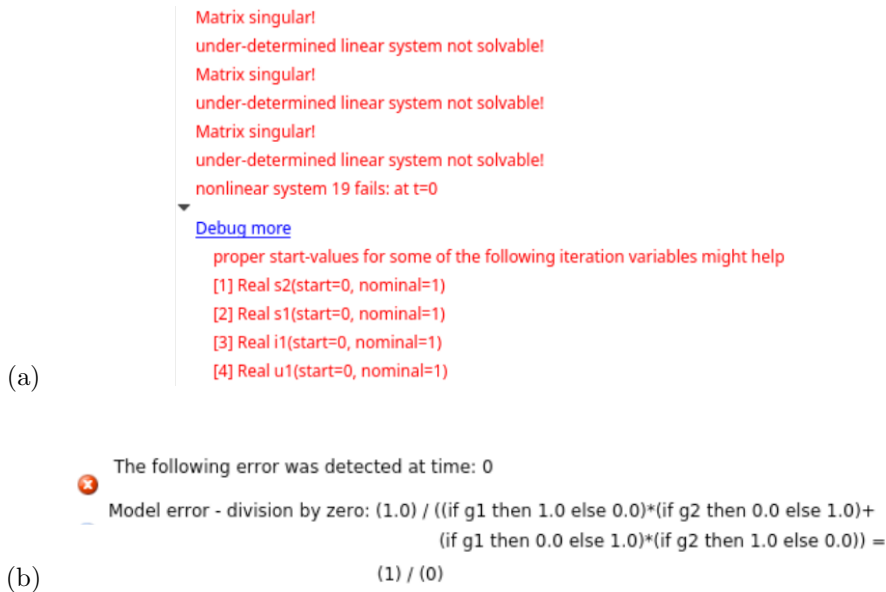


Figure 23: Error messages produced during simulation of the RLDC2 model with (a) OpenModelica and (b) Dymola, both at the initial time of the simulation.

the dependency graph of Figure 24 (resp. Figure 25); the same holds in the remaining modes. All equation blocks can be turned into efficient simulation code at compile time, so that the computational overhead due to a mode switching at runtime is minimized.

#### 10.4. Perspectives

As shown on the RLDC2 example above, our tool IsamDAE addresses the issue of performing the structural analysis of all modes of an mDAE while avoiding any explicit enumeration of these modes. Although the example detailed above is, for clarity purposes, a small one, IsamDAE is provably able to handle models with a very large number of modes, as detailed in [19, 20].

As such, the tool already provides the data structures and key algorithms required for addressing the second and third challenges invoked at the beginning of this section. This encourages us to implement, in the same tool, the structural analysis of mode changes, in the same ‘all-at-once’ fashion as for the modes themselves. This is indeed the main perspective of our work, as we believe it will be a crucial step towards the design of a mathematically sound compiler for multimode DAE systems.

## 11. Conclusion

We have proposed sound foundations for the compilation and code generation for physical systems modeling languages relying on DAE for their continuous dynamics; major instances are Modelica and Simscape. As its main contribution, our theory explicitly considers events of mode change and properly handles any multimode system having at most finite cascades of mode changes—sliding modes are not supported by our method.

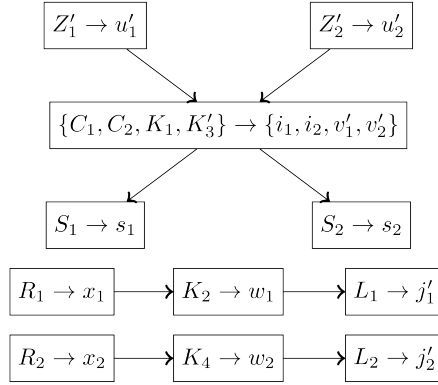


Figure 24: Block dependency graph for the RLDC2 example with both diodes passing (i.e.,  $\gamma_1 = \gamma_2 = \mathbf{T}$ ).

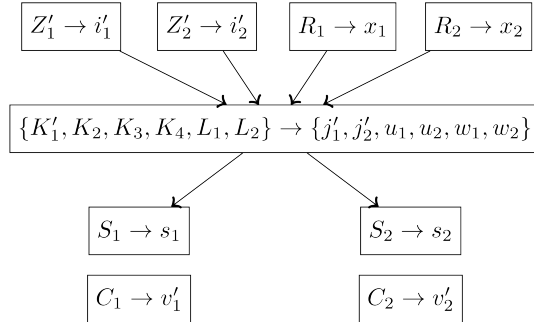


Figure 25: Block dependency graph for the RLDC2 example with both diodes blocking (i.e.,  $\gamma_1 = \gamma_2 = \mathbf{F}$ ).

A key step was the extension of structural analysis from single-mode to multimode systems. In this extension, the handling of mode changes was the main difficulty. We illustrated our approach on small examples that are not properly supported by existing tools (with the exception of [ModiaMath](#)). These examples exhibit mode-dependent index, which is not well supported by existing techniques.

For the clutch example, the most natural model consists in specifying the dynamics for the two modes “released” and “engaged”; then, the compiler automatically generates the restart conditions at mode changes. The “Westinghouse air brake” example motivated the consideration of assertions—not studied here but included in our tool under development. The Cup-and-Ball example is interesting in that, when the rope gets straight, inelastic or elastic impact can be hypothesized. This example pointed the issue of long vs. transient modes (modes having zero duration). Depending on whether the straight rope mode was assumed long (inelastic impact) or transient (elastic impact), the original model was accepted or rejected. For the latter, underspecification was detected and this information was returned to the user, who could refine his/her model by specifying an impact law. The so enriched model was accepted and code properly generated at mode change.

Generally, our method allows us to reject models on the basis of under/over-specification. Logico-numerical fixpoint equations (cases in which numerical variables

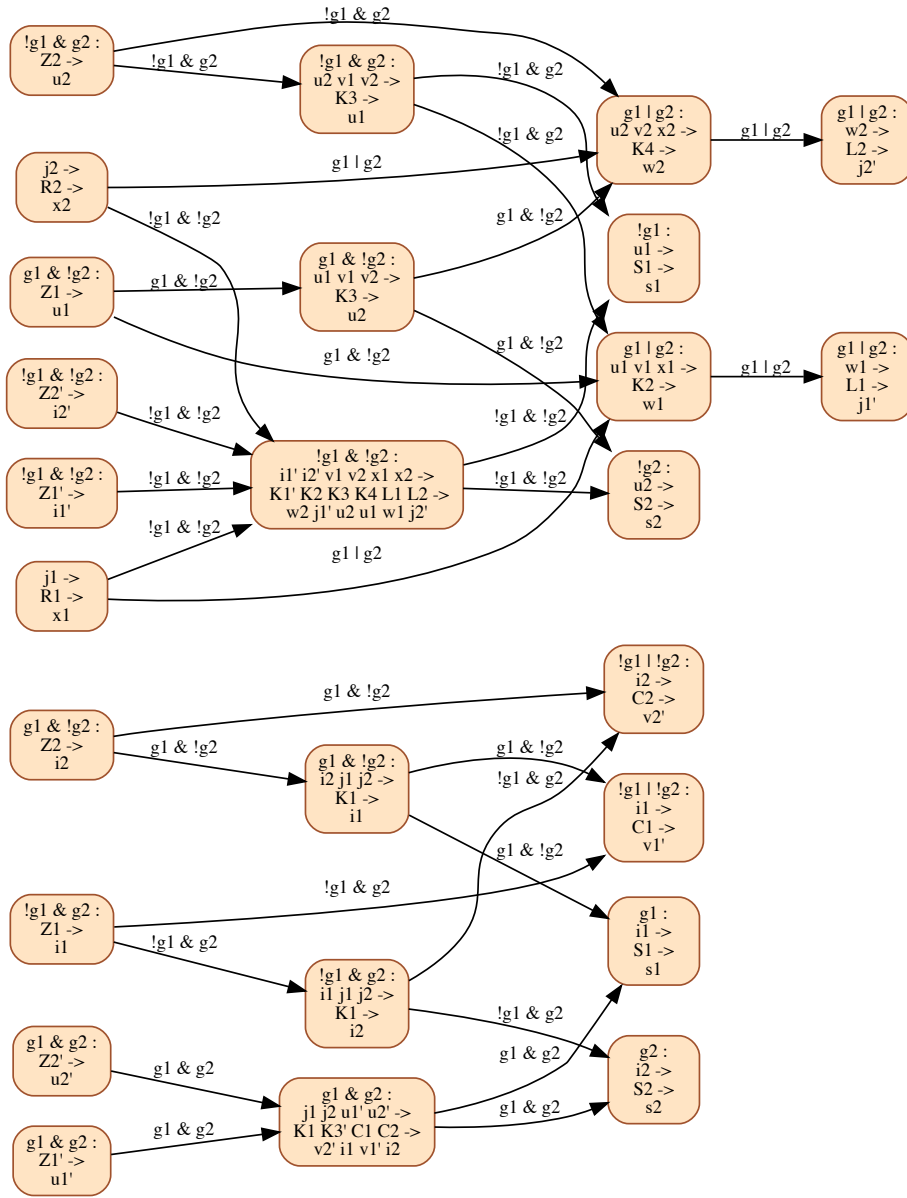


Figure 26: Block dependency graph of the RLDC2 model, generated by IsamDAE. Vertices are labeled  $p : R \rightarrow B \rightarrow W$ , where:  $p$  is a propositional formula defining in which modes the block is evaluated;  $R$  is the set of variables to read;  $B$  is the set of equations of the block;  $W$  is the set of variables to write. Edges are labeled by a propositional formula, defining in which modes the dependency applies—notation “!g” means “not g”.

determine the value of a Boolean condition that, in turn, guards the equation allowing to evaluate the same numerical variables) are prohibited by setting restrictions on the language syntax.

The structure of the code we generate is reminiscent of the so-called dummy derivatives method [37]: for both the continuous dynamics and the handling of transient modes or mode change events, we call an algebraic equation system solver for computing derivatives or next state values, as a function of the current state, while meeting algebraic constraints.

A first requirement on our theory was that it should support the analysis of any model related to any kind of physics—we aim at being physics agnostic. This support should hold regardless of the particular properties of the model, such as conditions ensuring existence and uniqueness of solutions. The motivation for this was that the user is not expected to check such conditions before submitting her/his model.

A second requirement was the ability to handle the conflicts that may result, at mode changes, between the dynamics of the previous mode and the consistency conditions implied by the new mode.

We do not see how these two requirements could be addressed without relying on nonstandard analysis. Mapping real time to a discrete time with infinitesimal step size allowed us to cast both (DAE-based) continuous dynamics, and the restart conditions at mode changes, to a uniform framework. This allowed us to properly solve the conflicts mentioned above, between the dynamics of the previous mode and the consistency conditions implied by the new mode. To make all of this doable, we developed a toolbox of new results on nonstandard analysis, mainly related to the needs of structural analysis. We proved that the result of our compilation remains independent from the particular scheme we used to map derivatives to the nonstandard domain.

From our experience in using nonstandard analysis, we can state that this is an extremely convenient tool to establish the mathematical foundations for the compilation of continuous-time systems modeling tools, for both the Simulink and Modelica classes.

A desirable objective of a work like ours is to prove that every execution scheme produced by the approach actually generates solutions of the given source mDAE system. As a reference, such soundness analyses were performed for all synchronous languages [12, 9, 10]. However, a particular difficulty holds: no mathematical definition of solutions of an mDAE system exists in full generality, that would serve as a reference. So far, we were only able to prove such a correctness result for the subclass of “semilinear multimode systems”, that may involve impulsive behaviors at mode changes.

In this work, we paid no specific attention to the efficiency of our structural analysis algorithm, and did not discuss how to generate scalable simulation code. We also only considered flat guards, whereas nested guards (naturally occurring, for instance, in the case of nested “if-then-else” constructs) are very useful in a practical formalism: flattening a model with nested guards is not desirable for getting efficient code. Work is in progress to address these issues by relying on the technique of Binary Decision Diagrams (BDD), and first results are reported in [19], where our IsamDAE tool under development is described.

### *Acknowledgements*

The authors wish to thank a number of people for this work. First, Hilding Elmqvist and Martin Otter introduced them to the subject in the early 2010’s; our collaboration then established for a couple of years, leading to both a first version of our approach and

the handling of semilinear mDAE systems; since then, Hilding and Martin have focused on their grand child of Modelica called Modia. We had extensive fruitful discussions with John Pryce on both structural analysis of DAEs, including the  $\Sigma$ -method}, and physical modeling in general; this helped us bring out the need for an explicit distinction between long and transient modes. Khalil Ghorbal participated to the first version of this approach and further interacted with the authors for this revised approach. Finally, Vincent Acary introduced us to nonsmooth systems techniques and solvers and we had frequent exchanges regarding the notions of solution for mDAE systems, revealing to us how strange these can be.

## References

- [1] Acary Vincent and Brogliato Bernard. *Numerical Methods for Nonsmooth Dynamical Systems. Applications in Mechanics and Electronics*, volume 35 of *Lecture Notes in Applied and Computational Mechanics*. Springer-Verlag, 2008.
- [2] R. Alur, C. Courcoubetis, T. Henzinger, and P. Ho. Hybrid automata: an algorithmic approach to the specification and verification of hybrid systems. In *Hybrid systems*, volume 736 of *Lecture Notes in Computer Science*, pages 209–229. Springer Verlag, 1993.
- [3] U. M. Ascher and L. R. Petzold. *Computer Methods for Ordinary Differential Equations and Differential-Algebraic Equations*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1st edition, 1998.
- [4] A. Benveniste, T. Bourke, B. Caillaud, B. Pagano, and M. Pouzet. A Type-based Analysis of Causality Loops in Hybrid Systems Modelers. *Nonlinear Analysis: Hybrid Systems*, 26:168–189, Nov. 2017.
- [5] A. Benveniste, T. Bourke, B. Caillaud, and M. Pouzet. Nonstandard semantics of hybrid systems modelers. *J. Comput. Syst. Sci.*, 78(3):877–910, 2012.
- [6] A. Benveniste, B. Caillaud, H. Elmqvist, K. Ghorbal, M. Otter, and M. Pouzet. Structural Analysis of Multi-Mode DAE Systems. Research Report RR-8933, Inria, July 2016.
- [7] A. Benveniste, B. Caillaud, H. Elmqvist, K. Ghorbal, M. Otter, and M. Pouzet. Structural analysis of multi-mode DAE systems. In *HSCC*, pages 253–263. ACM, 2017.
- [8] A. Benveniste, B. Caillaud, H. Elmqvist, K. Ghorbal, M. Otter, and M. Pouzet. Multi-mode DAE models - challenges, theory and implementation. In B. Steffen and G. J. Woeginger, editors, *Computing and Software Science - State of the Art and Perspectives*, volume 10000 of *Lecture Notes in Computer Science*, pages 283–310. Springer, 2019.
- [9] A. Benveniste, B. Caillaud, and P. L. Guernic. Compositionality in dataflow synchronous languages: Specification and distributed code generation. *Inf. Comput.*, 163(1):125–171, 2000.
- [10] A. Benveniste, P. Caspi, S. A. Edwards, N. Halbwachs, P. L. Guernic, and R. de Simone. The synchronous languages 12 years later. *Proceedings of the IEEE*, 91(1):64–83, 2003.
- [11] C. Berge. *The theory of graphs and its applications*. Wiley, 1962.
- [12] G. Berry. Constructive semantics of Esterel: From theory to practice (abstract). In *AMAST '96: Proceedings of the 5th International Conference on Algebraic Methodology and Software Technology*, page 225, London, UK, 1996. Springer-Verlag.
- [13] S. Bliudze. *Un cadre formel pour l'étude des systèmes industriels complexes: un exemple basé sur l'infrastructure de l'UMTS*. PhD thesis, Ecole Polytechnique, 2006.
- [14] S. Bliudze and D. Krob. Modelling of complex systems: Systems as dataflow machines. *Fundam. Inform.*, 91(2):251–274, 2009.
- [15] T. Bourke and M. Pouzet. Zélus: A synchronous language with ODEs. In *Hybrid Systems: Computation and Control (HSCC)*, pages 113–118, Philadelphia, USA, Apr. 2013. ACM.
- [16] K. E. Brenan, S. L. Campbell, and L. R. Petzold. *Numerical Solution of Initial Value Problems in Differential-Algebraic Equations*. SIAM, 1996.
- [17] D. Broman and J. G. Siek. Modelyze: a gradually typed host language for embedding equation-based modeling languages. Technical Report UCB/EECS-2012-173, EECS Department, University of California, Berkeley, Jun 2012.
- [18] R. E. Bryant. Graph-based algorithms for Boolean function manipulation. *IEEE Transactions on Computers*, C-35(8):677–691, aug 1986.
- [19] B. Caillaud, M. Malandain, and J. Thibault. Implicit structural analysis of multimode DAE systems. In *23rd ACM International Conference on Hybrid Systems: Computation and Control (HSCC 2020)*, Sydney, Australia, April 2020. to appear.
- [20] B. Caillaud, M. Malandain, and J. Thibault. Implicit Structural Analysis of Multimode DAE Systems. Research Report RR-9320, Inria Rennes - Bretagne Atlantique ; IRISA, Université de Rennes, Feb. 2020.
- [21] S. L. Campbell and C. W. Gear. The index of general nonlinear DAEs. *Numer. Math.*, 72:173–196, 1995.
- [22] H. Cartan. *Formes Différentielles*. Collection Méthodes. Hermann, 1967.
- [23] N. Cutland. *Nonstandard analysis and its applications*. Cambridge Univ. Press, 1988.
- [24] J. Dieudonné. *Fondements de l'analyse moderne*. Gauthier-Villars, 1965.
- [25] I. S. Duff, A. M. Erisman, and J. K. Reid. *Direct Methods for Sparse Matrices*. Numerical Mathematics and Scientific Computation. Oxford University Press, 1986.
- [26] N. Dunford and J. Schwartz. *Linear Operators, Part I, General Theory*. Wiley-Interscience, 1958.



- [27] H. Elmqvist, T. Henningsson, and M. Otter. Systems modeling and programming in a unified environment based on Julia. In T. Margaria and B. Steffen, editors, *Leveraging Applications of Formal Methods, Verification and Validation: Discussion, Dissemination, Applications*, pages 198–217, Cham, 2016. Springer.
- [28] H. Elmqvist, S.-E. Mattsson, and M. Otter. Modelica extensions for multi-mode DAE systems. In H. Tummescheit and K.-E. Arzen, editors, *Proc. of the 10th Int. Modelica Conference*, Lund, Sweden, Sept. 2014. Modelica Association.
- [29] I. Hasuo. Metamathematics for systems design - comprehensive transfer of formal methods techniques to cyber-physical systems. *New Generation Comput.*, 35(3):271–305, 2017.
- [30] W. P. M. H. Heemels, M. K. Camlibel, and J. M. Schumacher. On the dynamic analysis of piecewise-linear networks. *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, 49(3):315–327, Mar 2002.
- [31] C. Höger. Dynamic structural analysis for daes. In *Proceedings of the 2014 Summer Simulation Multiconference, SummerSim 2014, Monterey, CA, USA, July 6-10, 2014*, page 12. SCS/ ACM, 2014.
- [32] C. Höger. Elaborate control: variable-structure modeling from an operational perspective. In D. Zimmer and B. Bachmann, editors, *Proceedings of the 8th International Workshop on Equation-Based Object-Oriented Modeling Languages and Tools, EOOLT '17, Weßling, Germany, December 1, 2017*, pages 51–60. ACM, 2017.
- [33] Y. Iwasaki, A. Farquhar, V. Saraswat, D. Bobrow, and V. Gupta. Modeling time in hybrid systems: How fast is “instantaneous”? In *IJCAI*, pages 1773–1781, 1995.
- [34] D. Liberzon and S. Trenn. Switched nonlinear differential algebraic equations: Solution theory, lyapunov functions, and stability. *Automatica*, 48(5):954–963, 2012.
- [35] T. Lindström. An invitation to nonstandard analysis. In N. Cutland, editor, *Nonstandard Analysis and its Applications*, pages 1–105. Cambridge Univ. Press, 1988.
- [36] S.-E. Mattsson, M. Otter, and H. Elmqvist. Multi-Mode DAE Systems with Varying Index. In H. Elmqvist and P. Fritzson, editors, *Proc. of the 11th Int. Modelica Conference*, Versailles, France, Sept. 2015. Modelica Association.
- [37] S. E. Mattsson and G. Söderlind. Index reduction in Differential-Algebraic Equations using dummy derivatives. *Siam J. Sci. Comput.*, 14(3):677–692, 1993.
- [38] V. Mehrmann and L. Wunderlich. Hybrid systems of differential-algebraic equations – analysis and numerical solution. *Journal of Process Control*, 19(8):1218 – 1228, 2009. Special Section on Hybrid Systems: Modeling, Simulation and Optimization.
- [39] H. Nakamura, K. Kojima, K. Suenaga, and A. Igarashi. A nonstandard functional programming language. In B. E. Chang, editor, *Programming Languages and Systems - 15th Asian Symposium, APLAS 2017, Suzhou, China, November 27-29, 2017, Proceedings*, volume 10695 of *Lecture Notes in Computer Science*, pages 514–533. Springer, 2017.
- [40] C. Pantelides. The consistent initialization of differential-algebraic systems. *SIAM J. Sci. Stat. Comput.*, 9(2):213–231, 1988.
- [41] P. Pepper, A. Mehlhase, C. Höger, and L. Scholz. A compositional semantics for modelica-style variable-structure modeling. In *th International Workshop on Equation-Based Object-Oriented Modeling Languages and Tools*, 2011.
- [42] L. Petzold. Differential algebraic equations are not ODEs. *SIAM J. Sci. Stat. Comput.*, 3:367–384, 1982.
- [43] F. Pfeiffer. On non-smooth multibody dynamics. *Proceedings of the Institution of Mechanical Engineers, Part K: Journal of Multi-body Dynamics*, 226(2):147–177, 2012.
- [44] F. Pfeiffer and C. Glocker. *Multibody Dynamics with Unilateral Contacts*. Wiley, 2008.
- [45] A. Platzer. *Logical Foundations of Cyber-Physical Systems*. Springer Publishing Company, Incorporated, 1st edition, 2018.
- [46] A. Pothen and C. Fan. Computing the block triangular form of a sparse matrix. *ACM Trans. Math. Softw.*, 16(4):303–324, 1990.
- [47] D. Potop-Butucaru, S. A. Edwards, and G. Berry. *Compiling Esterel*. Springer, 2007.
- [48] J. D. Pryce. A simple structural analysis method for DAEs. *BIT*, 41(2):364–394, 2001.
- [49] A. Robinson. *Nonstandard Analysis*. Princeton Landmarks in Mathematics, 1996. ISBN 0-691-04490-2.
- [50] S. Schoeder, H. Ulbrich, and T. Schindler. Discussion of the Gear–Gupta–Leimkuhler method for impacting mechanical systems. In *Multibody System Dynamics*, volume 31, pages 477–495, 2013.
- [51] K. Suenaga, H. Sekine, and I. Hasuo. Hyperstream processing systems: nonstandard modeling of continuous-time signals. In R. Giacobazzi and R. Cousot, editors, *The 40th Annual ACM*

*SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '13, Rome, Italy - January 23 - 25, 2013*, pages 417–430. ACM, 2013.

- [52] S. Trenn. *Distributional Differential Algebraic Equations*. PhD thesis, Technischen Universität Ilmenau, 2009.
- [53] S. Trenn. Regularity of distributional differential algebraic equations. *MCSS*, 21(3):229–264, 2009.
- [54] D. Zimmer. *Equation-Based Modeling of Variable-Structure Systems*. PhD thesis, ETH Zürich, [No. 18924](#), 2010.