



# LUND UNIVERSITY

## Symbolic elimination in dynamic optimization based on block-triangular ordering

Magnusson, Fredrik; Åkesson, Johan

*Published in:*  
Optimization Methods and Software

*DOI:*  
[10.1080/10556788.2016.1270944](https://doi.org/10.1080/10556788.2016.1270944)

2018

*Document Version:*  
Peer reviewed version (aka post-print)

[Link to publication](#)

*Citation for published version (APA):*  
Magnusson, F., & Åkesson, J. (2018). Symbolic elimination in dynamic optimization based on block-triangular ordering. *Optimization Methods and Software*, 33(1), 92-119. <https://doi.org/10.1080/10556788.2016.1270944>

*Total number of authors:*  
2

*Creative Commons License:*  
Unspecified

### General rights

Unless other specific re-use rights are stated the following general rights apply:  
Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Read more about Creative commons licenses: <https://creativecommons.org/licenses/>

### Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

LUND UNIVERSITY

PO Box 117  
221 00 Lund  
+46 46-222 00 00

## Symbolic elimination in dynamic optimization based on block-triangular ordering

Fredrik Magnusson<sup>a\*</sup> and Johan Åkesson<sup>b</sup><sup>a</sup>*Department of Automatic Control, Lund University, SE-221 00 Lund, Sweden;*<sup>b</sup>*Modelon AB, Ideon Science Park, SE-223 70 Lund, Sweden**(Received 00 Month 20XX; final version received 00 Month 20XX)*

We consider dynamic optimization problems for systems described by differential-algebraic equations (DAEs). Such problems are usually solved by discretizing the full DAE. We propose techniques to symbolically eliminate many of the algebraic variables in a preprocessing step before discretization. These techniques are inspired by the causalization and tearing techniques often used when solving DAE initial value problems. Since sparsity is crucial for some dynamic optimization methods, we also propose a novel approach to preserving sparsity during this procedure.

The proposed methods have been implemented in the open-source JModelica.org platform. We evaluate the performance of the methods on a suite of optimal control problems solved using direct collocation. We consider both computational time and probability of solving the problem in a timely manner. We demonstrate that the proposed methods often are an order of magnitude faster than the standard way of discretizing the full DAE, and also significantly increase probability of successful convergence.

**Keywords:** dynamic optimization; differential-algebraic equations; block-triangular ordering; tearing; sparsity preservation; nonlinear programming; Modelica

*AMS Subject Classification:* 34A09; 49M37

### 1. Introduction

Dynamic optimization refers to optimization problems with differential equations as constraints. These problems occur in many different fields and contexts, including optimal control, parameter and state estimation, and design optimization. Examples of applications are minimization of material and energy consumption during setpoint transitions in power plants [39] and chemical processes [52], estimating occupancy and ambient air flow in buildings [61], and optimal experimental design for estimation of kinetic parameters in fed-batch processes [10].

There are many effective methods available for numerical solution of dynamic optimization problems [13, 53]. Most of these are applicable not only to dynamic systems described by explicit ordinary differential equations (ODE), but also to differential-algebraic equations (DAE). When dynamically simulating DAE systems—that is, solving initial value problems—one can either use specialized DAE solvers [17, 33], or transform the DAE to an equivalent ODE and then apply run-of-the-mill ODE solvers [20]. Key techniques in such transformations are block-triangular orderings to decompose the problem and tearing to reduce the size of the equation systems solved iteratively. While such transformation techniques are commonly applied when simulating DAEs, the conventional approach to solving DAE-constrained optimization problems is to discretize the full DAE [14, 15]. In this paper we consider the use of ODE transformation

---

\*Corresponding author. Email: fredrik.magnusson@control.lth.se

techniques on dynamic optimization problems. But rather than completing the transformation all the way to an explicit ODE, we will try to find the middle ground between the full implicit DAE and the underlying explicit ODE that maximizes computational efficiency and convergence robustness. We build upon previous work by the authors [46] and extend it with tearing, sparsity preservation and a thorough performance benchmark.

While these techniques can be applied in combination with any of the standard numerical methods for solving DAE-constrained optimization problems, the performance of some methods will be affected more than others. In this paper we focus on the use of direct collocation and how it is affected by these transformation techniques. The effect of these techniques is also dependent on the considered DAE. We will show the effect of symbolic transformation techniques on computational speed and convergence robustness for six different large-scale benchmark problems with different properties from various physical domains. The presented methods have been implemented in JModelica.org<sup>1</sup> [3], which is an open-source tool for analysis of dynamic systems described by the object-oriented, equation-based modelling language Modelica [30]. Modelica models tend to be large, sparse, and have considerably more algebraic than differential variables, making them a prime target for the considered methods. We will however show that the presented methods hold merit also for more typical textbook DAEs.

The contributions of this paper are the adaptation of causalization and tearing techniques to direct collocation for dynamic optimization and the development of a novel sparsity preservation technique. We also present a thorough benchmark of these techniques, showing how they improve both computational speed and convergence robustness.

The outline of the paper is as follows. Section 2 goes through a simple and illuminating example of how the ideas described in the paper are applied. Section 3 presents the theoretical background of dynamic optimization and the symbolic techniques we will employ, as well as JModelica.org and other related work. Section 4 presents how these techniques can be adapted to symbolically eliminate variables in dynamic optimization problems while preserving sparsity. Section 5 compares the effects of using and not using these techniques on six benchmark problems. Finally, Section 6 summarizes the paper and discusses potential future work.

Regarding notation, scalars and scalar-valued functions are denoted by regular italic letters  $x$ . Vectors and vector-valued functions are denoted by bold italic letters  $\mathbf{x}$ . The dimension of  $\mathbf{x}$  is denoted by  $n_x$  and component  $k$  of  $\mathbf{x}$  is denoted by  $x_k$ . Matrices are denoted by bold Roman letters  $\mathbf{A}$  and element  $(i, j)$  of  $\mathbf{A}$  is denoted by  $A_{i,j}$ . Row  $i$  and column  $j$  of  $\mathbf{A}$  are denoted by  $\mathbf{A}_{i,:}$  and  $\mathbf{A}_{:,j}$ , respectively. Systems of equations  $\mathbf{f} = \mathbf{0}$  that are to be considered as being parametrized by  $\mathbf{p}$  and having  $\mathbf{x}$  as unknowns are denoted by  $\mathbf{f}(\mathbf{x}; \mathbf{p}) = \mathbf{0}$ .

## 2. Illustrative example

To demonstrate the ideas of the symbolic eliminations we will use, we present a simple example. Steps that are unclear to the reader should be revisited after Section 3. Consider the DAE

$$\dot{x} + y_1 + y_2 - y_3 = 0, \quad (1a)$$

$$xy_3 + y_2 - \sqrt{x} - 2 = 0, \quad (1b)$$

$$2y_1y_2y_4 - \sqrt{x} = 0, \quad (1c)$$

$$y_1y_4 + \sqrt{y_3} - x - y_4 = 0, \quad (1d)$$

$$y_4 - \sqrt{y_5} = 0, \quad (1e)$$

$$y_5^2 - x = 0. \quad (1f)$$

---

<sup>1</sup><http://www.jmodelica.org>

A structural incidence matrix for (1) is (2).

	$\dot{x}$	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$
(1a)	1	1	1	1	0	0
(1b)	0	0	1	1	0	0
(1c)	0	*	*	0	*	0
(1d)	0	*	0	*	*	0
(1e)	0	0	0	0	1	*
(1f)	0	0	0	0	0	*

(2)

where zeroes indicate independence, ones indicate linear dependence, and asterisks indicate nonlinear dependence. Inspection of (1) reveals that we can eliminate  $y_4$  using (1e). That is, we use

$$y_4 = \sqrt{y_5} \quad (3)$$

to substitute  $y_4$  with  $\sqrt{y_5}$ , yielding

$$\dot{x} + y_1 + y_2 - y_3 = 0, \quad (4a)$$

$$xy_3 + y_2 - \sqrt{x} - 2 = 0, \quad (4b)$$

$$2y_1y_2\sqrt{y_5} - \sqrt{x} = 0, \quad (4c)$$

$$y_1\sqrt{y_5} + \sqrt{y_3} - x - \sqrt{y_5} = 0, \quad (4d)$$

$$y_5^2 - x = 0. \quad (4e)$$

Such eliminations can be identified by permuting (2) to a block-triangular matrix with blocks along the diagonal—henceforth referred to as diagonal blocks, despite them not being diagonal in general—of minimal size. Permuting (2) to such a form yields (5).

	$y_5$	$y_4$	$y_1$	$y_2$	$y_3$	$\dot{x}$
(1f)	*	0	0	0	0	0
(1e)	*	1	0	0	0	0
(1b)	0	0	0	1	1	0
(1c)	0	*	*	*	0	0
(1d)	0	*	1	0	*	0
(1a)	0	0	*	*	*	1

(5)

Note that some previously linear incidences have been reclassified as nonlinear and vice versa. We are no longer interested in the linearity of incidences outside of the diagonal blocks, and thus mark them all with asterisks. Within the diagonal blocks, we only care about linearity with respect to the variables in the block. The incidence of the equation-variable pair ((1d),  $y_1$ ) that was previously considered to be nonlinear is thus now considered to be linear. An elimination such as (3) can be found by identifying the diagonal blocks that are scalar and linear. The only algebraic variable that can be eliminated with this approach in this example is  $y_4$ .

We can however go further in our elimination procedure by eliminating  $y_1$  and  $y_2$  from (1d)

and (1b), respectively. That is, we identify

$$y_1 = \frac{x + \sqrt{y_5} - \sqrt{y_3}}{\sqrt{y_5}}, \quad (6a)$$

$$y_2 = 2 + \sqrt{x} - xy_3, \quad (6b)$$

and through substitution obtain the DAE

$$\dot{x} + \frac{x + \sqrt{y_5} - \sqrt{y_3}}{\sqrt{y_5}} + 2 + \sqrt{x} - xy_3 - y_3 = 0, \quad (7a)$$

$$2 \frac{x + \sqrt{y_5} - \sqrt{y_3}}{\sqrt{y_5}} (2 + \sqrt{x} - xy_3) \sqrt{y_5} - \sqrt{x} = 0, \quad (7b)$$

$$y_5^2 - x = 0. \quad (7c)$$

The eliminations (6) can be found by tearing the  $3 \times 3$  diagonal block. Selecting  $y_3$  as tearing variable and (1c) as tearing residual, we get the torn incidence matrix (8).

	$y_5$	$y_4$	$y_1$	$y_2$	$y_3$	$\dot{x}$
(1f)	*	0	0	0	0	0
(1e)	*	1	0	0	0	0
(1d)	0	*	1	0	*	0
(1b)	0	0	0	1	1	0
(1c)	0	*	*	*	0	0
(1a)	0	0	*	*	*	1

(8)

The eliminations (6) are then identified as the diagonal incidences in the upper left subblock of the torn block, which are feasible since this subblock is triangular (in this case it is even diagonal) and linear along the diagonal. Comparing (7) with (1), we have only 2 instead of 5 algebraic variables, but the residuals are more complicated.

### 3. Background

In this section we first present the class of DAEs and dynamic optimization problems that we consider. We then describe the open-source tool JModelica.org and in particular how it solves these problems. We review how ODE transformation techniques are employed when solving DAE simulation problems, in the context of which they are often referred to as causalization. We also describe well-established techniques for selecting pivot elements in direct, sparse linear solvers, which we later will utilize for sparsity preservation. Finally, we discuss related work.

#### 3.1 Differential-algebraic equations

We consider systems whose dynamics are described by an implicit DAE system. Specifically, we consider DAE systems of the form

$$\Phi(t, \dot{\xi}(t), \xi(t), \mathbf{v}(t), \mathbf{u}(t), \mathbf{p}) = \mathbf{0}, \quad (9)$$

where  $t$  is time—the sole independent variable— $\xi$  is the differential variable,  $\mathbf{v}$  is the algebraic variable,  $\mathbf{u}$  is the control variable, and  $\mathbf{p}$  is the parameters to be optimized. We assume that

the system is balanced; that is, that the dimension of the codomain of  $\Phi$  equals the sum of the dimensions of  $\xi$  and  $v$ .

To deal with high-index DAEs, we utilize JModelica.org's index reduction, which is based on the method of dummy derivatives [49] and reduces the index of the DAE to at most 1. The method of dummy derivatives prevents the numerical drift of straightforward reduction by differentiation and is applicable to a wide class of DAEs. The method however requires the selection of state variables, which is difficult to automate. In general it is not even possible to preserve solvability with a static choice of state variables [48], which however is out of scope for this paper.

We thus assume that there exists a static choice of state variables  $x$  out of the differential variables  $\xi$  that does not give rise to solvability issues, in particular implying that the local index of the DAE [17] is constant along all feasible trajectories. The index reduction then creates an equivalent—in the sense that there exists a bijection between the solution sets—low-index DAE

$$F(t, \dot{x}(t), x(t), y(t), u(t), p) = 0, \quad (10)$$

where  $y$  is the algebraic variables, consisting of the original algebraic variables  $v$ , those differential variables  $\xi$  not selected as states, and the newly created dummy derivatives. We also assume that (feasible) initial conditions  $x(0) = x_0$  are given and that  $F$  is twice continuously differentiable with respect to all of its arguments except the first.

### 3.2 Dynamic optimization

The problem considered throughout the paper is to

$$\text{minimize} \quad \phi(t_f, \xi(t_f), v(t_f), p) + \int_0^{t_f} \ell(t, \xi(t), v(t), u(t)) dt, \quad (11a)$$

$$\text{with respect to} \quad \xi : [0, t_f] \rightarrow \mathbb{R}^{n_\xi}, \quad v : [0, t_f] \rightarrow \mathbb{R}^{n_v}, \quad u : [0, t_f] \rightarrow \mathbb{R}^{n_u}, \\ t_f \in \mathbb{R}, \quad p \in \mathbb{R}^{n_p},$$

$$\text{subject to} \quad \Phi(t, \dot{\xi}(t), \xi(t), v(t), u(t), p) = 0, \quad x(0) = x_0, \quad (11b)$$

$$L \leq (\dot{\xi}(t), \xi(t), v(t), u(t), p) \leq U, \quad g(t, \dot{\xi}(t), \xi(t), v(t), u(t), p) \leq 0, \quad (11c)$$

$$G(\xi(t_f), v(t_f), p) = 0, \quad 0 \leq t_{f,L} \leq t_f \leq t_{f,U}, \quad (11d) \\ \forall t \in [0, t_f],$$

where the objective (11a) comprises the Mayer term  $\phi$  and Lagrange integrand  $\ell$ , (11b) is the system dynamics, (11c) is inequality constraints (where box constraints  $L$  and  $U$  are separated from the general nonlinear equalities  $g$ ), and (11d) is the terminal constraints  $G$  and nonnegative bounds on  $t_f$ . The bounds  $L$  and  $U$  can be infinite.

The problem will in general be nonconvex and we will not endeavour to find a global optimum. We will instead rely on first-order necessary optimality conditions to find a local optimum. Just like for  $F$ , we assume that  $\phi$ ,  $\ell$ ,  $g$ , and  $G$  are twice continuously differentiable in order to apply techniques based on Newton's method to find a solution to first-order optimality conditions.

The methods and software framework we present in this paper can be applied on the wider class of problems that is considered by Magnusson and Åkesson [45], which in particular includes implicit initial conditions and point constraints in the interior of the time horizon. However, (11) captures the essentials and the generalization of the presented methods to the wider class is straightforward, and so we restrict ourselves to (11) for the sake of brevity.

We mainly consider the use of direct collocation to solve (11). The main idea of direct collo-

cation is to discretize (11) by dividing the time horizon into a finite number of elements and then within each element approximate the system trajectories by polynomials defined by collocation points [14, 15, 45]. The result is a nonlinear program (NLP) whose solution approximates the solution of (11). This approach leads to the full DAE being exposed to the NLP solver. This is in contrast to sequential approaches based on shooting, in which the DAE (or at least parts of it) are hidden in embedded integrator calls. The methods discussed in this paper are thus more relevant for simultaneous approaches, such as direct collocation, where they have a major impact on the size and structure of the NLP, whereas they often only affect the NLP function evaluation times when using sequential approaches. Nevertheless, these methods certainly have the potential to be beneficial also when using sequential approaches, since the DAE simulation in shooting methods is often the computational bottleneck and a source of numerical problems.

### 3.3 *JModelica.org*

JModelica.org [3] is an open-source tool designed for analysis of models described in the acausal modelling language Modelica [30]. Modelica [30] is a standardized language for which several different tools exist. JModelica.org historically has a strong focus on dynamic optimization with its support for the language extension Optimica [2] and an efficient dynamic optimization framework [45]. Today it also has a strong support for FMI [16], allowing inter-tool exchange of models for dynamic simulation.

While acausal modelling is convenient for the modeller, the resulting mathematical model is often unwieldy in its raw form and requires symbolic processing before it can be simulated efficiently using numerical methods. In Section 3.4 we review techniques that often are used to alleviate these issues. While these techniques often are crucial for handling object-oriented, acausal models, we will in Section 5 see that they hold merit also when applied to examples of specialized, non-hierarchical models in the context of dynamic optimization.

The dynamic optimization framework in JModelica.org leverages the JModelica.org compiler to first perform structural analysis and some symbolic transformations and then transfer the resulting flattened problem to CasADi Interface [41]. The symbolic representation of CasADi Interface—which is based on CasADi [4] for efficient computation of derivatives using algorithmic differentiation—is then connected to algorithms based on direct collocation for solving dynamic optimization problems by transcribing them into NLPs. The NLPs are then finally solved by IPOPT [60] or WORHP [18].

It is further worth noting that while the overarching ideas of this paper are widely applicable, some of the detailed algorithmic choices discussed in Section 4 have been adapted to the dynamic optimization framework of JModelica.org, that is, the use of IPOPT, CasADi and the collocation algorithm in JModelica.org [45].

### 3.4 *Causalization*

The process of transforming an implicit DAE

$$\Phi(t, \dot{\xi}(t), \xi(t), v(t), u(t), p) = \mathbf{0} \quad (12)$$

to an equivalent explicit ODE

$$\dot{x}(t) = f(t, x(t), u(t), p) \quad (13)$$

is called causalization [20]. We divide this procedure into the following steps.

- (1) Alias elimination
- (2) Variability propagation
- (3) Index reduction
- (4) Matching
- (5) Block-lower triangular ordering
- (6) Tearing

Alias elimination and variability propagation are conceptually simple and serve to eliminate algebraic variables described by trivial equations. Alias elimination identifies variables occurring in equations of the form  $x \pm y = 0$ , which are ubiquitous in object-oriented acausal modelling, and eliminates one of them. Variability propagation identifies algebraic equations which are independent of time, such as  $y = p + 1$ , where  $p$  is a parameter, allowing the corresponding algebraic variables to be eliminated by solving the static equations.

Conceptually, an explicit ODE allows the computation of  $\dot{\mathbf{x}}(t)$  given the current known variables  $t, \mathbf{x}(t), \mathbf{u}(t)$ , and  $\mathbf{p}$ . After having performed index reduction as discussed in Section 3.1, the remaining steps of the ODE transformation problem can thus be considered equivalent to solving the square system

$$\mathbf{F}(\mathbf{z}; t, \mathbf{x}, \mathbf{u}, \mathbf{p}) = 0, \quad (14)$$

where  $\mathbf{z} = (\dot{\mathbf{x}}, \mathbf{y})$ . Note that we have now dropped the dependence on time for  $\dot{\mathbf{x}}, \mathbf{x}$ , and  $\mathbf{u}$ , since we with this perspective only consider them as elements of  $\mathbb{R}^n$  rather than functions.

While solving (14) can be done in a straightforward numerical manner by applying e.g. Newton's method<sup>2</sup>, such an approach may be inefficient and may also be practically challenging due to the need of having a sufficiently good initial guess of the solution for convergence. The idea of causalization and this paper is to apply symbolic transformations to the problem to allow for more efficient subsequent numerical solution.

The next step is to find a perfect matching between each scalar component of  $\mathbf{F}$  and each scalar component of  $\mathbf{z}$ , meaning that each equation is matched to a single variable and vice versa. Such a matching exists if and only if the DAE is low index. Finding perfect matchings is a well-studied graph theoretical problem with several available efficient algorithms. JModelica.org uses the Hopcroft-Karp [35] algorithm, which has the best known worst case performance and often performs well in practice.

The two final steps, block-lower triangular (BLT) ordering and tearing, are discussed below.

### 3.4.1 Block-lower triangular ordering

Consider the structural incidence matrix of the DAE, whose elements are defined by

$$\text{struct } F_{i,j} = \begin{cases} 0 \text{ or } *, & \text{if } \nabla_{z_j} F_i \equiv 0, \\ 1 \text{ or } *, & \text{if } \nabla_{z_j} F_i \not\equiv 0 \text{ and } \nabla_{z, z_j}^2 F_i \equiv \mathbf{0}, \\ *, & \text{otherwise,} \end{cases} \quad (15)$$

that is, 0 denotes incidences that do not depend on unknowns, 1 denotes incidences that depend affinely—henceforth called linearly—on unknowns, and  $*$  denotes any kind of incidence. The main step of the causalization procedure is to permute the DAE structural incidence matrix to a BLT form. The result is that the equations and variables of the DAE have been sorted so that the

<sup>2</sup>It is dubious whether the result really can be considered to be an explicit ODE, since a closed-form expression for  $\mathbf{f}$  has not been found (in fact, one may not even exist). However, typical numerical ODE solvers will not be able to tell the difference as long as  $\mathbf{f}$  can be evaluated.



DAE can be described by

$$\mathbf{F}^1(\mathbf{z}^1; \mathbf{v}^1) = \mathbf{0}, \quad (16a)$$

$$\mathbf{F}^2(\mathbf{z}^2; \mathbf{v}^2) = \mathbf{0}, \quad (16b)$$

$$\vdots$$

$$\mathbf{F}^m(\mathbf{z}^m; \mathbf{v}^m) = \mathbf{0}, \quad (16c)$$

where  $\mathbf{F}^i$  corresponds to a diagonal block (colloquially known as an algebraic loop if it is not scalar valued),  $m$  is the number of such blocks,  $\mathbf{z}^i$  is the unknown variables of  $\mathbf{F}^i$ , and

$$\mathbf{v}^i := (\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^{i-1}, t, \mathbf{x}, \mathbf{u}, \mathbf{p}) \quad (17)$$

is the known variables and unknown variables of preceding blocks. In other words, we have found permutation matrices  $\mathbf{P}$  and  $\mathbf{Q}$  such that

$$\mathbf{P} \text{struct}(\mathbf{F})\mathbf{Q} = \begin{bmatrix} \mathbf{F}^{1,1} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{F}^{2,1} & \mathbf{F}^{2,2} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{F}^{3,1} & \mathbf{F}^{3,2} & \mathbf{F}^{3,3} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{F}^{m,1} & \mathbf{F}^{m,2} & \mathbf{F}^{m,3} & \dots & \mathbf{F}^{m,m} \end{bmatrix}, \quad (18)$$

where  $\mathbf{F}^{i,j} := \text{struct} \mathbf{F}^i(\mathbf{z}^j; \mathbf{z}^1, \dots, \mathbf{z}^{j-1}, \mathbf{z}^{j+1}, \dots, \mathbf{z}^i, t, \mathbf{x}, \mathbf{u}, \mathbf{p})$ .

This form allows the sequential treatment of each diagonal block, allowing us to solve multiple small systems rather than a single large. Furthermore, it enables specialized treatment of each diagonal block, allowing the exploitation of equation structure that may exist within a diagonal block but not in the full system  $\mathbf{F}$ , such as linearity.

We are thus interested in the BLT form that has the maximal number of diagonal blocks. This form is unique in the sense that the number of diagonal blocks and their respective sizes and the variables and equations within a diagonal block are unique, but the ordering of diagonal blocks and orderings of variables and equations within the blocks are in general not unique. A trivial consequence of this form is that the diagonal blocks are irreducible.

Finding a BLT ordering that is optimal in this sense is equivalent to finding the strongly connected components of the directed graph defined by the matching found in step 4 as follows. For each matched equation-variable pair, create a vertex. For each non-zero struct  $F_{i,j}$ , create an edge from the vertex corresponding to equation  $i$  to the vertex corresponding to variable  $j$ . This will create loops on all of the vertices, which may be removed.

Tarjan's algorithm [56] is widely regarded as the most efficient algorithm for finding the strongly connected components of a directed graph, with linear complexity in the number of vertices and edges. It also has the added benefit of not only identifying the strongly connected components, but also topologically sorting them so that all edges out of a component lead to preceding components, giving us the sought BLT ordering.

Block-triangular orderings are useful in contexts other than DAE systems, such as solving sparse, unsymmetric linear equations [26].

### 3.4.2 Tearing

All that remains in the computation of  $\dot{\mathbf{x}}$  is the solution of each  $\mathbf{F}^i$ . In general this will require iterative numerical methods. In this paper we will however not proceed down this route at this

stage. We will instead only rely on symbolic techniques, without involving iterative methods, in our computation of  $\dot{\mathbf{x}}$ . Consequently, we will not go all the way to an explicit ODE. The details of this are discussed in Section 4. An important part of this is the use of tearing [8, 25, 27, 38]. Tearing is a method for solving sparse systems of equations that lack significant structure, such as being triangular or block diagonal, by ordering the variables and equations to get a partitioning of the system

$$\mathbf{0} = \mathbf{F}^i(\mathbf{z}^i; \mathbf{v}^i) = \begin{bmatrix} \bar{\mathbf{F}}^i(\bar{\mathbf{z}}^i; \hat{\mathbf{z}}^i, \mathbf{v}^i) \\ \hat{\mathbf{F}}^i(\hat{\mathbf{z}}^i; \bar{\mathbf{z}}^i, \mathbf{v}^i) \end{bmatrix} \quad (19)$$

such that the first partition  $\bar{\mathbf{F}}^i$  is highly structured, allowing this part of the system to be solved efficiently on its own and then utilized in the solution of the full system by block elimination.

Tearing can be done in different ways, for different applications, with different goals in mind. In this paper we apply tearing to the diagonal blocks of the BLT form. We seek a partitioning such that one part of the diagonal block is triangular and linear along the diagonal. That is, we find permutations  $\mathbf{P}^i$  and  $\mathbf{Q}^i$  so that

$$\mathbf{P}^i \text{ struct}(\mathbf{F}^i) \mathbf{Q}^i = \left[ \begin{array}{cc|c} & \bar{\mathbf{z}}^i & \hat{\mathbf{z}}^i \\ \hline 1 & & \\ & 1 & \mathbf{0} \\ & & \ddots \\ * & & 1 \\ \hline & * & * \end{array} \right] \begin{array}{l} \bar{\mathbf{F}}^i \\ \hat{\mathbf{F}}^i \end{array}, \quad (20)$$

where  $\bar{\mathbf{z}}^i$  is called the causalized variables,  $\bar{\mathbf{F}}^i$  the causalized equations (due to the lower-triangular structure),  $\hat{\mathbf{z}}^i$  the tearing variables, and  $\hat{\mathbf{F}}^i$  the tearing residuals. This form is appealing because it allows for symbolic elimination of the causalized variables in terms of the tearing variables; that is,

$$\bar{\mathbf{z}}^i = (\bar{\mathbf{F}}^i)^{-1}(\mathbf{0}; \hat{\mathbf{z}}^i, \mathbf{v}^i) =: \bar{\mathbf{H}}^i(\hat{\mathbf{z}}^i, \mathbf{v}^i). \quad (21)$$

The fact that  $\bar{\mathbf{F}}^i$  is triangular and linear along the diagonal allows for efficient and numerically stable computation of  $(\bar{\mathbf{F}}^i)^{-1}$  through forward substitution. The remaining equations to be solved for  $\hat{\mathbf{z}}^i$  numerically, typically using Newton's method, are given by

$$\bar{\mathbf{H}}^i(\hat{\mathbf{z}}^i) := \hat{\mathbf{F}}^i(\hat{\mathbf{z}}^i; \bar{\mathbf{H}}^i(\hat{\mathbf{z}}^i, \mathbf{v}^i), \mathbf{v}^i) = \mathbf{0}. \quad (22)$$

The benefits of this approach over solving the full diagonal block (19) numerically is that (22) has fewer variables and equations, equal to the number of tearing variables, often allowing it to be solved more efficiently. But a perhaps more important benefit is that a sufficiently good initial guess for convergence is only needed for  $\hat{\mathbf{z}}^i$ , whereas solving the full system would also require a sufficiently good guess for  $\bar{\mathbf{z}}^i$ . For these reasons, it is desirable to find a partitioning that minimizes the number of tearing variables and residuals. Unfortunately, this problem is NP-hard [8]. While in many cases it is tractable to solve the problem to optimality, it is common to instead

apply heuristics to find near-optimal partitionings. The tearing algorithm used in JModelica.org, and in this paper, is based on the heuristics described by Meijer [50].

The use of tearing is however a double-edged sword. While the symbolic solution of (21) and numerical solution of (22) are numerically stable, the tearing residuals (22) may be ill-conditioned even if the full system (19) is not. Furthermore, even if the tearing residuals and full system are well-conditioned, the numerical computation of  $\bar{z}^i$  through (21) after having computed the solution of (22) may be numerically unstable. Another potential drawback is that (22) is significantly more dense than (19), which may cause it to actually be more expensive to solve if sparsity is exploited in the computations. This topic is discussed further below. In conclusion, there is more to the choice of tearing variables and residuals than just minimizing their number.

### 3.5 Pivot selection in direct, sparse linear solvers

When solving a nonlinear system of equations with Newton's method, the symbolic elimination of variables is computationally equivalent to a priori selection of pivot elements in the linear equation solver in each Newton iteration. Since pivots are selected based on numerical values to prevent numerical instability (caused by error growth due to numerical round-off), a priori selection based entirely on structural, as opposed to numerical, information can lead to numerical instability, as mentioned in Section 3.4.2 and discussed by Duff et al. [25].

Another potential issue of blindly eliminating variables is that the reduced system may be more dense. As opposed to dense direct linear solvers, sparse direct linear solvers do not only select pivots to minimize error growth, but also to minimize fill-in. The typical approach is to select as pivot element the element that causes the least amount of estimated fill-in in the matrix factors while also being bigger (in magnitude) than a fraction of the largest element in the same column. The reason that only an estimate of the fill-in is used is due to the computational intractability of computing the fill-in caused on a global level (the final matrix factors). The most widely used estimate of fill-in is the Markowitz criterion [47]. When LU-factorizing a (sub)matrix  $\mathbf{M}$ , the Markowitz criterion selects as the next pivot the element  $M_{i,j}$  that minimizes

$$(\text{nnz } \mathbf{M}_{i,:} - 1) \cdot (\text{nnz } \mathbf{M}_{:,j} - 1), \quad (23)$$

typically out of those elements satisfying

$$|M_{i,j}| \geq u \max_l |M_{l,j}|, \quad (24)$$

where  $0 < u \leq 1$  is the pivot tolerance. Another estimate is that of local minimum fill-in, also proposed by Markowitz [47], which is a better but more expensive estimate. Local minimum fill-in computes the actual fill-in in each stage, whereas the Markowitz criterion estimates the fill-in by the amount of incidences in the pivot row and column, i.e., it does not take into account that some of the incidences already occur in the remaining equations and thus do not cause additional fill-in. The sparsity-preservation techniques of Section 4.3 are related to these ideas.

### 3.6 Related work

The standard approach of discretizing the full DAE relies on fill-reducing orderings, such as nested dissection [31], for efficient numerical linear algebra and numerical pivoting for stability. While our proposed approach of using causalization techniques has been used in the context of dynamic optimization before [7, 29], in this paper we consider the effects they have on the

solution procedure. We also do not complete the causalization all the way to an explicit ODE, but rather choose to keep some algebraic variables and implicit equations for efficiency.

Safdarnejad et. al [55] consider the use of BLT decompositions for the purpose of more efficiently generating initial guesses to (11) by solving square problems (with fixed degrees of freedom), and in particular identifying infeasibilities. They do however not consider BLT decompositions for the actual solution of (11). Fletcher [28] considers the use of block-triangular ordering and tearing for efficient and sparsity-preserving orderings for implicit LU factorization tailored for linear programming.

The main purpose of causalization is the efficient treatment of algebraic equations. There are other approaches to achieving this. One technique that is common when using direct multiple shooting is the elimination of all algebraic variables in the shooting nodes by linearizing the consistency conditions (algebraic equations) in each iteration, which is possible for semi-explicit index-one DAEs [23]. Another related method is the use of reduced space methods [21] which primarily work in the null space of the NLP equality constraints, which can be beneficial when there are few degrees of freedom, as is typically the case in dynamic optimization problems. Since these techniques are applied to the NLP rather than the dynamic problem (11), they could be combined with the ideas presented in this paper with potential benefits, but such possibilities are not considered further in this paper.

Another possibility when employing direct local collocation is the use of parallelization by exploiting the arrowhead structure of the KKT system that arises due to the temporal decoupling of the finite elements [59]. This has great potential speedups when there are significantly more algebraic than differential variables, a sufficiently large amount of finite elements is used, and a large number of processor cores are available. Unlike the other approaches outlined in this section, this approach is ill-suited to be combined with the ideas of this paper, as its efficiency is dependent on there being a large amount of algebraic variables exposed to the collocation discretization.

## 4. Symbolic elimination for dynamic optimization

In Section 3.4 we reviewed how implicit DAEs can be transformed to explicit ODEs. While almost all of the steps were symbolic in nature, in the end we were left with the algebraic loops  $F^i = 0$  (which can be further symbolically reduced through the use of tearing), which in general will require numerical iterative methods to solve. However, going all the way to an explicit ODE is not necessary when applying direct collocation in the context of dynamic optimization, as collocation methods are perfectly capable of dealing with low-index (and to some extent also high-index) implicit DAEs. In this section we describe how the causalization techniques can be applied to (11) to symbolically eliminate many of the algebraic variables.

We will describe various approaches to symbolic elimination which can be combined in different ways, giving rise to different *schemes* of symbolic elimination. In Section 5 we will compare the performance of these schemes to each other, and also to the trivial elimination scheme of not performing any eliminations at all, which can be considered to be the standard approach when employing direct collocation to dynamic optimization problems.

**SCHEME 0** *Do not eliminate any algebraic variables.*

### 4.1 Elimination in scalar, linear diagonal blocks

We start by applying the first 5 steps of the causalization procedure of Section 3.4 to the DAE in (11b) in the standard manner, and postpone the discussion on the use of tearing. We next identify

the diagonal blocks  $F^i$  of the BLT form that are both scalar and linear; that is, struct  $F^i = 1$ . The corresponding variables  $z^i$  of these blocks, which are scalar, are considered for elimination. By elimination we mean that all occurrences of  $z^i$  in (11) are substituted by  $(F^i)^{-1}(\mathbf{0}; v^i)$  using (16). Note that the computation of  $(F^i)^{-1}$  in this case is just division by a single, scalar closed-form expression<sup>3</sup>. After this elimination step, the variable  $z^i$  and the equation to which it has been matched are removed from the problem.

In a first effort, we choose to eliminate all such variables (whose matched equations are scalar and linear with respect to the variable) that are unbounded, that is, those for which the corresponding elements of  $L$  and  $U$  are  $-\infty$  and  $+\infty$ , respectively. No other variables are eliminated.

**SCHEME 1** *Eliminate those algebraic variables that are unbounded, belong to scalar diagonal blocks, and whose incidence in the diagonal block is linear.*

A consequence of eliminating bounded variables is that the linear bound is transformed into a nonlinear inequality constraint (it is moved from the left part of (11c) to the right part). When using IPOPT to solve the discretized version of (11), there is no obvious benefit in eliminating bounded variables. This is because IPOPT transforms each nonlinear inequality constraint into a nonlinear equality constraint and a variable bound, through the introduction of a slack variable. So even if we eliminate a bounded variable, IPOPT would essentially undo the elimination, leading to no difference in the dimension of the KKT system that is solved. But while the size of the KKT system is unaffected, the structure is not. One situation when the change in structure is potentially harmful is when the NLP functions are undefined if the algebraic variable bound is not satisfied. Since IPOPT iterates always are feasible with respect to the variable bounds, but not necessarily the nonlinear inequality constraints, eliminating bounded algebraic variables can thus cause evaluation errors which otherwise would not have occurred. For these reasons, we do not consider the elimination of bounded variables.

There is a problem with the approach of only eliminating unbounded variables, since in Optima, many of the variables are bounded. This is because bounds are defined by the `min` and `max` Modelica attributes, which variables inherit from their types. For example, all variables declared as temperatures are bounded below by 0 [K], even though these bounds will never be active at feasible points. Because of this, we force the user to specify which variables are expected to be actively bounded, and consider all other variables as unbounded for the purpose of elimination, even if they have finite `min` and `max` values.

We only consider the elimination of algebraic variables, and not differential variables. While the techniques we use to find closed-form expressions for computing algebraic variables also can be used to find closed-form expressions for the derivatives of differential variables (which indeed is what is done when doing full causalization to obtain an explicit ODE), finding these expressions is not sufficient for the elimination of the derivatives of differential variables. In order to eliminate the derivatives of differential variables, we would also have to eliminate the corresponding differential variables, which would require the solution of differential equations. Finding closed-form solutions to differential equations is rarely possible for the cases of interest, so this possibility is not considered further in this paper.

There is however another potential use of the closed-form expressions for the derivatives of the differential variables without eliminating them. All occurrences of a differential variable derivative, except the ones in the matched diagonal block equations, can be substituted by the closed-form expression. While this will never affect NLP size, it will affect the NLP sparsity in a way that depends on the simultaneous discretization method. The collocation algorithm in

<sup>3</sup>In the general case, care is needed to ensure that this scalar expression is non-zero [9, 50]. However, since this is always the case for the problems we consider in Section 5, we henceforth ignore this issue for the sake of simplicity.

JModelica.org creates NLP variables for all differential variable derivatives. So from a sparsity perspective it is never good to substitute the differential variable derivatives in JModelica.org. However, another common approach to direct collocation is to eliminate the differential variable derivatives by using the collocation equations. In this case it is not obvious whether such substitutions would be beneficial. Comparing the number of incidences from a sparsity perspective would be easy enough, but since the closed-form expressions found via the causalization tend to be highly nonlinear whereas the collocation equations are linear, a pure sparsity perspective is probably too narrow to be useful. But since this is a non-issue in JModelica.org, we do not consider this further.

While it is known that performing eliminations along the lines of SCHEME 1 can be beneficial for dynamic optimization, tools usually require such eliminations to be manually identified and performed by the user. Such work can be tedious and, depending on the modelling language, lead to convoluted and model code and inefficient computations. The techniques of Section 3.4 enable efficient automation of SCHEME 1.

## 4.2 Elimination in non-scalar diagonal blocks

We next consider the elimination of variables in non-scalar diagonal blocks. If the diagonal blocks are linear—that is, struct  $F^i \in \{0, 1\}^{n^i \times n^i}$ , where  $n^i$  is the dimension of  $z^i$ —a conceivable approach would be to invoke a numerical factorization algorithm to solve for  $z^i$ . This would however lead to crippling inefficiencies in the current dynamic optimization framework in JModelica.org<sup>4</sup>, and is thus not considered further in this paper.

Another possibility of treating linear, non-scalar blocks is symbolic factorization. While this is readily supported by JModelica.org, experiments on the problems in Section 5.1 have shown that symbolic QR factorization all too often leads to numerical issues caused by instability due to the lack of numerical pivoting for practical use, and is thus not considered further in this paper.

There is however another approach that allows us to eliminate some of the variables in non-scalar diagonal blocks, which even extends to nonlinear blocks: tearing. By tearing the diagonal blocks as described in Section 3.4.2, we can then symbolically eliminate the causalized variables in the torn blocks by forward substitution.

There are two issues with the choice of causalized variables by the JModelica.org compiler. The first is that it can choose to causalize differential variable derivatives, which is useful when the causalization goal is to get an explicit ODE. But as per the discussion in Section 4.1, we are not interested in eliminating differential variable derivatives. Likewise, we are not interested in eliminating (actively) bounded variables. Our approach is thus to besides using the tearing variables and residuals selected by the compiler, we additionally add all differential variable derivatives and bounded variables as tearing variables and their respectively matched equations as tearing residuals. It would have been better to force the compiler to make these choices, allowing it to potentially make better choices in choosing the remaining tearing variables and residuals, but this has not been implemented.

**SCHEME 2** *Eliminate those algebraic variables that are unbounded, belong to scalar diagonal blocks, and whose incidence in the diagonal block is linear. Also eliminate causalized, unbounded algebraic variables in torn, non-scalar diagonal blocks.*

<sup>4</sup>The current dynamic optimization framework in JModelica.org symbolically represents the DAE using CasADi MX (Matrix eXpression) or SX (Scalar eXpression) [4] graphs. The embedding of numerical factorization algorithms requires the use of MX graphs. Since representing the full DAE using MX graphs is highly inefficient in terms of computational speed, an efficient implementation of this would require a sophisticated mixture of SX and MX graphs to represent the DAE in JModelica.org, or an extension of CasADi's SX graphs to allow function calls. Both of these possibilities are beyond the scope of this paper.

### 4.3 Sparsity-preserving elimination

In our attempt to eliminate as many variables as possible using only symbolic techniques, we may actually end up eliminating too many from a computational efficiency perspective. As briefly mentioned before, elimination tends to yield more dense systems. The increased density can cause crippling slowdowns when performing sparse computations, in particular when solving sparse linear systems (the linearized KKT system in each iteration of an NLP solver).

We thus propose to avoid some of the eliminations performed by SCHEME 1 & 2, by analysing the effect each elimination has on the DAE structural incidence matrix in a seventh step of the causalization procedure presented in Section 3.4. We define a density measure  $\mu$  of a causalized block variable<sup>5</sup>  $\bar{z}_j^i$  which measures how much denser the resulting NLP will become if it is eliminated using (21). Unfortunately, any useful measure will depend on which other eliminations are performed. This makes it intractable to consider all possible combinations. However, due to the causalization, we can consider the blocks and the variables within the blocks in sequence to devise a greedy algorithm which only considers the current block. Thus, for a given block, whether to perform the eliminations in preceding blocks has already been decided, and we only consider the situation in which no eliminations are performed in succeeding blocks. This simplified approach is analogous to how linear solvers use the Markowitz criterion or local minimum fill-in as discussed in Section 3.5.

We thus want to estimate how much the number of nonzero elements in the KKT system will increase if  $\bar{z}_j^i$  is eliminated. We perform the elimination if

$$\mu(i, j) \leq \mu_{\text{tol}}, \quad (25)$$

where  $\mu_{\text{tol}}$  is a user-provided tolerance. Rather than considering the sparsity of the KKT matrix, we instead only consider the number of nonzeros in the structural incidence matrix of the DAE—which typically is the dominant part of the KKT system in dynamic optimization—for simplicity and computational efficiency during symbolic preprocessing. An estimative measure for this is

$$\mu(i, j) = \left( -2 + \sum_{\alpha \in \{\bar{x}, x, y, u, p\}} \sum_{k=1}^{n_\alpha} \mathcal{I}(\nabla_{\alpha_k} \bar{F}_j^i) \cdot d(\alpha_k) \right) \cdot \left( -1 + \text{nnz } \nabla_{\bar{z}_j^i} \mathbf{F} \right), \quad (26)$$

where  $\mathcal{I}$  is the indicator function (mapping zero functions to 0 and all other functions to 1),  $d(\alpha_k)$  is the number of post-elimination dependencies of  $\alpha_k$ , and  $\text{nnz } \nabla_{\bar{z}_j^i} \mathbf{F}$  is the number of nonzero incidences of  $\bar{z}_j^i$  in the DAE. A recursive expression for  $d$  is

$$d(\alpha_k) = \begin{cases} 1, & \text{if } \alpha_k \text{ is not eliminated,} \\ \sum_{\beta \in \{\bar{x}, x, y, u, p\}} \sum_{l=1}^{n_\beta} \mathcal{I}(\nabla_{\beta_l} \bar{h}_{\alpha_k}) \cdot d(\beta_l), & \text{otherwise,} \end{cases} \quad (27)$$

where  $\bar{h}_{\bar{z}_j^i} = \bar{H}_j^i$  (see (21)). Note that  $d(\alpha_k)$  not only depends on  $\alpha_k$  but also which variables have been chosen for elimination, which changes over the course of the sparsity analysis. The measure (26) is similar to the Markowitz criterion (23), with some differences. It considers the whole system simultaneously rather than a single stage of Gaussian elimination, and also takes into account that the full system not only depends on  $\mathbf{z}$  but also  $\mathbf{v}$ . Another difference is that the Markowitz criterion estimates fill-in, whereas (26) estimates the increase in number of nonzeros.

<sup>5</sup>To allow uniform treatment of scalar, linear and torn diagonal blocks, we consider a scalar, linear diagonal block equation and its corresponding variable to be causalized.

For example, the elimination of an alias variable can cause a lot of fill-in, but will not increase the number of nonzeros in the DAE. Thus, the value given by (23) of an alias variables is equal to the number of incidences in the subsequent equations, whereas the value given by (26) is 0.

To demonstrate (25)–(27), consider the torn incidence (8) of the previous example with  $\mu_{\text{tol}} = 3$ . The variables that are eligible for elimination are  $y_4, y_1$ , and  $y_2$ . At the start of the procedure, we have  $d(\alpha) = 1$  for all  $\alpha$ . Since the first block  $F^1$  has no eligible variables, we move on to  $F^2$  where we find the pair  $((1e), y_4)$ —noting that  $\bar{F}_1^2$  is the residual for (1e)—for which we compute

$$\begin{aligned} \mu(2, 1) &= \left( -2 + \sum_{\alpha \in \{\bar{x}_1, x_1, y\}} \sum_{k=1}^{n_\alpha} \mathcal{I}(\nabla_{\alpha_k} \bar{F}_1^2) \cdot d(\alpha_k) \right) \cdot (-1 + \text{nnz } \nabla_{\bar{z}_1^2} F) \\ &= \left( -2 + \sum_{\alpha \in \{\bar{x}_1, x_1, y_1, y_2, y_3, y_4, y_5\}} \mathcal{I}(\nabla_{\alpha} \bar{F}_1^2) \cdot d(\alpha) \right) \cdot (-1 + \text{nnz } \nabla_{y_4} F) \\ &= \left( -2 + \mathcal{I}(1) \cdot d(y_4) + \mathcal{I}\left(-\frac{1}{2\sqrt{y_5}}\right) \cdot d(y_5) \right) \cdot (-2 + \text{nnz}(0, 0, 2y_1y_2, y_1 - 1, 1, 0)), \\ &= (-2 + 2) \cdot (-1 + 3) = 0 \leq \mu_{\text{tol}} = 3, \end{aligned} \quad (28)$$

and thus proceed to eliminate  $y_4$  and compute  $d(y_4) = 1$ . In the next block we find the eligible pairs  $((1d), y_1)$  and  $((1b), y_2)$ . For the first pair we compute

$$\mu(3, 1) = (-2 + 4) \cdot (3 - 1) = 4 > \mu_{\text{tol}}, \quad (29)$$

and thus do not eliminate  $y_1$ . Moving on to the second and final pair, we compute

$$\mu(3, 2) = (-2 + 3) \cdot (3 - 1) = 2 \leq \mu_{\text{tol}}, \quad (30)$$

and thus eliminate  $y_2$ , for which we compute  $d(y_2) = 2$ .

The measure (26) is locally suboptimal in the same sense as the Markowitz criterion: It does not take into account that not all incidences will cause fill-in; it is a worst case estimate. To address this, we can instead of (26) use the measure

$$\mu(i, j) = \sum_{l=i}^m \sum_{j=\begin{cases} j+1, & \text{if } l=i, \\ 1, & \text{otherwise} \end{cases}}^{n^l} \mathcal{I}(\nabla_{\bar{z}_j^l} F_j^l) \left( -1 + \sum_{\alpha \in \{\bar{x}, x, y, u, p\}} \sum_{k=1}^{n_\alpha} (1 - \mathcal{I}(\nabla_{\alpha_k} F_j^l)) \mathcal{I}(\nabla_{\alpha_k} F_j^l) d(\alpha_k) \right), \quad (31)$$

where  $F_j^i$  is element  $j$  in block  $i$  with the ordering obtained after tearing; that is,

$$F_j^i = \begin{cases} \bar{F}_j^i, & \text{if } j \leq \bar{n}_i, \\ \bar{F}_{j-\bar{n}_i}^i, & \text{if } j > \bar{n}_i. \end{cases} \quad (32)$$

This measure is analogous to local minimum fill-in in the same way as (26) is analogous to (23).

As discussed in Section 3.5, Markowitz is generally regarded as the best criterion for general-purpose pivot selection, with local minimum fill-in being a lot more computationally expensive while usually only yielding slightly better results. However, for our purposes, we only perform the sparsity-preservation analysis once offline. Furthermore, the difference in computation times between using (26) and (31) is usually negligible compared to the other offline computations (model compilation, BLT analysis, collocation discretization, algorithmic differentiation graph



construction, and so on). So while both measures have been implemented in JModelica.org, we will henceforth only consider the usage of (31) due to its slightly better expected performance.

The usage of sparsity preservation can be used both with and without tearing, yielding two new families of schemes, parametrized by the density tolerance  $\mu_{\text{tol}}$ .

**SCHEME 3 $_{\mu_{\text{tol}}}$**  *Eliminate those algebraic variables that are unbounded, belong to scalar diagonal blocks, whose incidence in the diagonal block is linear, and whose density measure is smaller than or equal to  $\mu_{\text{tol}}$ .*

**SCHEME 4 $_{\mu_{\text{tol}}}$**  *Eliminate those algebraic variables that are unbounded, belong to scalar diagonal blocks, whose incidence in the diagonal block is linear, and whose density measure is smaller than or equal to  $\mu_{\text{tol}}$ . Also eliminate causalized, unbounded algebraic variables in torn, non-scalar diagonal blocks whose density measure is smaller than or equal to  $\mu_{\text{tol}}$ .*

The choice of density tolerance is important. A tolerance of 0 means that we only perform eliminations that do not increase the number of nonzeros in the other equations. A tolerance of  $-\infty$  means that no eliminations are performed at all, and we thus preserve the original DAE. A tolerance of  $\infty$  means that we eliminate all causalized variables. We thus get **SCHEME 0** = **SCHEME 3 $_{-\infty}$**  = **SCHEME 4 $_{-\infty}$** , **SCHEME 1** = **SCHEME 3 $_{\infty}$** , and **SCHEME 2** = **SCHEME 4 $_{\infty}$** .

The schemes to be evaluated in the next section are summarized in Table 1.

Table 1. Summary of techniques used in schemes.

SCHEME	BLT	Tearing	Sparsity preservation
0			
1	✓		
2	✓	✓	
3 $_{\mu_{\text{tol}}}$	✓		✓
4 $_{\mu_{\text{tol}}}$	✓	✓	✓

## 5. Benchmark

To evaluate the various proposed schemes, we will consider six different optimal control problems. Since six problems is too small a test suite to draw any general conclusions regarding computational speed, and in particular convergence robustness, we will generate a larger number of problem instances by randomizing the initial state. We will primarily focus on local collocation, but will in the end also briefly consider global collocation. Code for reproducing most of the results in this section is available in [44].

### 5.1 Problems

The systems considered in the six problems are a car, a combined-cycle power plant (CCPP), a double pendulum, a kinematic loop with four bars, a distillation column, and a heat recovery steam generator (HRSG). All of the problems are encoded in Modelica and Optimica, and all of them except the HRSG (whose model is proprietary) are available as a part of JModelica.org's suite of examples. This section briefly describes each problem. [43, Chapter 5] contains a more detailed description of each problem.

#### 5.1.1 Car

The problem is to find the time-minimal manoeuvre for a car in a 90-degree turn with a high initial velocity, making it challenging to simply stay on the road. A single-track chassis model is

used, where the two wheels on each axle are lumped together. This problem, and several other similar ones, has been developed and studied by Berntorp et al. [11, 12, 46].

Unlike most of the other problems, the Modelica code for this problem is atypical in its being flat, rather than hierarchically object-oriented, leading to it consisting of relatively few lines of code and not very many algebraic variables and no trivial algebraic equations resulting from component connections. It thus serves to demonstrate the effects of symbolic elimination on more typical representations of DAEs, rather than the verbose equations that result when using Modelica model libraries. It is also the only problem we consider that has a free time horizon.

### 5.1.2 CCPP

The next case considers warm startup of combined-cycle power plants (CCPP). This problem has become highly industrially relevant during the last decade, due to an increasing need to improve power-generation flexibility due to the unpredictable output of renewable energy sources such as solar and wind power. The model was first developed by Casella et al. [19] but has since been further developed. The objective is a quadratic penalty from reference values, with an important upper bound on the turbine thermal stress.

This problem is the closest thing to a standard problem of dynamic optimization in the Modelica community, having been used several times for benchmark purposes [5, 6, 40, 51, 59].

### 5.1.3 Double pendulum

One of the elementary examples of multibody mechanics in the Modelica Standard Library (MSL) is a damped double pendulum: `Modelica.Mechanics.Examples.Elementary.DoublePendulum`. While a double pendulum is a conceptually simple fourth-order system, the MSL model is an index-three DAE with more than a hundred variables and an algebraic loop and is thus computationally non-trivial. This model is used to formulate the optimal control problem of inverting the full pendulum with quadratic penalties on the states and torque.

### 5.1.4 Fourbar1

Another high-index example from MSL multibody mechanics is `Modelica.Mechanics.MultiBody.Examples.Loops.Fourbar1`, which consists of 4 bars connected by 6 revolute joints and a prismatic joint. This system has three large algebraic loops, one of which is nonlinear. Based on this model, we consider the optimization problem to control the translation along the prismatic joint `j2` in one end of the mechanism by applying a bounded torque to the revolute joint `j1` in the other end by minimizing a quadratic cost function.

### 5.1.5 HRSG

The next problem concerns startup of a Heat Recovery Steam Generator (HRSG). The model has been developed and studied as part of an industrial research project and a series of industrial Master's theses [1, 54, 57]. While the model is moderate in size, it contains high-fidelity modelling of thermodynamic properties of media, causing high NLP function evaluation times. The objective is a quadratic penalty on the deviation from the desired steady state, with bounds on input rates and temperature gradients in spatially discretized walls.

The torn BLT decomposition obtained with SCHEME 4<sub>5</sub> is shown in Figure 1. The BLT-ordered incidence in Figure 1 is typical for Modelica models in that it is sparse, has a large amount of scalar, linear diagonal blocks, and a small amount of large algebraic loops with mostly linear incidences.

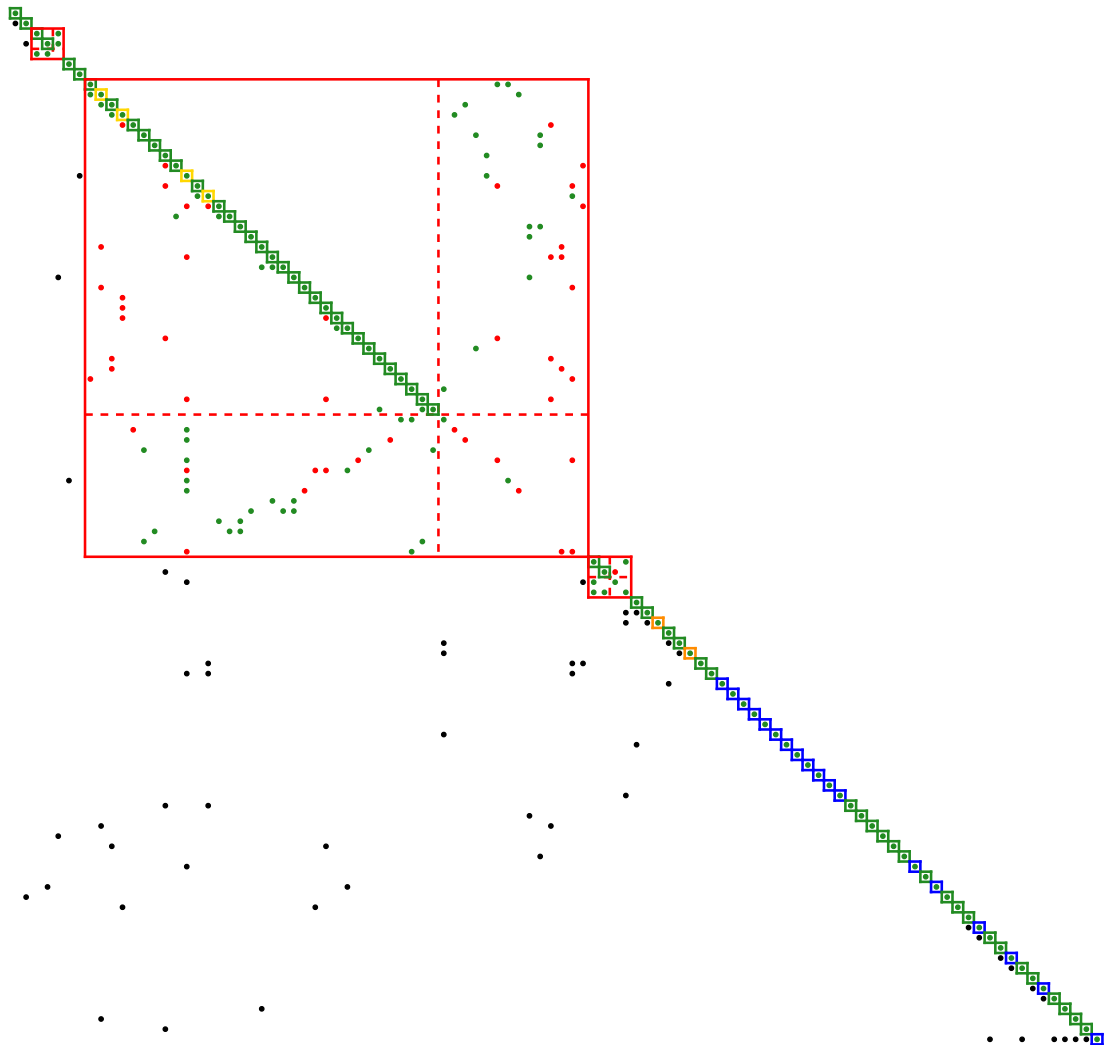


Figure 1. Torn, sparsity-preserving BLT decomposition of HRSG model. Linear incidences are marked by green dots, and nonlinear incidences are marked by red dots. Since we do not distinguish between linear and nonlinear incidences outside of the diagonal blocks, such incidences are marked by black dots. Torn blocks are marked by red edges. Variables, and their respective matched equations, that have been user-specified as actively bounded (and thus are not eliminated) are marked by orange edges. Differential variable derivatives (which are not eliminated), and their respective matched equations, are marked by blue edges. Variable-equation pairs along the diagonal that are not sparsity preserving—that is, do not satisfy (25)—are marked by yellow edges. The remaining variable-equation pairs along the diagonal are the ones used for elimination, which are marked by green edges.

### 5.1.6 Distillation column

The final problem concerns optimal control of a binary distillation column, which separates methanol from n-propanol and has 40 trays. The model was developed by Diehl [22] and the Modelica implementation was based on the MATLAB implementation by Hedengren [34]. The considered scenario is a short reflux breakdown during steady state, with the objective to steer back to the desired steady state, using quadratic costs on the deviation of two tray temperatures and input signals from the high-purity steady state. The Modelica implementation of this scenario and model has been previously used for benchmarking dynamic optimization algorithms [40, 45]. However, the model was then incorrectly implemented, leading to 125 differential variables, instead of the correct 40. We here use a correct implementation. Just like the Car model of Section 5.1.1, this model is implemented in a flat rather than hierarchical manner.

## 5.2 Benchmark setup

To generate a large test suite, we start with the solution to each of the six *nominal* problems described in Section 5.1. We then randomly perturb the initial state  $\mathbf{x}_0$  and solve the perturbed problem using the nominal solution as initial guess. The purpose of this approach is to emulate the setting of Model Predictive Control (MPC) [42], without actually doing MPC, which would have had the drawback of introducing correlation between the problems solved in each sample point.

The nominal value of each initial state variable  $x_{0,i}$  is multiplicatively perturbed by independently, identically, normally distributed random variables with standard deviation  $\sigma$  to yield the new initial state variable  $\bar{x}_{0,i}$ , that is

$$\bar{x}_{0,i} = \nu_i \cdot x_{0,i}, \quad \nu_i \sim \mathcal{N}(1, \sigma^2), \quad i = 1 \dots n_x. \quad (33)$$

The standard deviation  $\sigma$  is hand-picked for each problem to make the corresponding instances suitably difficult. The resulting perturbed problem, henceforth referred to as an *instance*, may be infeasible. To counteract this, we make sure that the initial state satisfies all the problem bounds and path inequality constraints. To satisfy the bounds on state variables, we project  $\bar{x}_{0,i}$  inside its bounds, and thus modify the first equation of (33) to

$$\bar{x}_{0,i} = \max(\min(\nu_i \cdot x_{0,i}, x_{0,i} + 0.9(x_{U,i} - x_{0,i})), x_{0,i} - 0.9(x_{0,i} - x_{L,i})), \quad (34)$$

where  $x_{U,i}$  is the element of the bounds  $\mathbf{U}$  that corresponds to  $x_i$ .<sup>6</sup> Projecting inside the feasible region of general path inequality constraints and algebraic variable bounds is more difficult. To satisfy these, we use JModelica.org's FMI-based DAE initialization algorithm to compute the value of  $\bar{\mathbf{y}}_0$  that corresponds to  $\bar{\mathbf{x}}_0$  and then check if they satisfy the constraints with a safety margin analogous to the one in (34); that is, the projected distance to the boundary should decrease by no more than 90%. If they do not, we discard the problem instance and generate a new one and repeat.

We solve each instance with all schemes and use the values  $\{5, 10, 20, 30, 40\}$  for  $\mu_{\text{tol}}$ . The problem sizes are detailed in Table 2, where we also specify the number of elements  $n_e$  and collocation points  $n_c$  used in each element. Each scheme is only allowed a certain CPU time  $t_{\text{max}}$  [s] amount of time for each problem instance, after which we regard the scheme as having failed. The time  $t_{\text{max}}$  has been chosen as approximately the average observed solution time of the slowest scheme over all instances for a given problem plus 5 estimated standard deviations. For some problems, some schemes yield identical results, as indicated in the table. For example, Car has no algebraic loops, and so SCHEME 1 and SCHEME 2 are the same for this problem.

All problems—except Distillation column—are solved with JModelica.org revision [8915], IPOPT 3.12.5, and the linear solver MA57 [36] with ordering by MeTiS [37]. The Distillation column problem is instead solved with JModelica.org revision [9153] and IPOPT 3.12.6, in order to use the correct model implementation discussed in Section 5.1.6. The acceptable NLP tolerance in IPOPT is set equal to the NLP tolerance of  $10^{-8}$ . To put a focus on convergence robustness rather than speed, the automatic scaling of MA57 is enabled and the pivot tolerance (see (24)) is increased from  $10^{-8}$  to  $10^{-4}$ . The IPOPT barrier parameter strategy is changed to adaptive, to avoid the issue of selecting the initial value of the barrier parameter and also because the authors favour this strategy. Other than the above exceptions, all the default values of the options of all the algorithms are used.

<sup>6</sup>The factor 0.9 in (34) is used instead of 1.0 in order to project strictly inside the interior of the bounds, which is needed to satisfy the Linear Independence Constraint Qualification (LICQ) as the collocation algorithm in JModelica.org introduces  $\mathbf{x}_0$  as an optimization variable in order to treat general implicit initial equations.

Table 2. Benchmark problems and the considered schemes for each problem, where  $t_{\max}$  [s] is the allotted CPU time for each instance,  $n_y$  the number of non-eliminated algebraic variables,  $n$  the number of NLP variables divided by 1000,  $\text{nnz } J$  the number of nonzero elements in the NLP constraint Jacobian divided by 1000,  $\text{nnz } H$  the number of nonzero elements in the Hessian of the NLP Lagrangian divided by 1000.

Problem	$n_x$	$n_u$	$n_e$	$n_c$	$\sigma$	$t_{\max}$	Schemes	$n_y$	$n$	$\text{nnz } J$	$\text{nnz } H$
Car	13	3	60	3	0.1	30	0	23	9.7	37.4	9.4
							1, 2, 3 <sub>20</sub> , 3 <sub>30</sub> , 3 <sub>40</sub> , 4 <sub>20</sub> , 4 <sub>30</sub> , 4 <sub>40</sub>	4	6.3	30.4	10.0
							3 <sub>5</sub> , 3 <sub>10</sub> , 4 <sub>5</sub> , 4 <sub>10</sub>	5	6.4	30.0	9.2
CCPP	10	1	40	4	0.3	40	0	123	23.6	73	11.6
							1, 3 <sub>10</sub> , 3 <sub>20</sub> , 3 <sub>30</sub> , 3 <sub>40</sub>	3	4.3	19.7	4.2
							2, 4 <sub>10</sub> , 4 <sub>20</sub> , 4 <sub>30</sub> , 4 <sub>40</sub>	2	4.1	19.3	4.2
							3 <sub>5</sub>	6	4.7	21.3	4.2
							4 <sub>5</sub>	5	4.6	20.1	4.2
Dbl. pend.	4	1	100	3	0.3	50	0	124	40.4	123	24
							1, 3 <sub>10</sub> , 3 <sub>20</sub> , 3 <sub>30</sub> , 3 <sub>40</sub>	16	7.9	25	5
							2, 4 <sub>20</sub> , 4 <sub>30</sub> , 3 <sub>40</sub>	2	3.7	13	5
							3 <sub>5</sub>	17	8.2	27	5
							4 <sub>10</sub>	3	4.0	15	5
							4 <sub>5</sub>	5	4.6	17	5
Fourbar1	2	1	60	3	0.03	30	0	452	82.8	334	192
							1, 3 <sub>30</sub> , 3 <sub>40</sub>	246	45.5	215	160
							2	23	5.2	50	24
							3 <sub>20</sub>	247	45.7	216	161
							3 <sub>10</sub>	249	46.1	217	161
							3 <sub>5</sub>	255	47.2	220	162
							4 <sub>40</sub>	29	6.3	63	29
							4 <sub>30</sub>	30	6.5	63	32
							4 <sub>20</sub>	46	9.3	88	58
							4 <sub>10</sub>	85	16.4	127	107
							4 <sub>5</sub>	114	21.7	142	128
HRSG	18	3	25	5	0.3	60	0	84	15.9	56	6
							1, 3 <sub>5</sub> , 3 <sub>10</sub> , 3 <sub>20</sub> , 3 <sub>30</sub> , 3 <sub>40</sub>	56	12.4	47	6
							2, 4 <sub>30</sub> , 4 <sub>40</sub>	19	7.7	42	7
							4 <sub>20</sub>	20	7.9	41	6
							4 <sub>10</sub>	22	8.1	39	5
							4 <sub>5</sub>	23	8.2	39	5
Dist. col.	40	2	20	3	0.2	40	0	1083	72.1	258	96
							1	291	23.8	223	153
							2	85	11.2	269	238
							3 <sub>40</sub>	294	24.0	148	75
							3 <sub>30</sub>	296	24.1	143	68
							3 <sub>20</sub>	298	24.2	136	62
							3 <sub>10</sub>	305	24.6	132	57
							3 <sub>5</sub>	331	26.2	130	53
							4 <sub>40</sub>	90	11.5	128	83
							4 <sub>30</sub>	93	11.7	114	66
							4 <sub>20</sub>	97	12.0	102	55
							4 <sub>10</sub>	108	12.6	69	38
							4 <sub>5</sub>	138	14.4	86	33

### 5.3 Results

In this section we present the results of the benchmark described in Sections 5.1 and 5.2. We generate 1000 instances whose initial states are in the interior of the constraints for each problem. We will use the same table headers throughout this section, with the following meanings. A scheme is considered to have *succeeded* on a problem instance if it converges to within tolerance of a solution within the maximum CPU time. A problem instance is considered *valid* if at least one scheme succeeds on it. Success [1] is the ratio of success of a scheme on the valid instances. For the instances on which all schemes succeeded, Time [s] is the average solution CPU time,  $\sigma_t$  [s] is the sample standard deviation of the solution time, and Iter [1] is the average number of iterations needed by a scheme.

#### 5.3.1 Car

The results for 1000 instances of the Car problem are shown in Table 3. On 68.0% of the instances, all schemes succeed. 18.6% of the instances are invalid. The most significant difference between the schemes is that SCHEME 0 is approximately 35% slower than the others and slightly less robust. The dominant reason for failure is maximum CPU time, with some cases of convergence to a point of local infeasibility. A likely reason for the high percentage of invalid instances is that the nominal solution is close to the boundaries of feasibility, due to the high initial car velocity. The perturbed problems can thus end up being infeasible with significant probability, which is also why we chose the relatively small value of  $\sigma = 0.1$ .

Table 3. Results for Car.

SCHEME	Success	Time	$\sigma_t$	Iter
0	89.8%	8.9	4.5	104.0
1	97.2%	6.8	4.9	94.3
3 <sub>5</sub>	95.2%	6.4	5.0	103.9

#### 5.3.2 CCPP

The results for 1000 instances of the CCPP problem are shown in Table 4. On 70.1% of the instances, all schemes succeed. 20.9% of the instances are invalid. We see that SCHEME 0 is an order of magnitude slower and also sometimes fails unlike the other schemes. There is little difference between the other schemes, which is to be expected, since tearing and sparsity preservation only makes minor changes to the problem as seen in Table 2. On almost all of the instances in which all schemes fail, all schemes except 0 report local infeasibility. SCHEME 0 only fails due to maximum CPU time.

Table 4. Results for CCPP.

SCHEME	Success	Time	$\sigma_t$	Iter
0	88.6%	13.6	5.3	101.9
1	100.0%	0.9	0.6	46.3
2	100.0%	0.9	0.5	46.4
3 <sub>5</sub>	100.0%	1.0	0.7	46.4
4 <sub>5</sub>	99.9%	1.0	0.6	46.5

#### 5.3.3 Double pendulum

The results for 1000 instances of the Double pendulum problem are shown in Table 5. On 99.6% of the instances, all schemes succeed. All instances are valid. We see that SCHEME 1 reduces

the solution time by an order of magnitude compared to SCHEME 0, which is further halved by SCHEME 3. The overzealous sparsity preservation that is needed to have an impact on the number of variables is not beneficial.

Table 5. Results for Double pendulum.

SCHEME	Success	Time	$\sigma_t$	Iter
0	99.6%	15.5	7.0	99.9
1	100.0%	2.1	1.0	72.4
2	100.0%	0.8	0.4	58.3
3 <sub>5</sub>	100.0%	2.2	1.0	71.1
4 <sub>10</sub>	100.0%	0.9	0.4	58.7
4 <sub>5</sub>	100.0%	1.0	0.5	58.9

#### 5.3.4 Fourbar1

The results for 1000 instances of the Fourbar1 problem are shown in Table 6. On 56.1% of the instances, all schemes succeed. 34.0% of the instances are invalid. The dominant reasons for failure are restoration failure in IPOPT and maximum CPU time. We see that SCHEME 2 performs better than all the other schemes, and that despite significant density in the problem after applying tearing, sparsity preservation actually does more harm than good. One might have suspected this when inspecting Table 2, as Fourbar1 is the only problem where sparsity preservation significantly increases nnz  $J$  rather than decrease it. Also, the Lagrangian Hessian is denser for this problem than the others, which we have neglected in the sparsity preservation.

Table 6. Results for Fourbar1.

SCHEME	Success	Time	$\sigma_t$	Iter
0	86.8%	8.9	2.4	15.6
1	87.1%	4.6	1.3	15.4
2	99.8%	1.4	0.4	15.2
3 <sub>20</sub>	87.1%	4.6	1.3	15.4
3 <sub>10</sub>	87.1%	4.6	1.3	15.4
3 <sub>5</sub>	86.8%	5.0	1.4	15.4
4 <sub>40</sub>	93.5%	1.6	0.5	15.3
4 <sub>30</sub>	95.2%	1.5	0.4	15.2
4 <sub>20</sub>	90.8%	1.7	0.5	15.4
4 <sub>10</sub>	85.8%	2.3	0.7	15.4
4 <sub>5</sub>	95.6%	3.2	0.9	15.4

#### 5.3.5 HRSG

The results for 1000 instances of the HRSG problem are shown in Table 7. On 72.2% of the instances, all schemes succeed. 23.1% of the instances are invalid. The dominant reasons for failure are local infeasibility and maximum CPU time. We see that SCHEME 1 only offers a slight improvement over SCHEME 0, but the use of tearing and sparsity preservation improves both robustness and solution times.

Table 7. Results for HRSG.

SCHEME	Success	Time	$\sigma_t$	Iter
0	96.0%	13.5	8.2	66.0
1	96.9%	11.1	7.2	68.8
2	99.3%	8.9	4.1	48.5
4 <sub>20</sub>	99.2%	7.7	4.4	49.8
4 <sub>10</sub>	99.1%	6.7	4.0	49.3
4 <sub>5</sub>	99.2%	6.8	4.2	49.7

### 5.3.6 Distillation column

The state variable  $\mathbf{x}$  for this problem consists entirely of molar fractions. Although  $\mathbf{x}$  is formally unbounded, each component lies in the interval  $[0, 1]$  by definition. For the purposes of (34), we thus impose the bounds  $\mathbf{0} \leq \mathbf{x} \leq \mathbf{1}$  to adhere to physics. These bounds are not enforced during the numerical solution of the optimization problem. SCHEME 0 is unable to solve a single instance of this problem. The results for the remaining schemes for 1000 instances of the Distillation column problem are thus shown in Table 8. On 97.2% of the instances, all remaining schemes succeed. All of the instances are valid. On the valid instances, the sole reason for failure is maximum CPU time. While tearing slightly improves robustness, it also slightly increases the solution time in the absence of sparsity preservation. Sparsity preservation however greatly reduces the solution time, especially in combination with tearing.

Table 8. Results for Distillation column.

SCHEME 0 failed all instances.

SCHEME	Success	Time	$\sigma_t$	Iter
1	98.8%	10.2	1.5	16.3
2	99.9%	10.7	1.0	14.3
3 <sub>40</sub>	99.5%	4.3	1.2	16.5
3 <sub>30</sub>	99.4%	4.7	0.5	16.2
3 <sub>20</sub>	99.4%	6.0	1.0	16.3
3 <sub>10</sub>	99.3%	4.1	1.1	16.3
3 <sub>5</sub>	99.6%	4.8	0.6	15.8
4 <sub>40</sub>	100.0%	3.6	0.3	14.4
4 <sub>30</sub>	100.0%	3.0	0.3	14.5
4 <sub>20</sub>	100.0%	2.6	0.2	14.4
4 <sub>10</sub>	99.9%	2.2	0.2	14.4
4 <sub>5</sub>	100.0%	3.4	0.3	14.2

## 5.4 Performance profiles

In Section 5.3 we saw that Scheme 1 on average outperforms Scheme 0 for each considered problem, and that Schemes 2, 3 <sub>$\mu_{\text{tol}}$</sub> , and 4 <sub>$\mu_{\text{tol}}$</sub>  usually yields further improvements. To illustrate the aggregated results, Figure 2a shows the performance profile [24] for all schemes on the valid portion of the 6000 problem instances. The performance  $\rho_s(\tau)$  of a scheme  $s$  is defined as the ratio of instances in which  $s$  solved the problem no slower than a factor  $\tau$  of the fastest scheme of that instance. In particular,  $\rho_s(1)$  is the ratio of instances in which  $s$  was the fastest scheme, and  $\rho_s(\infty)$  is the ratio of instances in which  $s$  succeeded.

Without looking too closely at the overwhelming data of Figure 2a, we can see that SCHEME 0 is outperformed by SCHEME 1, which is outperformed by both SCHEME 2 and SCHEME 3 <sub>$\mu_{\text{tol}}$</sub> , both of which are outperformed by SCHEME 4 <sub>$\mu_{\text{tol}}$</sub> . There is however no clear best value for  $\mu_{\text{tol}}$  in Figure 2a, and average performance as a function of  $\mu_{\text{tol}}$  is far from convex, suggesting that a larger suite of problems than that of Section 5.1 is needed to pinpoint a suitable default value for  $\mu_{\text{tol}}$ . Suitable values however seem to lie in the range of  $[10, 30]$ , and so 15 is chosen as the default value of JModelica.org.

For a clearer comparison of the various schemes without excessive focus on  $\mu_{\text{tol}}$ , we generate a new performance profile where we only consider  $\mu_{\text{tol}} = 30$ , which was the density tolerance yielding the best robustness. The result is shown in Figure 2b, where we clearly see the average superiority of SCHEME 4 both in terms of speed and robustness, although SCHEME 2 actually has slightly better robustness. In particular, we see that SCHEME 4<sub>30</sub> is an order of magnitude faster than SCHEME 0 in approximately half of the considered instances. Fourbar1 was the only problem where sparsity preservation significantly deteriorated robustness, which leads to SCHEME 2 being the most robust scheme.



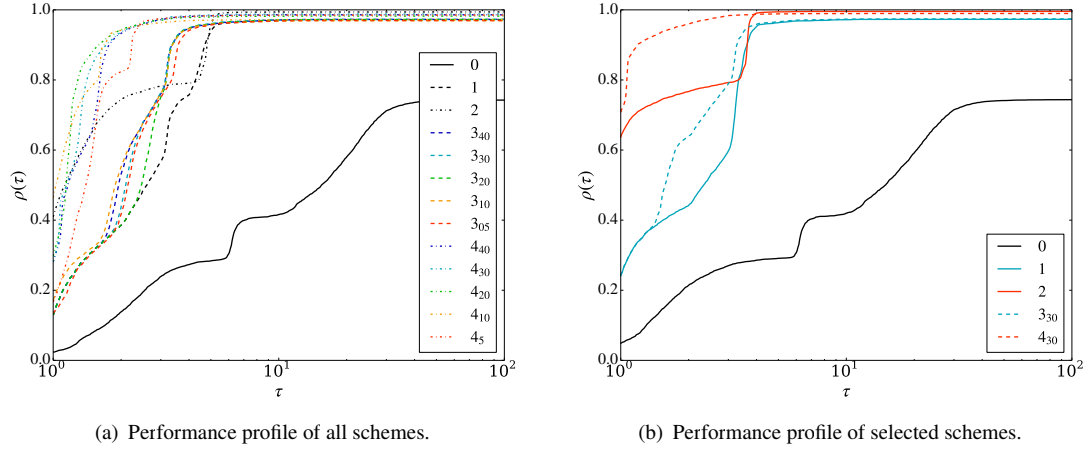


Figure 2. Performance profiles for the valid portion of 6000 problem instances. Subfigure (a) shows all considered schemes with 5 different values of  $\mu_{\text{tol}}$ . Subfigure (b) shows the same results but only for  $\mu_{\text{tol}} = 30$ .

### 5.5 Computation times

Regarding computation times, we have so far focused on the respective solution times of the schemes. For many applications, such as online MPC, these are the only times that matter. For other applications, the computation time of the full toolchain is important, in which case the proposed techniques add additional computational steps during preprocessing compared to SCHEME 0. The main parts of the offline computation times are model compilation, BLT analysis, collocation discretization, algorithmic differentiation graph construction (including first- and second-order derivatives). The main parts of the online computation times are NLP function evaluations and KKT matrix factorization. While SCHEME 4<sub>30</sub> adds additional computation time compared to SCHEME 0 in the form of BLT analysis, it also results in a smaller DAE which may reduce computation in subsequent offline steps, allowing the additional computation time to be regained before we even reach the online computation.

Table 9 compares the offline and online computation times of SCHEME 0 and SCHEME 4<sub>30</sub> for the six benchmark problems. The online times are the averages observed in Section 5.3. We see that for all except one problem, SCHEME 4<sub>30</sub> has a lower total time. We also see that the offline times for SCHEME 4<sub>30</sub> become significantly larger for problems with many DAE variables. Since the implementation of SCHEME 4<sub>30</sub> is just a prototype with little regard for efficient offline computations (other than remaining in the realm of tractability), its implementation can probably be optimized to scale better. The bottleneck for Distillation column lies in the construction of the structural incidence matrix: Identifying the nonzero incidences and determining which of those are linear, which thus affects all schemes except 0.

### 5.6 Global collocation

We have so far only evaluated the schemes when combined with local collocation. Global collocation (also known as pseudospectral) methods, where a large number of collocation points and a small number of elements are used instead of vice versa, have become increasingly popular the last decade. We will thus also consider these methods, but only briefly for the sake of brevity and also avoiding the inadvertent comparison between the performances of the two flavours of collocation. One important difference between local and global collocation to keep in mind is that the NLPs resulting from global collocation tend to be smaller, due to their potential spectral

Table 9. Offline and online computation times [s].

Problem	SCHEME	Offline	Online	Total
Car	0	2.7	8.9	11.6
	4 <sub>30</sub>	2.7	6.8	9.5
CCPP	0	3.5	13.6	17.1
	4 <sub>30</sub>	3.6	0.9	4.5
Double pendulum	0	7.8	15.5	23.3
	4 <sub>30</sub>	7.1	0.8	7.9
Fourbar1	0	31.3	8.9	40.2
	4 <sub>30</sub>	43.4	1.5	44.9
HRSG	0	10.7	13.5	24.2
	4 <sub>30</sub>	12.5	8.9	21.4
Distillation column	0	14.7	$\infty$	$\infty$
	4 <sub>30</sub>	119.2	3.0	122.2

(exponential) convergence rate, but also more dense, due to their global temporal coupling.

We will limit our evaluation of global collocation to the Distillation column problem. Instead of using  $n_e = 20$  and  $n_c = 3$ , we use  $n_e = 1$  and  $n_c = 25$ . Just as we observed for local collocation, SCHEME 0 is unable to solve the problem. The results for the remaining schemes are shown in Table 10. All of the instances are valid. The results are similar to those obtained with local collocation, see Table 8. Two notable differences are that every SCHEME (except 0) solves every instance and that smaller values of  $\mu_{\text{tol}}$  are preferable. This indicates that it is crucial to preserve what little sparsity there is in an NLP resulting from global collocation. However, if the problem dimensions allow it, it may be more efficient to combine SCHEME 2 with dense rather than sparse numerical linear algebra, which is a possibility we do not consider further in this paper.

Table 10. Results for Distillation column with global collocation.

SCHEME	Success	Time	$\sigma_t$	Iter
1	100.0%	9.4	1.4	9.1
2	100.0%	8.9	1.1	8.6
3 <sub>40</sub>	100.0%	3.1	0.5	9.1
3 <sub>30</sub>	100.0%	3.5	0.6	9.1
3 <sub>20</sub>	100.0%	2.6	0.4	9.1
3 <sub>10</sub>	100.0%	2.4	0.4	9.1
3 <sub>5</sub>	100.0%	2.7	0.4	9.1
4 <sub>40</sub>	100.0%	6.0	0.8	8.6
4 <sub>30</sub>	100.0%	3.6	0.5	8.6
4 <sub>20</sub>	100.0%	2.6	0.4	8.6
4 <sub>10</sub>	100.0%	3.8	0.5	8.6
4 <sub>5</sub>	100.0%	2.0	0.2	8.6

## 6. Conclusion

DAE-constrained optimization problems are usually solved by exposing the full DAE to a discretization method. In this paper we considered ways of preprocessing the DAE by symbolically eliminating most of the algebraic variables using techniques based on block-triangular orderings, tearing, and sparsity preservation, resulting in 4 different schemes. These schemes have been implemented in the open-source JModelica.org platform, which solves dynamic optimization problems where the system dynamics are described using the Modelica language. We evaluated these schemes on 6 different optimal control problems when combined with direct local collocation, and found that the scheme that utilizes all of the proposed techniques performed

the best on average, often being an order of magnitude faster than the conventional scheme of exposing the full DAE to the discretization method.

We found that a suitable value of the density tolerance  $\mu_{\text{tol}}$  seems to be 15 when employing local collocation, although it is problem dependent and more data is needed to pinpoint a suitable default value with certainty. It also seems that smaller tolerances are preferable when there are more collocation points per element. We observed that the proposed techniques lend themselves well to typical Modelica models, due to the hierarchical modelling approach, but that they also show potential for flat DAEs.

Sparsity preservation was beneficial for the problems where it made a significant difference, excepting Fourbar1. A possible and tractable refinement of the proposed sparsity preservation procedure that may remedy this is to not only consider the sparsity of the NLP Jacobian, but also the Hessian of the NLP Lagrangian.

While block-triangular ordering and sparsity preservation guarantees preservation of numerical stability, tearing does not. Although numerical instability due to tearing does not appear to have been an issue for the considered problems, the authors are confident that there are industrially relevant problems where JModelica.org selects numerically troublesome tearing variables and residuals. This is a well-known drawback of tearing, which is difficult to address under typical circumstances. However, in the context of dynamic optimization a decent initial guess of the full solution is often required in order to solve the problem. Utilizing this initial guess, it should be tractable to design a numerical tearing algorithm which guarantees numerical stability (under the assumption that the initial guess is sufficiently close to the solution) using techniques such as those considered by Westerberg et al. [32, 58].

The used tearing algorithm first selects tearing variables and residuals to obtain causalized equations that are triangular and linear along the diagonal, and then selects additional tearing variables and residuals in order to not eliminate variables that are bounded, differential, or cause too much fill-in. Considering all of these criteria simultaneously, rather than sequentially, would enable fewer tearing variables and residuals to be selected.

## Acknowledgements

This work was supported by the Swedish Research Council through the LCCC Linnaeus Center. Fredrik Magnusson is a member of the eLLIIT Excellence Center at Lund University. The authors sincerely thank Ali Baharev for his valuable comments.

## References

- [1] M. Åberg, *Optimisation-friendly modelling of thermodynamic properties of media*, M.Sc. thesis, Department of Automatic Control, Lund University, Sweden, 2016.
- [2] J. Åkesson, *Optimica—An extension of Modelica supporting dynamic optimization*, in *Proceedings of the 6th International Modelica Conference*, Bielefeld, Germany, 2008, pp. 57–66.
- [3] J. Åkesson, K.E. Årzén, M. Gäfvert, T. Bergdahl, and H. Tummescheit, *Modeling and optimization with Optimica and JModelica.org—languages and tools for solving large-scale dynamic optimization problems*, *Comput. Chem. Eng.* 34 (2010), pp. 1737–1749.
- [4] J. Andersson, *A general-purpose software framework for dynamic optimization*, Ph.D. thesis, Arenberg Doctoral School, KU Leuven, 2013.
- [5] J. Andersson, J. Åkesson, F. Casella, and M. Diehl, *Integration of CasADi and JModelica.org*, in *Proceedings of the 8th International Modelica Conference*, Dresden, Germany, 2011, pp. 218–231.
- [6] M. Axelsson, F. Magnusson, and T. Henningsson, *A framework for nonlinear model predictive control in JModelica.org*, in *Proceedings of the 11th International Modelica Conference*, Paris, France, 2015, pp. 301–310.
- [7] B. Bachmann, L. Ochel, V. Ruge, M. Gebremedhin, P. Fritzson, V. Nezhadali, L. Eriksson, and M. Sivertsson,

- Parallel multiple-shooting and collocation optimization with OpenModelica*, in *Proceedings of the 9th International Modelica Conference*, Munich, Germany, 2012, pp. 659–668.
- [8] A. Baharev, H. Schichl, and A. Neumaier, *Tearing systems of nonlinear equations – I. A survey*. Submitted for publication. Available at [http://reliablecomputing.eu/baharev\\_tearing\\_survey.pdf](http://reliablecomputing.eu/baharev_tearing_survey.pdf).
  - [9] A. Baharev, H. Schichl, and A. Neumaier, *Tearing systems of nonlinear equations – II. A practical exact algorithm*. Submitted for publication. Available at [http://reliablecomputing.eu/baharev\\_tearing\\_exact\\_algorithm.pdf](http://reliablecomputing.eu/baharev_tearing_exact_algorithm.pdf).
  - [10] M. Baltes, R. Schneider, C. Sturm, and M. Reuss, *Optimal experimental design for parameter estimation in unstructured growth models*, *Biotech. Progr.* 10 (1994), pp. 480–488.
  - [11] K. Berntorp and F. Magnusson, *Hierarchical predictive control for ground-vehicle maneuvering*, in *American Control Conference, 2015*, Chicago, IL, 2015, pp. 2771–2776.
  - [12] K. Berntorp, B. Olofsson, K. Lundahl, and L. Nielsen, *Models and methodology for optimal trajectory generation in safety-critical road-vehicle manoeuvres*, *Vehicle System Dynamics* 52 (2014), pp. 1304–1332.
  - [13] J.T. Betts, *Survey of numerical methods for trajectory optimization*, *J. Guid. Contr. Dynam.* 21 (1998), pp. 193–207.
  - [14] J.T. Betts, *Practical Methods for Optimal Control and Estimation Using Nonlinear Programming*, 2nd ed., SIAM, Philadelphia, PA, 2010.
  - [15] L.T. Biegler, *Nonlinear Programming: Concepts, Algorithms, and Applications to Chemical Processes*, MOS-SIAM, Philadelphia, PA, 2010.
  - [16] T. Blochwitz, M. Otter, M. Arnold, C. Bausch, C. Clauß, H. Elmqvist, A. Junghanns, J. Mauss, M. Monteiro, T. Neidhold, D. Neumerkel, H. Olsson, J.V. Peetz, and S. Wolf, *The functional mockup interface for tool independent exchange of simulation models*, in *Proceedings of the 8th International Modelica Conference*, Dresden, Germany, 2011, pp. 105–114.
  - [17] K. Brenan, S. Campbell, and L. Petzold, *Numerical Solution of Initial-Value Problems in Differential-Algebraic Equations*, *Classics in Applied Mathematics*, Vol. 14, SIAM, Philadelphia, PA, 1996.
  - [18] C. Büskens and D. Wassel, *The ESA NLP solver WORHP*, in *Modeling and Optimization in Space Engineering*, G. Fasano and J.D. Pintér, eds., Vol. 73, Springer, New York, NY, 2013, pp. 85–110.
  - [19] F. Casella, F. Donida, and J. Åkesson, *Object-oriented modeling and optimal control: A case study in power plant start-up*, in *18th IFAC World Congress*, Milano, Italy, 2011, pp. 9549–9554.
  - [20] F.E. Cellier and E. Kofman, *Continuous System Simulation*, Springer, New York, NY, 2006.
  - [21] A.M. Cervantes, A. Wächter, R.H. Tütüncü, and L.T. Biegler, *A reduced space interior point strategy for optimization of differential algebraic systems*, *Comput. Chem. Eng.* 24 (2000), pp. 39–51.
  - [22] M. Diehl, *Real-time optimization for large scale nonlinear processes*, Ph.D. thesis, Heidelberg University, 2002.
  - [23] M. Diehl, H.G. Bock, J.P. Schlöder, R. Findeisen, Z. Nagy, and F. Allgöwer, *Real-time optimization and nonlinear model predictive control of processes governed by differential-algebraic equations*, *J. Proc. Cont.* 12 (2002), pp. 577–585.
  - [24] E.D. Dolan and J.J. Moré, *Benchmarking optimization software with performance profiles*, *Math. Program.* 91 (2002), pp. 201–213.
  - [25] I.S. Duff, A. Erisman, and J.K. Reid, *Direct Methods for Sparse Matrices*, Clarendon Press, Oxford, United Kingdom, 1986.
  - [26] I.S. Duff and J.K. Reid, *The design of MA48: a code for the direct solution of sparse unsymmetric linear systems of equations*, *ACM Trans. Math. Software* 22 (1996), pp. 187–226.
  - [27] H. Elmqvist and M. Otter, *Methods for tearing systems of equations in object-oriented modeling*, in *European Simulation Multiconference*, Barcelona, Spain, 1994, pp. 326–332.
  - [28] R. Fletcher, *Block triangular orderings and factors for sparse matrices in LP*, in *Numerical Analysis 1997*, D. Griffiths, G. Watson, and D. Higham, eds., Longman, Harlow, United Kingdom, 1998, pp. 91–110.
  - [29] R. Franke, *Formulation of dynamic optimization problems using Modelica and their efficient solution*, in *Proceedings of the 2nd International Modelica Conference*, Oberpfaffenhofen, Germany, 2002, pp. 315–322.
  - [30] P. Fritzson, *Principles of Object-Oriented Modeling and Simulation with Modelica 3.3*, 2nd ed., Wiley-IEEE Press, Piscataway, NJ, 2015.
  - [31] A. George, *Nested dissection of a regular finite element mesh*, *SIAM J. Numer. Anal.* 10 (1973), pp. 345–363.
  - [32] P.K. Gupta, A.W. Westerberg, J.E. Hendry, and R.R. Hughes, *Assigning output variables to equations using linear programming*, *AIChE J.* 20 (1974), pp. 397–399.
  - [33] E. Hairer and G. Wanner, *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*, 2nd ed., Springer-Verlag, Berlin, Germany, 1996.
  - [34] J.D. Hedengren, *A nonlinear model library for dynamics and control*, 2008. Available online: [http://www.hedengren.net/research/Publications/Cache\\_2008/NonlinearModellibrary.pdf](http://www.hedengren.net/research/Publications/Cache_2008/NonlinearModellibrary.pdf).
  - [35] J.E. Hopcroft and R.M. Karp, *An  $n^{5/2}$  algorithm for maximum matchings in bipartite graphs*, *SIAM J. Comput.* 2 (1973), pp. 225–231.

- [36] HSL, *A collection of Fortran codes for large scale scientific computation*, software available at <http://www.hsl.rl.ac.uk>.
- [37] G. Karypis and V. Kumar, *A fast and high quality multilevel scheme for partitioning irregular graphs*, SIAM J. Sci. Comput. 20 (1998), pp. 359–392.
- [38] G. Kron, *Diakoptics: The Piecewise Solution of Large-Scale Systems*, Vol. 2, MacDonald, London, United Kingdom, 1963.
- [39] K. Krüger, R. Franke, and M. Rode, *Optimization of boiler start-up using a nonlinear boiler model and hard constraints*, Energy 29 (2004), pp. 2239–2251.
- [40] E. Lazutkin, A. Geletu, S. Hopfgarten, and P. Li, *An analytical hessian and parallel-computing approach for efficient dynamic optimization based on control-variable correlation analysis*, Ind. Eng. Chem. Res. 54 (2015), pp. 12086–12095.
- [41] B. Lennernäs, *A CasADi based toolchain for JModelica.org*, M.Sc. thesis, Department of Automatic Control, Lund University, Sweden, 2013.
- [42] J.M. Maciejowski, *Predictive Control with Constraints*, Prentice-Hall, Upper Saddle River, NJ, 2002.
- [43] F. Magnusson, *Numerical and symbolic methods for dynamic optimization*, Ph.D. thesis, Lund University, 2016.
- [44] F. Magnusson, *Symbolic elimination performance profiles* (2016), Available at [https://github.com/Bleevoe/blt\\_performance\\_profiles](https://github.com/Bleevoe/blt_performance_profiles), GitHub repository.
- [45] F. Magnusson and J. Åkesson, *Dynamic optimization in JModelica.org*, Processes 3 (2015), pp. 471–496.
- [46] F. Magnusson, K. Berntorp, B. Olofsson, and J. Åkesson, *Symbolic transformations of dynamic optimization problems*, in *Proceedings of the 10th International Modelica Conference*, Lund, Sweden, 2014, pp. 1027–1036.
- [47] H.M. Markowitz, *The elimination form of the inverse and its application to linear programming*, Manag. Sci. 3 (1957), pp. 255–269.
- [48] S.E. Mattsson, H. Olsson, and H. Elmqvist, *Dynamic selection of states in Dymola*, in *Modelica Workshop 2000 Proceedings*, Lund, Sweden, 2000, pp. 61–67.
- [49] S.E. Mattsson and G. Söderlind, *Index reduction in differential-algebraic equations using dummy derivatives*, SIAM J. Sci. Comput. 14 (1993), pp. 677–692.
- [50] P. Meijer, *Tearing differential algebraic equations*, M.Sc. thesis, Centre for Mathematical Sciences, Lund University, Sweden, 2011.
- [51] P. Parini, *Object oriented modeling and dynamic optimization of energy systems with application to combined cycle power plant start-up*, M.Sc. thesis, Dipartimento di Elettronica, Informazione e Bioingegneria, Politecnico di Milano, Italy, 2015.
- [52] A. Prata, J. Oldenburg, A. Kroll, and W. Marquardt, *Integrated scheduling and dynamic optimization of grade transitions for a continuous polymerization reactor*, Comput. Chem. Eng. 32 (2008), pp. 463–476.
- [53] A.V. Rao, *A survey of numerical methods for optimal control*, in *Proceedings of the 2009 AAS/AIAA Astrodynamics Specialists Conference*, Pittsburgh, PA, 2009.
- [54] H. Runvik, *Modelling and start-up optimization of a coal-fired power plant*, M.Sc. thesis, Department of Automatic Control, Lund University, Sweden, 2014.
- [55] S.M. Safdarnejad, J.D. Hedengren, N.R. Lewis, and E.L. Haseltine, *Initialization strategies for optimization of dynamic systems*, Comput. Chem. Eng. 78 (2015), pp. 39–50.
- [56] R.E. Tarjan, *Depth-first search and linear graph algorithms*, SIAM J. Comput. 1 (1972), pp. 146–160.
- [57] M. Thelander Andrén and C. Wedding, *Development of a solution for start-up optimization of a thermal power plant*, M.Sc. thesis, Department of Automatic Control, Lund University, Sweden, 2015.
- [58] A.W. Westerberg and F.C. Edie, *Computer-aided design, part 1 enhancing convergence properties by the choice of output variable assignments in the solution of sparse equation sets*, Chem. Eng. J. 2 (1971), pp. 9–16.
- [59] D.P. Word, J. Kang, J. Åkesson, and C.D. Laird, *Efficient parallel solution of large-scale nonlinear dynamic optimization problems*, Comput. Optim. and Appl. 59 (2014), pp. 667–688.
- [60] A. Wächter and L.T. Biegler, *On the implementation of a primal-dual interior point filter line search algorithm for large-scale nonlinear programming*, Math. Program. 106 (2006), pp. 25–57.
- [61] V.M. Zavala, *Inference of building occupancy signals using moving horizon estimation and Fourier regularization*, J. Proc. Cont. 24 (2014), pp. 714–722.