

Indirect methods for optimal control of hybrid PDE-dynamical / switching systems using relaxation

Falk M. Hante ^a and Christian Kuchler ^a

^a Humboldt-Universität zu Berlin, Unter den Linden 6, 10099 Berlin, Germany

ARTICLE HISTORY

Compiled December 11, 2023

ABSTRACT

We propose a novel algorithmic approach to computationally solve optimal control problems governed by linear evolution-type PDEs including a state-dependent control-regime switching mechanism. We introduce an equivalent mixed-integer formulation featuring vanishing constraints arising by methods of disjunctive programming. We embed the problem into the class of equilibrium constraints by introduction of an additional slack variable. Based on theoretical results associated with Sum-Up-Rounding strategies, we proceed with the solution of the related relaxed formulation by an indirect approach. In order to obtain a computationally tractable optimality system, we apply a Moreau-Yosida type penalty approach of the vanishing constraints. After the theoretical discussion, we introduce and exert the algorithmic framework founded on a semismooth Newton method. Finally, we communicate computational experiments based on our approach.

AMS CLASSIFICATION

49K20 49M05

KEYWORDS

Hybrid systems, Implicit switching, Vanishing constraints, Equilibrium constraints, Moreau-Yosida penalization, Semismooth Newton

1. Introduction

Various technical and economical processes can be modeled as optimal control problems. Applications displaying implicit switching behaviour for instance cover safety circuits for heating processes, where the heating process is interrupted once a certain temperature treshold is reached, c.f. [29]. Another example is presented by bacteria growth within a petri dish, c.f. [30], where the transition of the bacteria from an active to dormant state or vice versa, is regulated by the overall cell concentration within the dish. Finally, we mention applications in gas networks where the transport through the network is optimised, while flap-valves open/close state-dependently to prevent flow reversal for example on compressor in- and outlets, c.f. [14].

Apparently the paramount assignment in the implicit switching framework is a suitable resolution of the involved switching rule since it introduces a non-linear coupling between the involved control and obtained state into the system. For ordinary

Email: falk.hante@hu-berlin.de

Email: christian.kuchler@hu-berlin.de

differential equations (ODEs) one can for instance proceed along the path of disjunctive programming proposed by Bock et al., c.f. [4], and replace the implicit switching rule by a combination of explicit switching variables and vanishing constraints. After discretisation the authors obtain a mixed-integer optimal control problem (MIOCP), that can be solved by relaxed partial outer convexification and an adapted rounding scheme, c.f. [25].

Our goal in this paper is to extend and apply this approach also in the presence of dynamics governed by an abstract semigroup setting. This framework covers important applications such as hyperbolic or parabolic partial differential equations (PDEs). In contrast to the available techniques we proceed with an indirect approach. However, the reformulation based on disjunctive programming does not immediately enable the characterisation of minimisers via necessary first order optimality conditions. The core challenges can be traced back to the appearance of binary multipliers in the disjunctive formulation in combination with vanishing constraints (VCs), which result from the resolution of the implicit switching rule. Therefore we proceed similar to [13] and derive a surrogate model that enables us to derive necessary optimality conditions. In [13] we applied a time transformation method, c.f. [11] and [22], to resolve the binary multipliers on a reference interval, where the choice of the discrete control is fixed on given subdomains. In the present paper we pursue a different approach by a direct relaxation of the appearing binary multipliers. Both approaches then still contain the aforementioned vanishing constraints. We treat those constraints equally by embedding them into framework of equilibrium constraints (ECs) by introduction of a slack variable and application of a path following approach afterwards.

The core contribution of our paper is an indirect numerical approach towards a relaxed, but vanishing-constrained surrogate optimal control problem. Along this path we utilise a technical result based on rounding schemes designed to construct binary controls, c.f. [24] and [25], which allows us to consider a surrogate relaxed formulation that still guarantees certain approximation properties towards solutions of the original disjunctive formulation. Secondly, we combine this theoretical result with an embedding of vanishing constraints, [1] and [17], into the framework of ECs, [9], [10] and [23], by introduction of an additional slack variable together with a Moreau-Yosida type penalty approach, c.f. [18], towards the resulting mixed control-state constraints. The obtained setting then permits the characterisation of candidates for optimality of the surrogate model by a set of first order criteria as well as efficient numerical treatment.

This article is organised as follows. In Section 2 the setting in combination with an example is introduced. Section 3 starts with the reformulation of the presented problem as a MIOCP with VCs by means of disjunctive programming. On the basis of available approximation results for the involved rounding strategies a relaxation approach is pursued. In the process the appearing VCs are embedded into the framework of ECs. Afterwards the section is concluded by the introduction of the penalised formulation and the derivation of the corresponding necessary optimality conditions. In Section 4 we present our algorithm based on a semismooth Newton scheme together with a globalization heuristic to numerically solve the set of necessary optimality conditions derived at the end of Section 3. In Section 5 we report on the performance of our algorithm in numerical experiments. In Section 6 the core innovations are summarised. Also future research branches are mentioned.

2. Problem formulation

For $T > 0$ we consider the following problem involving the generator A of a C_0 -semigroup on a real space X . We denote the space of the involved control by U . The problem under consideration is

$$\min_{y,u} J(y, u) = \frac{1}{2} \|y - y_{des}\|_{L^2([0,T];X)}^2 + \frac{\gamma_u}{2} \|u\|_{L^2([0,T];U)}^2 \quad (1a)$$

s.t.

$$\dot{y}(t) = Ay(t) + f_{d(t)}(u(t)) \quad t \in [0, T], \quad (1b)$$

$$y(0) = y_0, \quad (1c)$$

$$d(t) = C(y(t)) \quad t \in [0, T], \quad (1d)$$

$$C = R \circ S. \quad (1e)$$

In this subsequent discussion, we impose the following assumptions.

Assumption 1.

- (1) The spaces X and U are Hilbert spaces.
- (2) The mapping $f_d : U \rightarrow X, u \mapsto f_d(u)$ is linear and continuous, i.e., $f_d \in \mathcal{L}(U, X)$ for all $d \in [D]$.

In the provided setting $C : X \rightarrow [D] := \{1, \dots, D\}$ denotes the mode function for a fixed number of modes $D \in \mathbb{N}$. The input for C is the current evaluation of the state variable $y(t)$ for any $t \in [0, T]$. The resulting mode function $d(t)$ in turn determines the inhomogeneity $f_{d(t)}$ in (1b). Therefore the choice of the current control mode $d(t)$ is a function of the current state evaluation $y(t)$, i.e., the current mode implicitly depends on the evaluation of the state. Hence the choice of the mode $d(t)$ is unavailable as an independent optimisation variable without further technical assumptions on C . In particular the switches of $d(t)$ cannot be formulated explicitly.

Equation (1e) is a technical assumption to provide additional structure towards the involved switching surfaces, c.f. [3] and [4]. We further assume that $S : X \rightarrow \mathbb{R}$ is a linear and continuous function and assert $R : \mathbb{R} \rightarrow [D]$ is a piecewise constant function. We also postulate that the inverse image of each mode $d \in [D]$ is a half closed interval, i.e., there exist real numbers $a_d < b_d$ such that

$$R^{-1}(d) =]a_d, b_d], \quad \forall d \in [D].$$

Here we apply the convention that $b_d := \infty$ if $R^{-1}(d)$ is not bounded from above respectively $a_d := -\infty$ if $R^{-1}(d)$ is not bounded from below.

In addition we assume that the desired state y_{des} satisfies $y_{des} \in L^2([0, T]; X)$. We denote by $y \in C^0([0, T]; X)$ a mild solution to (1b) – (1e) in the following sense, c.f. [26].

Let $I = (t_i, t_f) \subset [0, T]$ be an open interval such that $C|_I \equiv d \in [D]$ and denote by $(T(t))_{t \geq 0}$ the one-parameter semigroup generated by A . Then for any $t \in \bar{I}$ the evaluation of the solution $y(t)$ to (1b) – (1e) on I is defined by the following variation of constants formula

$$y(t) = T(t - t_i)y(t_i) + \int_{t_i}^t T(t - s)f_d(s, u(s)) ds. \quad (2)$$

We introduce the following abbreviations $C^k([0, T]; X) := C^k_{[0, T]}(X)$ for $k \in \mathbb{N}$ and $L^p([0, T]; X) := L^p_{[0, T]}(X)$ with $p \in [1, \infty]$.

In order to show that the setting concerns non-trivial dynamics, we consider the following example. We will revisit it also in the numerical discussion.

Example 2.1. Let $D = 2$ and $\Omega_d \subset \Omega$ for $d \in [2]$ be bounded domains such that $\Omega_1 \cap \Omega_2 = \emptyset$. Let χ_M denote the characteristic function for a Lebesgue measurable set $M \subset \mathbb{R}^n$, i.e.,

$$\chi_M : \mathbb{R}^n \rightarrow \mathbb{R}, x \mapsto \begin{cases} 1, & \text{if } x \in M, \\ 0, & \text{if else.} \end{cases}$$

We consider the optimisation problem

$$\begin{aligned} \min_{y, u} J(y, u) &= \frac{1}{2} \|y - y_{des}\|_{L^2_{[0, T]}(X)}^2 + \frac{\gamma_u}{2} \|u\|_{L^2_{[0, T]}(U)}^2 \\ \text{s.t.} & \\ \dot{y}(t) &= \Delta y(t) + u(t) \chi_{\Omega_{d(t)}} & t \in [0, T], \\ y(t)|_{\partial\Omega} &= 0 & t \in [0, T], \\ y(0) &= y_0, \\ d(t) &= \begin{cases} 1, & \text{if } \int_{\Omega} y(x, t) dx \leq \sigma, \\ 2, & \text{if } \int_{\Omega} y(x, t) dx > \sigma, \end{cases} & t \in [0, T]. \end{aligned}$$

In this example we attempt a best approximation towards y_{des} by solutions of the heat equation denoted by y involving the Laplace operator with respect to $x \in \Omega$ denoted by Δy . The right side of the equation is the control u restricted to disjoint control domains Ω_d for $d \in [2]$. The switching mechanism is realised by the evaluation of the space integral $\int_{\Omega} y(x, t) dx$ at any time $t \in [0, T]$ against a selected threshold σ .

In order to bring this into the form (1) we set $Ay = \Delta y$ for $y \in D(A) = H_0^1(\Omega) \cap H^2(\Omega)$. We consider the spaces $U = L^2(\Omega)$ and $X = L^2(\Omega)$ and define the mappings

$$\begin{aligned} f_d &: L^2([0, T]; L^2(\Omega)) \rightarrow L^2([0, T]; L^2(\Omega)), u \mapsto u \cdot \chi_{\Omega_d}, \\ C &: L^2(\Omega) \rightarrow [2], y \mapsto \begin{cases} 1, & \text{if } \int_{\Omega} y(x) dx \leq \sigma, \\ 2, & \text{if } \int_{\Omega} y(x) dx > \sigma, \end{cases} \\ S &: L^2(\Omega) \rightarrow \mathbb{R}, y \mapsto \int_{\Omega} y(x) dx, \\ R &: \mathbb{R} \rightarrow [2], y \mapsto \begin{cases} 1, & \text{if } y \leq \sigma, \\ 2, & \text{if } y > \sigma. \end{cases} \end{aligned}$$

Furthermore the assertions posed under Assumption 1 hold and the linear operator A creates a C_0 -semigroup $(T(t))_{t \geq 0}$, c.f., [2, Chapter 2.10.1].

3. Reformulation

In this section we establish reformulations of (1) to a computationally more tractable setting. In order to simplify the considerations and similarly as in [4], we include certain technical assumptions on S .

Assumption 2. *Let $t_s \in \mathbb{R}$ be such that*

$$S(y(t_s)) \in \bigcup_{d \in [D]} \{a_d, b_d\}. \quad (3)$$

We set

$$y^-(t_s) := \lim_{t \downarrow t_s} y(t) \in X, \quad y^+(t_s) := \lim_{t \uparrow t_s} y(t) \in X,$$

for the limit of the mild solution of the state equation (1b) - (1d) from the left and from the right at t_s . Furthermore the derivatives from the left and from the right

$$S'_-(t_s) := \lim_{t \uparrow t_s} \frac{\partial S}{\partial t}(y^-(t_s)), \quad S'_+(t_s) := \lim_{t \downarrow t_s} \frac{\partial S}{\partial t}(y^+(t_s)),$$

are supposed to exist. The problem (1) satisfies the transversality assumption if the evaluation of $S'_+(t_s) \cdot S'_-(t_s) > 0$ for all $t_s \in \mathbb{R}$, that fulfil (3), holds.

The assertions posed in Assumption 2 prohibit the state trajectories to slide tangential with respect to the switching surfaces. This simplifies the discussion of the switching behavior as sliding mode solutions in the sense of Filippov [8] are ruled out. The observed switching behavior is then called consistent.

Aside from the behavior on the switching manifold, another important aspect of switched processes is the number of switching points, in particular the avoidance of arcs, which display accumulation of switches, i.e. so-called Zeno-behavior.

Assumption 3. *We assume that system (1) possesses only a finite number of switches for each admissible control and state pair (u, y) .*

Remark 1. The Assumptions 2 and 3 guarantee that the mild solution of the form (2) is well-posed in the classical sense. For conditions imposing Assumption 3, e.g., for certain hyperbolic PDEs, see [12] and [27]. However, solution of hybrid dynamical systems can be defined also in a broader sense, i.e., in a set-valued sense with possible branching of solutions at switching surfaces when Assumption 2 fails [31]. This is typically the case even in the most simplest practical examples [29]. The subsequent reformulations then remain applicable, but require some care with respect to their interpretation and adaptations. Continuous dependency is then to be replaced with upper-semicontinuity of the solution set [31]. If desired also sliding on switching manifolds can be included by adding an extra sliding mode into the admissible right-hand sides. The subsequent reformulations then represent one of all possible solution branches. This is in particular to be remembered whenever solutions are approximated numerically and if one passes back from the reformulation to the original problem formulation. We demonstrate this for our numerical realization in Section 5 on a particular example.

We apply the ideas provided in [4] and originating from disjunctive programming to

reformulate (1) by introducing binary multipliers and vanishing constraints. For $\delta \geq 0$ we consider the following problem

$$\min_{y,u,\omega} J(y, u) = \frac{1}{2} \|y - y_{des}\|_{L^2_{[0,T]}(X)}^2 + \frac{\gamma u}{2} \|u\|_{L^2_{[0,T]}(U)}^2 \quad (4a)$$

s.t.

$$\dot{y}(t) = Ay(t) + \sum_{d=1}^D \omega_d(t) f_d(u(t)) \quad \text{a.e. } t \in [0, T], \quad (4b)$$

$$y(0) = y_0, \quad (4c)$$

$$\omega_d(t) \in \{0, 1\} \quad d \in [D], \text{ a.e. } t \in [0, T], \quad (4d)$$

$$1 = \sum_{d=1}^D \omega_d(t) \quad \text{a.e. } t \in [0, T], \quad (4e)$$

$$-\delta \leq \omega_d(t)(S(y(t)) - a_d) \quad d \in [D], \text{ a.e. } t \in [0, T], \quad (4f)$$

$$-\delta \leq \omega_d(t)(b_d - S(y(t))) \quad d \in [D], \text{ a.e. } t \in [0, T]. \quad (4g)$$

Under validity of Assumption 2 the formulations (1) and (4) with choice $\delta = 0$ are in deed equivalent. The necessity to relax the constraints in (4f) – (4g) by $\delta > 0$ in order to enable an approximation of the optimal cost for the relaxed problem by solutions of (4) is demonstrated by an example first mentioned by Cesari, c.f. [5], and picked up in [20, Section 4]. The advantage of formulation (4) is the removal of implicit switching mechanism, but this comes at the expense of introducing binary variables and vanishing constraints. In a next step we relax the binary multipliers $\omega_d \in \{0, 1\}$ to $\alpha_d \in [0, 1]$. We obtain the following relaxed formulation of (4)

$$\min_{y,u,\alpha} J(y, u) = \frac{1}{2} \|y - y_{des}\|_{L^2_{[0,T]}(X)}^2 + \frac{\gamma u}{2} \|u\|_{L^2_{[0,T]}(U)}^2 \quad (5a)$$

s.t.

$$\dot{y}(t) = Ay(t) + \sum_{d=1}^D \alpha_d(t) f_d(u(t)) \quad \text{a.e. } t \in [0, T], \quad (5b)$$

$$y(0) = y_0, \quad (5c)$$

$$\alpha_d(t) \in [0, 1] \quad d \in [D], \text{ a.e. } t \in [0, T], \quad (5d)$$

$$1 = \sum_{d=1}^D \alpha_d(t) \quad \text{a.e. } t \in [0, T], \quad (5e)$$

$$0 \leq \alpha_d(t)(S(y(t)) - a_d) \quad d \in [D], \text{ a.e. } t \in [0, T], \quad (5f)$$

$$0 \leq \alpha_d(t)(b_d - S(y(t))) \quad d \in [D], \text{ a.e. } t \in [0, T]. \quad (5g)$$

We require the notion for a sequence to possess vanishing integrality gap, c.f. [24].

Definition 3.1. Let $(\varphi_k)_{k \in \mathbb{N}} \subset L^\infty_{[0,T]}(\mathbb{R})$ be a bounded sequence such that $\varphi_k(t) := \int_0^t \varphi_k(t) dt$ satisfies

$$\|\varphi_k\|_\infty \rightarrow 0, \quad k \rightarrow \infty.$$

Then we call $(\varphi_k)_{k \in \mathbb{N}}$ a sequence of vanishing integrality gap.

Definition 3.1 is closely linked with certain rounding strategies.

Definition 3.2. Let $0 = t_0 < \dots < t_N = T$ be a rounding grid of $[0, T]$ with maximum discretisation width $\Delta_N := \max_{i \in [N]} t_i - t_{i-1}$. We abbreviate $G_N := \cup_{i \in \{0, \dots, N\}} t_i$.

For a function $\alpha \in L^\infty_{[0, T]}(\mathbb{R}^D)$, we define a binary-valued piecewise constant function $\omega : [0, T] \rightarrow \{0, 1\}^D$ iteratively for $i \in [N]$ as

$$\begin{aligned} \varphi_0 &:= 0_{\mathbb{R}^D}, \\ \gamma_i &:= \varphi_{i-1} + \int_{t_{i-1}}^{t_i} \alpha(t) dt, \\ \omega_{i,j} &:= \begin{cases} 1 & j = \arg \max\{\gamma_{i,k} \mid k \in F_i\}, \\ 0 & \text{else.} \end{cases} & \forall j \in [D], \\ \omega_j|_{(t_i, t_{i+1})} &:= \omega_{i,j} & \forall j \in [D], \\ \varphi_i &:= \int_{t_0}^{t_i} \alpha(t) - \omega(t) dt. \end{aligned}$$

In case there exist several indices $j \in [D]$ such that $\gamma_{i,j} = \max\{\gamma_{i,k} \mid k \in F_i\}$, the tie is to be broken for instance by the choice of the minimum among those indices.

The rounding strategies are equipped with different labels depending on the particular choice of indices admissible to rounding, F_i , on each interval.

Definition 3.3. We consider in particular the following rounding strategies associated with a rounding grid G_N , where the admissible index sets

$$F_i := [D], \quad (\text{SUR-SOS})$$

$$F_i := \{j \in [D] \mid \int_{t_i}^{t_{i+1}} \alpha_j(t) dt > 0\}, \quad (\text{SUR-SOS-VC})$$

are selected $\forall i \in \{0, \dots, N-1\}$.

Lemma 3.4. Let $\alpha \in L^\infty_{[0, T]}(\mathbb{R}^D)$ with $\alpha(t) \in [0, 1]^D$ and $\sum_{d=1}^D \alpha_d(t) = 1$ f.a.e. $t \in [0, T]$ and a rounding grid G_N be given. Then the rounding strategies in Definition 3.3 produce $\omega \in L^\infty_{[0, T]}(\mathbb{R}^D)$ with $\omega(t) \in \{0, 1\}^D$ and $\sum_{d=1}^D \omega_d(t) = 1$ f.a.e. $t \in [0, T]$ such that there exists a constant $C > 0$ satisfying

$$\sup_{t \in [0, T]} \left\| \int_0^t \alpha(t) - \omega(t) dt \right\|_\infty \leq C \Delta_N. \quad (6)$$

Proof. The proof of the statement for (SUR-SOS) can be found in [28] and with respect to (SUR-SOS-VC) we refer to [25]. \square

The consequence of Lemma 3.4 can be summarised as follows. Let $(G_k)_{k \in \mathbb{N}}$ be sequence of rounding grids such that $\lim_{k \rightarrow \infty} \Delta_k = 0$. Furthermore suppose that the family $(\alpha_k)_{k \in \mathbb{N}}$ satisfies the assumptions posed in Lemma 3.4 and $(\omega_k)_{k \in \mathbb{N}}$ is generated from $(\alpha_k)_{k \in \mathbb{N}}$ by application of an arbitrary rounding strategy from Definition 3.2.

Then according to (6) the components of $(\varphi_k)_{k \in \mathbb{N}}$ with $\varphi_{k,d} := \alpha_{k,d} - \omega_{k,d}$ are of vanishing integrality gap for each $d \in [D]$, c.f. Definition 3.1.

Remark 2. As stated in Lemma 3.4 both involved presented strategies (SUR-SOS) and (SUR-SOS-VC) satisfy estimate (6) for the obtained integrality gap. However, the involved constant C displays different asymptotic behavior $\mathcal{O}(D)$ for $D \rightarrow \infty$. For the scheme (SUR-SOS) the behavior $\mathcal{O}(D)$ was for instance shown in [28] and later even improved to $\mathcal{O}(\log(D))$ in [20]. In [20] the asymptotic behavior $\mathcal{O}(D)$ with respect to (SUR-SOS-VC) has been proved. However, the asymptotic property $\mathcal{O}(\log(D))$ was ruled out.

Next we aim to discuss the approximation properties of solutions to (5) with respect to (4). We formulate a first lemma regarding the feasibility of the obtained binary control and approximation properties of the associated state in absence of the vanishing constraints.

Lemma 3.5. [24, Proposition 2.3]. *Let a trajectory $y \in C_{[0,T]}^0(X)$ together with controls $u \in L_{[0,T]}^2(U)$ and $\alpha \in L_{[0,T]}^\infty(\mathbb{R}^D)$ be feasible for (5). Let $(\omega_k)_{k \in \mathbb{N}} \in L_{[0,T]}^\infty(\mathbb{R}^D)$ be binary valued functions, such that $\varphi_{k,d} := \alpha_d - \omega_{k,d}$ are of vanishing integrality gap for all $d \in [D]$. Then for every $\delta > 0$ there exists a mild solution $y^\delta \in C_{[0,T]}^0(X)$ constructed by $u \in L_{[0,T]}^2(U)$ and $\omega^\delta \in L_{[0,T]}^\infty(\mathbb{R}^D)$ such that for $(y^\delta, u, \omega^\delta)$ the estimate*

$$\|y(u, \alpha) - y^\delta(u, \omega^\delta)\|_{C_{[0,T]}^0(X)} < \delta, \quad (7)$$

holds.

The previous Lemma 3.5 implies that $(y^\delta, u, \omega^\delta)$ satisfies (4b) - (4e). We proceed with a statement on the approximation property in the presence of vanishing constraints.

Lemma 3.6. [25, Theorem 2.1 (3)]. *Let a trajectory $y \in C_{[0,T]}^0(X)$ together with controls $u \in L_{[0,T]}^2(U)$ and $\alpha \in L_{[0,T]}^\infty(\mathbb{R}^D)$ be feasible for (5). Let $(\omega_k)_{k \in \mathbb{N}} \in L_{[0,T]}^\infty(\mathbb{R}^D)$ be binary-valued functions such that $\varphi_{k,d} = \alpha_d - \omega_{k,d}$ are of vanishing integrality gap for all $d \in [D]$. Then for every $\delta > 0$ there exists a mild solution $y^\delta \in C_{[0,T]}^0(X)$ constructed by $u \in L_{[0,T]}^2(\mathbb{R}^D)$ and $\omega^\delta \in L_{[0,T]}^\infty(\mathbb{R}^D)$ such that for $(y^\delta, u^\delta, \omega^\delta)$ the constraints (4f) - (4g) are satisfied.*

We can now combine the last two lemmata to obtain the following result.

Theorem 3.7. *Let a trajectory $y \in C_{[0,T]}^0(X)$ together with controls $u \in L_{[0,T]}^2(U)$ and $\alpha \in L_{[0,T]}^\infty(\mathbb{R}^D)$ be feasible for (5). Let $\delta > 0$ and $\varepsilon > 0$ be chosen arbitrarily. Then for a family of rounding grids with gridwidth $\Delta_k \rightarrow 0$, the binary control ω_k obtained from the relaxed control α by the application of the rounding method (SUR-SOS-VC) satisfies the following properties: There exists an $k_0 \in \mathbb{N}$ such that for all $k \geq k_0$ it holds that*

- i) The triple $(y(u, \omega), u, \omega)$ is feasible for (4).*
- ii) The cost function satisfies $|J(y(u, \alpha), u) - J(y(u, \omega_k), u)| < \varepsilon$.*

Proof. The proof for i) is given by the combination of Lemma 3.5 and Lemma 3.6. For ii) we argue as follows. The continuity of J at $(y(u, \alpha), u)$ ensures for every $\varepsilon > 0$ the

existence of $\delta_J > 0$ such that $|J(y(u, \alpha), u) - J(\tilde{y}, u)| < \varepsilon$ holds for all $\tilde{y} \in C_{[0, T]}^0(X)$, satisfying $\|y(u, \alpha) - \tilde{y}\|_{C_{[0, T]}^0(X)} < \delta_J$. The claim follows with the results of Lemma 3.5 and Lemma 3.6 in combination with the choice $\hat{\delta} := \min\{\delta_J, \delta\}$. \square

Theorem 3.7 provides a theoretical justification to solve the relaxed formulation (5) instead of (4), as from a solution $(y(u, \alpha), u, \alpha)$ to (5), we can construct by application of (SUR-SOS-VC) for any threshold $\delta > 0$ a binary feasible solution $(y(u, \omega_k), u, \omega_k)$ to (4) with arbitrary $\varepsilon > 0$ deviation of the corresponding costs $J(y(u, \alpha), u)$ towards $J(y(u, \omega_k), u)$. Therefore we continue with the discussion of solutions for the relaxed formulation. We proceed along the path of indirect methods and hence require suitable optimality conditions for (5). The conditions (5f) and (5g) in combination (5d) present vanishing constraints in a Banach space. These mathematical programs together with equilibrium constraints are currently still subject to research, c.f. [21] and [32], and a unified approach is still unavailable. Therefore we advance with a penalty approach as in [13].

At first we transform the vanishing constraints into equilibrium constraints by the addition of a slack variable $s(t) := (s_a(t), s_b(t))^T \in \mathbb{R}^{2D}$ and penalise the resulting constraint afterwards by $\gamma_{EC} > 0$. We set

$$J_{EC}(\alpha, s) := \frac{1}{2} \sum_{k \in \bar{\mathfrak{s}}} \sum_{d \in [D]} \int_0^T \varphi_{FB}(\alpha_d(t), s_{d,k}(t))^2 dt$$

for $\bar{\mathfrak{s}} := \{a, b\}$. Here φ_{FB} denotes a non-linear complementarity function (NCP function). We select

$$\varphi_{FB} : \mathbb{R}^2 \rightarrow \mathbb{R}, (a, b) \rightarrow \sqrt{a^2 + b^2} - a - b. \quad (8)$$

This function is known as the Fischer-Burmeister function. We obtain the following intermediate formulation

$$\min_{y, u, \alpha, s} J(y, u, \alpha, s) \quad (9a)$$

$$= \frac{1}{2} \|y - y_{des}\|_{L_{[0, T]}^2(X)}^2 + \frac{\gamma_u}{2} \|u\|_{L_{[0, T]}^2(U)}^2 + \gamma_{EC} J_{EC}(\alpha, s)$$

s.t.

$$\dot{y}(t) = Ay(t) + \sum_{d=1}^D \alpha_d(t) f_d(u(t)) \quad \text{a.e. } t \in [0, T], \quad (9b)$$

$$y(0) = y_0, \quad (9c)$$

$$0 \leq \alpha_d(t) \quad d \in [D], \text{ a.e. } t \in [0, T], \quad (9d)$$

$$1 = \sum_{d=1}^D \alpha_d(t) \quad \text{a.e. } t \in [0, T], \quad (9e)$$

$$-s_{d,a}(t) \leq (S(y(t)) - a_d) \quad d \in [D], \text{ a.e. } t \in [0, T], \quad (9f)$$

$$-s_{d,b}(t) \leq (b_d - S(y(t))) \quad d \in [D], \text{ a.e. } t \in [0, T], \quad (9g)$$

$$0 \leq s_{d,k}(t) \quad k \in \bar{\mathfrak{s}}, d \in [D], \text{ a.e. } t \in [0, T]. \quad (9h)$$

Similar to [13, Lemma 3.], we can observe the following connection between admissible points for the problems (5) and (9).

Lemma 3.8.

i) Let (y, u, α) be feasible for Problem (5), then (y, u, α, s) is admissible for Problem (9) with evaluation $J_{EC}(\alpha, s) = 0$, where we initialise $s \in C_{[0,T]}^0(\mathbb{R}^{2D})$ as

$$s_{d,a}(t) = \max\{0, a_d - S(y(t))\}, \quad (10a)$$

$$s_{d,b}(t) = \max\{0, S(y(t)) - b_d\}. \quad (10b)$$

ii) If (y, u, α, s) is admissible to (9) and $J_{EC}(\alpha, s) = 0$, then (y, u, α) is feasible for (5).

Proof. We prove each statement individually.

Let (y, u, α) be feasible for (5). The conditions (5b) – (5e) are equivalent to the postulates (9b) – (9e). Hence we remain with the discussion of (9f) – (9h). We first note that the solution of (9b) – (9c) satisfies $y \in C_{[0,T]}^0(X)$. As a consequence $S(y) \in C_{[0,T]}^0(\mathbb{R})$ holds. Furthermore the functions $s_{d,a}$ and $s_{d,b}$ in $C_{[0,T]}^0(\mathbb{R})$ are well posed. By construction the properties (9f) – (9h) for s_a and s_b are satisfied. We are left with the contribution in the cost function $J_{\gamma_{EC}}(\alpha, s) = 0$. By the fundamental lemma of calculus of variations and nonnegative arguments of $J_{EC}(\alpha, s) = 0$ is equivalent to $\varphi_{FB}(\alpha_d(t), s_{d,k}(t)) = 0$ for all $k \in \bar{s}$, $d \in [D]$ and f.a.e. $t \in [0, T]$. By the properties of the NCP functions the last statement can equivalently be restated as

$$\varphi_{FB}(\alpha_d(t), s_{d,k}(t)) = 0 \iff \alpha_d(t) \geq 0, s_{d,k}(t) \geq 0, \alpha_d(t) \cdot s_{d,k}(t) = 0.$$

This means we have to check the complementarity condition for each α_d and $s_{d,k}$. W.l.o.g. we only consider modes d such that $\alpha_d(t) > 0$ for a fixed $t \in [0, T]$. By (5f) – (5g) we conclude $S(y(t)) - a_d \geq 0$ and $b_d - S(y(t)) \geq 0$. Evaluation of (10a) – (10b) yields $s_{d,k}(t) = 0$. This completes the proof of the first statement.

Let now (y, u, α, s) be admissible to (9) with $J_{EC}(\alpha, s) = 0$. The conditions (9b) – (9e) imply (5b) – (5e). We remain with the discussion of (5f) – (5g). In the previous paragraph we established

$$\begin{aligned} J_{EC}(\alpha, s) = 0 &\iff \\ \alpha_d(t) \geq 0, s_{d,k}(t) \geq 0, \alpha_d(t)s_{d,k}(t) = 0, &\quad \forall k \in \bar{s}, d \in [D], \text{ a.e. } t \in [0, T]. \end{aligned}$$

The expressions (5f) – (5g) are restricting only if $\alpha_d(t) > 0$. Therefore w.l.o.g we consider just this case. By the complementarity condition this implies $s_{d,k}(t) = 0$. Hence (9f) – (9g) read as (5f) – (5g). This concludes the second statement and completes the proof. \square

Unfortunately formulation (9) still includes mixed control-state constraints (9f) – (9g) and is therefore impractical for the immediate derivation of an optimality system. We apply a Moreau-Yosida type regularization for these constraints. We conclude the

following formulation. We abbreviate

$$J_{MY}(y, s) := \frac{1}{2} \sum_{d=1}^D \int_0^T ((a_d - S(y(t)) - s_{d,a}(t))^+)^2 dt \\ + \frac{1}{2} \sum_{d=1}^D \int_0^T ((S(y(t)) - b_d - s_{d,b}(t))^+)^2 dt$$

and add regularization terms with respect to the control variables α and s together with nonnegative coefficients γ_α and γ_s . In the case $\gamma_s = 0$, we include the postulate of a uniform upper bound $S^\dagger > 0$ on $s(t)$ a.e. into the upcoming problem formulation.

$$\min_{y, u, \alpha, s} J(y, u, \alpha, s) \tag{11a}$$

$$= \frac{1}{2} \|y - y_{des}\|_{L^2_{[0,T]}(X)}^2 + \frac{\gamma_u}{2} \|u\|_{L^2_{[0,T]}(U)}^2 \\ + \gamma_{EC} J_{EC}(\alpha, s) + \gamma_{MY} J_{MY}(y, s) \\ + \frac{\gamma_\alpha}{2} \int_0^T \|\alpha(t)\|_{\mathbb{R}^D}^2 dt + \frac{\gamma_s}{2} \int_0^T \|s(t)\|_{\mathbb{R}^{2D}}^2 dt$$

s.t.

$$\dot{y}(t) = Ay(t) + \sum_{d=1}^D \alpha_d(t) f_d(u(t)) \quad \text{a.e. } t \in [0, T], \tag{11b}$$

$$y(0) = y_0, \tag{11c}$$

$$0 \leq \alpha_d(t) \quad d \in [D], \text{ a.e. } t \in [0, T], \tag{11d}$$

$$1 = \sum_{d=1}^D \alpha_d(t) \quad \text{a.e. } t \in [0, T], \tag{11e}$$

$$0 \leq s_{d,k}(t) \quad k \in \bar{s}, d \in [D], \text{ a.e. } t \in [0, T], \tag{11f}$$

$$0 \leq S^\dagger - s_{d,k}(t) \quad k \in \bar{s}, d \in [D], \text{ a.e. } t \in [0, T]. \tag{11g}$$

The formulation (11) now consists of the state equation together with box and equality constraints on the involved controls (u, α, s) .

Remark 3. The differentiability of the functions

$$J_1(y) = \frac{1}{2} \|y - y_{des}\|_{L^2_{[0,T]}(X)}^2, \\ J_2(u) = \frac{\gamma_u}{2} \|u\|_{L^2_{[0,T]}(U)}^2,$$

can be deduced from their well-posedness as function from X respectively U into \mathbb{R} . For a mild solution y we always obtain $C^0_{[0,T]}(X) \subset L^2_{[0,T]}(X)$. Hence the mapping J_k for $k \in [2]$ can be interpreted as a composition of continuously differentiable mappings.

We consider the following Lagrange function to derive first order optimality con-

ditions, e.g., [19]. The inputs satisfy $y \in C^1_{[0,T]}(X)$, $u \in L^2_{[0,T]}(U)$, $\alpha \in L^2_{[0,T]}(\mathbb{R}^D)$, $s \in L^2_{[0,T]}(\mathbb{R}^{2D})$, $p \in C^1_{[0,T]}(X^*)$, $\lambda \in L^2_{[0,T]}(\mathbb{R})$, $\rho_\alpha \in L^2_{[0,T]}(\mathbb{R}^D)$, $\rho_s \in L^2_{[0,T]}(\mathbb{R}^{2D})$, $\zeta_s \in L^2_{[0,T]}(\mathbb{R}^{2D})$. We then consider.

$$\begin{aligned}
& (y, u, \alpha, s, p, \lambda, \rho_\alpha, \rho_s, \zeta_s)^T \mapsto \mathcal{L}(y, u, \alpha, s, p, \lambda, \rho_\alpha, \rho_s, \zeta_s) \\
& = J(y, u, \alpha, s) \\
& - \int_0^T \langle p(t), \dot{y}(t) - Ay(t) - \sum_{d=1}^D \alpha_d(t) f_d(u(t)) \rangle_{X^*, X} dt - \langle p(0), y(0) - y_0 \rangle_{X^*, X} \\
& + \int_0^T \lambda(t) \left(\sum_{d=1}^D \alpha_d(t) - 1 \right) dt - \sum_{d=1}^D \int_0^T (\rho_\alpha)_d(t) \alpha_d(t) dt \\
& - \sum_{k \in \bar{s}} \sum_{d=1}^D \int_0^T (\rho_s)_{d,k}(t) s_{d,k}(t) dt - \sum_{k \in \bar{s}} \sum_{d=1}^D \int_0^T (\zeta_s)_{d,k}(t) (S^\dagger - s_{d,k}(t)) dt.
\end{aligned}$$

Here $p \in C^1_{[0,T]}(X^*)$ is supposed to be a classical solution of

$$-\dot{p}(t) = A^* p(t) + f^p(\hat{y}, \hat{s})(t), \quad (12a)$$

$$p(T) = 0, \quad (12b)$$

where we define the expression in the right side of (12a) by

$$\begin{aligned}
& f^p(\hat{y}, \hat{s})(t) = \Phi_{X,R}^{-1}(\hat{y}(t) - y_{des}(t)) \\
& - \gamma_{MY} \sum_{d=1}^D \left((a_d - S(\hat{y}(t)) - \hat{s}_{d,a}(t))^+ - (S(\hat{y}(t)) - b_d - \hat{s}_{d,b}(t))^+ \right) S'(\hat{y}(t)).
\end{aligned} \quad (13)$$

Here we utilise the inverse mapping of the Riesz representation $\Phi_{X,R}^{-1} : X \rightarrow X^*$, $y \mapsto y^*(\cdot) := \langle y, \cdot \rangle_X$. Since $S : X \rightarrow \mathbb{R}$ is supposed to be linear and continuous, we conclude that $S'(\hat{y}(t)) = S \in X^*$ for all $\hat{y}(t) \in X$ holds. Let now $(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}) \in C^1_{[0,T]}(X) \times L^2_{[0,T]}(U) \times L^2_{[0,T]}(\mathbb{R}^D) \times L^2_{[0,T]}(\mathbb{R}^{2D})$ as in the definition of (12) be fixed. By combination of the previous calculations we deduce the following KKT-system to characterise a candidate for optimality.

Theorem 3.9. *Let $(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}) \in C^1_{[0,T]}(X) \times L^2_{[0,T]}(U) \times L^2_{[0,T]}(\mathbb{R}^D) \times L^2_{[0,T]}(\mathbb{R}^{2D})$ be a local minimiser of (11). Furthermore we denote by $p \in C^1([0, T]; X^*)$ the classical solution to (12). Then there exist Lagrange multipliers $(\hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s) \in L^2_{[0,T]}(\mathbb{R}) \times L^2_{[0,T]}(\mathbb{R}^D) \times L^2_{[0,T]}(\mathbb{R}^{2D}) \times L^2_{[0,T]}(\mathbb{R}^{2D})$ such that the system of optimality conditions is fulfilled f.a.e. $t \in [0, T]$:*

$$\dot{\hat{y}}(t) = A\hat{y}(t) + \sum_{d=1}^D \alpha_d(t) f_d(\hat{u}(t)), \quad \hat{y}(0) = y_0, \quad (14a)$$

$$0 = \frac{\partial \mathcal{L}}{\partial y}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s), \quad (14b)$$

$$0 = \frac{\partial \mathcal{L}}{\partial u}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s), \quad (14c)$$

$$0 = \frac{\partial \mathcal{L}}{\partial \alpha_d}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s), \quad (14d)$$

$$0 = \frac{\partial \mathcal{L}}{\partial s_{d,k}}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s), \quad (14e)$$

$$0 = \sum_{d=1}^D \alpha_d(t) - 1, \quad (14f)$$

$$0 \leq \alpha_d(t), \quad 0 \leq (\rho_\alpha)_d(t), \quad 0 = \alpha_d(t)(\rho_\alpha)_d(t), \quad (14g)$$

$$0 \leq s_{d,k}(t), \quad 0 \leq (\rho_s)_{d,k}(t), \quad 0 = s_{d,k}(t)(\rho_s)_{d,k}(t), \quad (14h)$$

$$0 \leq S^\dagger - s_{d,k}(t), \quad 0 \leq (\zeta_s)_{d,k}(t), \quad 0 = (S^\dagger - s_{d,k}(t))(\zeta_s)_{d,k}(t). \quad (14i)$$

Proof. The provided problem can be interpreted as an instance of optimisation problems in Banach spaces. The associated optimality conditions for box constraints are formulated in, e.g., [16, Section 1.7.2.2 Corollary 1.2]. \square

We evaluate the corresponding optimality conditions (14b) - (14i) to obtain pointwise formulations. We start our calculations with $\frac{\partial \mathcal{L}}{\partial y}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s)$ in the direction $y \in C_{[0,T]}^1(X)$. We apply integration by parts and apply the definition of the adjoint state p .

$$\begin{aligned} 0 &= \frac{\partial \mathcal{L}}{\partial y}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s)(y) \quad (15) \\ &= \int_0^T \langle \hat{y}(t) - y_{des}(t), y(t) \rangle_X dt \\ &\quad - \gamma_{MY} \sum_{d=1}^D \int_0^T (a_d - S(\hat{y}(t)) - s_{d,a})^+ \langle S'(\hat{y}(t)), y(t) \rangle_{X^*, X} dt \\ &\quad + \gamma_{MY} \sum_{d=1}^D \int_0^T (S(\hat{y}(t)) - b_d - s_{d,a})^+ \langle S'(\hat{y}(t)), y(t) \rangle_{X^*, X} dt \\ &\quad - \int_0^T \langle p(t), \dot{y}(t) - Ay(t) \rangle_{X^*, X} dt - \langle p(0), y(0) \rangle_{X^*, X} \end{aligned}$$

$$\begin{aligned}
&= \int_0^T \langle \Phi_{X,R}^{-1}(\hat{y}(t) - y_{des}(t)), y(t) \rangle_{X^*,X} dt \\
&\quad - \gamma_{MY} \sum_{d=1}^D \int_0^T (a_d - S(\hat{y}(t)) - s_{d,a})^+ \langle S'(\hat{y}(t)), y(t) \rangle_{X^*,X} dt \\
&\quad + \gamma_{MY} \sum_{d=1}^D \int_0^T (S(\hat{y}(t)) - b_d - s_{d,a})^+ \langle S'(\hat{y}(t)), y(t) \rangle_{X^*,X} dt \\
&\quad - \int_0^T \langle -\dot{p}(t) - A^*p(t), y(t) \rangle_{X^*,X} dt - \langle p(T), y(T) \rangle_{X^*,X}.
\end{aligned}$$

This yields the variational formulation of (12a) – (12b). We proceed with the calculation of $\frac{\partial \mathcal{L}}{\partial u}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s)$ for an element $u \in L^2_{[0,T]}(U)$

$$\begin{aligned}
0 &= \frac{\partial \mathcal{L}}{\partial u}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s)(u) \\
&= \gamma_u \int_0^T \langle \hat{u}(t), u(t) \rangle_U dt + \sum_{d=1}^D \alpha_d(t) \int_0^T \langle p(t), f'_d(\hat{u}(t))u(t) \rangle_{X^*,X} dt.
\end{aligned} \tag{16}$$

Since the choice of u in (16) is arbitrary, we conclude by the fundamental lemma of calculus of variations, that f.a.e. $t \in [0, T]$

$$\gamma_u(\hat{u}(t)) + \sum_{d=1}^D \alpha_d(t) \Phi_{U,R}(f'_d(\hat{u}(t))^*p(t)) = 0_U. \tag{17}$$

We continue with $\frac{\partial \mathcal{L}}{\partial \alpha_d}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s)$ for $d \in [D]$ in the direction of $\alpha_d \in L^2_{[0,T]}(\mathbb{R})$

$$\begin{aligned}
0 &= \frac{\partial \mathcal{L}}{\partial \alpha_d}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s)(\alpha_d) \\
&= \gamma_{EC} \sum_{k \in \bar{s}} \int_0^T \varphi_{FB}(\hat{\alpha}_d(t), \hat{s}_{d,k}(t)) \frac{\partial \varphi_{FB}}{\partial \alpha}(\hat{\alpha}_d(t), \hat{s}_{d,k}(t)) \alpha_d(t) dt + \gamma_\alpha \int_0^T \hat{\alpha}_d(t) \alpha_d(t) dt \\
&\quad + \int_0^T \langle p(t), f_d(u(t)) \rangle_{X^*,X} \alpha_d(t) dt + \int_0^T \lambda(t) \alpha_d(t) dt - \int_0^T (\rho_\alpha)_d(t) \alpha_d(t) dt.
\end{aligned} \tag{18}$$

Repetition of the previous arguments yields f.a.e. $t \in [0, T]$ and for all $d \in [D]$

$$\gamma_{EC} \sum_{k \in \bar{s}} \varphi_{FB}(\hat{\alpha}_d(t), \hat{s}_{d,k}(t)) \frac{\partial \varphi_{FB}}{\partial \alpha}(\hat{\alpha}_d(t), \hat{s}_{d,k}(t)) + \gamma_\alpha \hat{\alpha}_d(t) + \lambda(t) - (\rho_\alpha)_d(t) = 0. \tag{19}$$

We conclude our derivative calculations by $\frac{\partial \mathcal{L}}{\partial s_{d,k}}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, p, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s)$ for $d \in [D]$

and $k \in \bar{s}$ along $s_{d,k} \in L^2_{[0,T]}(\mathbb{R})$

$$0 = \frac{\partial \mathcal{L}}{\partial s_{d,a}}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, \hat{p}, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s)(s_{d,a}) \quad (20)$$

$$\begin{aligned} &= \gamma_{EC} \int_0^T \varphi_{FB}(\hat{\alpha}_d(t), \hat{s}_{d,a}(t)) \frac{\partial \varphi_{FB}}{\partial s}(\hat{\alpha}_d(t), \hat{s}_{d,a}(t)) s_{d,a}(t) dt \\ &+ \gamma_s \int_0^T \hat{s}_{d,a}(t) s_{d,a}(t) dt - \gamma_{MY} \int_0^T (a_d - S(\hat{y}(t)) - \hat{s}_{d,a}(t))^+ s_{d,a}(t) dt \\ &- \int_0^T (\rho_s)_{d,a}(t) s_{d,a}(t) dt + \int_0^T (\zeta_s)_{d,a}(t) s_{d,a}(t) dt, \end{aligned}$$

$$0 = \frac{\partial \mathcal{L}}{\partial s_{d,b}}(\hat{y}, \hat{u}, \hat{\alpha}, \hat{s}, \hat{\lambda}, \hat{\rho}_\alpha, \hat{\rho}_s, \hat{\zeta}_s)(s_{d,b}) \quad (21)$$

$$\begin{aligned} &= \gamma_{EC} \int_0^T \varphi_{FB}(\hat{\alpha}_d(t), \hat{s}_{d,b}(t)) \frac{\partial \varphi_{FB}}{\partial s}(\hat{\alpha}_d(t), \hat{s}_{d,b}(t)) s_{d,b}(t) dt \\ &+ \gamma_s \int_0^T \hat{s}_{d,b}(t) s_{d,b}(t) dt - \gamma_{MY} \int_0^T (S(\hat{y}(t)) - b_d - \hat{s}_{d,b}(t))^+ s_{d,b}(t) dt \\ &- \int_0^T (\rho_s)_{d,b}(t) s_{d,b}(t) dt + \int_0^T (\zeta_s)_{d,b}(t) s_{d,b}(t) dt. \end{aligned}$$

By a variational argument we conclude the pointwise equality for all $d \in [D]$ and f.a.e. $t \in [0, T]$

$$0 = \gamma_{EC} \varphi_{FB}(\hat{\alpha}_d(t), \hat{s}_{d,a}(t)) \frac{\partial \varphi_{FB}}{\partial s}(\hat{\alpha}_d(t), \hat{s}_{d,a}(t)) + \gamma_s \hat{s}_{d,a}(t) \quad (22)$$

$$- \gamma_{MY} (a_d - S(\hat{y}(t)) - \hat{s}_{d,a}(t))^+ - (\rho_s)_{d,a}(t) + (\zeta_s)_{d,a}(t),$$

$$0 = \gamma_{EC} \varphi_{FB}(\hat{\alpha}_d(t), \hat{s}_{d,b}(t)) \frac{\partial \varphi_{FB}}{\partial s}(\hat{\alpha}_d(t), \hat{s}_{d,b}(t)) + \gamma_s \hat{s}_{d,b}(t) \quad (23)$$

$$- \gamma_{MY} (S(\hat{y}(t)) - b_d - \hat{s}_{d,b}(t))^+ - (\rho_s)_{d,b}(t) + (\zeta_s)_{d,b}(t).$$

The system of Theorem 3.9 and the deduced pointwise formulations (17) – (23) present a system of necessary conditions to characterise the minimiser of the surrogate problem (11). For a suitable choice of the penalty parameters γ_s, γ_{EC} and γ_{MY} we expect local minimisers of (5) to provide good solutions for problem (11). With the help of the presented rounding strategies, we afterwards attempt to construct qualified points of (4) and ultimately (1). The algorithmic solution of the derived conditions is subject of the next section.

4. Algorithmic design

In this section we introduce our algorithmic approach to solve (11) by utilising (14a) – (14i). Similar to [13], we utilise a semismooth Newton method, [16]. For that purpose we reformulate the equilibrium conditions (14g) – (14i) on $\alpha_d, s_{d,k}$ and their associated multipliers $(\rho_\alpha)_d, (\rho_s)_{d,k}, (\zeta_s)_{d,k}$ as nonsmooth equality constraints via application of a suitable NCP function. We again select φ_{FB} , c.f. (8). Furthermore consider $M : V \rightarrow$

W , where we set

$$V := C_{[0,T]}^1(X) \times C_{[0,T]}^1(X^*) \times L_{[0,T]}^2(U) \times L_{[0,T]}^2(\mathbb{R}^D) \times L_{[0,T]}^2(\mathbb{R}^{2D}) \\ \times L_{[0,T]}^2(\mathbb{R}) \times L_{[0,T]}^2(\mathbb{R}^D) \times L_{[0,T]}^2(\mathbb{R}^{2D}) \times L_{[0,T]}^2(\mathbb{R}^{2D}).$$

Similarly we define the image space by

$$W := C_{[0,T]}^1(X) \times C_{[0,T]}^1(X^*) \times L_{[0,T]}^2(U) \times L_{[0,T]}^2(\mathbb{R}^D) \times L_{[0,T]}^2(\mathbb{R}^{2D}) \\ \times L_{[0,T]}^2(\mathbb{R}) \times L_{[0,T]}^2(\mathbb{R}^D) \times L_{[0,T]}^2(\mathbb{R}^{2D}) \times L_{[0,T]}^2(\mathbb{R}^{2D}).$$

Then M describes the nonlinear mapping, which contains the previously discussed conditions, (14a) – (14i), in lexicographical order, i.e., solutions of the optimality system are exactly the roots of M . We define M via

$$M : V \rightarrow W, (y, p, u, \alpha, s, \lambda, \rho_\alpha, \rho_s, \zeta_s)^T \mapsto M(y, p, u, \alpha, s, \lambda, \rho_\alpha, \rho_s, \zeta_s).$$

We denote by ∂M the generalised gradient according to Clarke, c.f. [6].

The proposed algorithm consists of two integral parts. In the inner while loop the nonlinear root problem with respect to M is solved by the semismooth Newton method. In the subsequent block the parameters of the path following method are adapted according to the ratio of the violation of the equilibrium constraints and the satisfaction of the state constraints. The procedure is performed in an outer loop until the residuum of optimality system derived from Theorem 3.9 is less than a given threshold Tol_N . After successful termination also selected tolerances for Tol_{EC} and Tol_{MY} are achieved.

Algorithm 1. Let $x = (y, p, u, \alpha, s, \lambda, \rho_\alpha, \rho_s, \zeta_s)^T \in V$. Set $k = 0$ and initialise $x_0 \in V$ together with $\gamma_{EC}, \gamma_{MY} > 0$. Select $\text{Tol}_N, \text{Tol}_{EC}, \text{Tol}_{MY} > 0$ in combination with $\delta_{EC}, \delta_{MY} > 1$.

```

while  $\|M(x_k)\|_W \geq \text{Tol}_N$  or  $J_{EC}(\alpha_k, s_k) \geq \text{Tol}_{EC}$  or  $J_{MY}(y_k, s_k) \geq \text{Tol}_{MY}$  do
  while  $\|M(x_k)\|_W \geq \text{Tol}_N$  do
    | Select an element  $N \in \partial M(x_k)$ ;
    | Solve  $Nd_k = -M(x_k)$ ;
    | Update  $x_{k+1} = x_k + d_k$ ;
  end
  if  $J_{EC}(\alpha_{k+1}, s_{k+1}) \geq J_{MY}(y_{k+1}, s_{k+1})$  then
    | Update  $\gamma_{EC} = \gamma_{EC} \cdot \delta_{EC}$ ;
  else
    | Update  $\gamma_{MY} = \gamma_{MY} \cdot \delta_{MY}$ ;
  end
  Increment  $k = k + 1$ ;
end

```

The presented scheme is typically only locally convergent. We attempt to heuristically increase the region of convergence by suitable modification of the involved step-size. For that purpose, we expand the cost function for given parameters $\gamma_\# > 0$, $\# \in \{y, \{\alpha, \Sigma\}, \{\alpha, LB\}, \{s, LB\}, \{s, UB\}\}$, towards the merit function, which apart from

cost also takes feasibility aspects into account.

$$\begin{aligned}
M(y, u, \alpha, s) &:= J(y, u, \alpha, s) \\
&+ \gamma_y \left\| y - y_0 - \int_0^\cdot Ay(s) + \sum_{d=1}^D \alpha_d(s) f_d(u(s)) ds \right\|_{L^1_{[0,T]}(X)} \\
&+ \gamma_\Sigma \left\| 1 - \sum_{d=1}^D \alpha_d \right\|_{L^1_{[0,T]}(\mathbb{R})} + \gamma_{\alpha, LB} \sum_{d=1}^D \left\| (-\alpha_d)^+ \right\|_{L^1_{[0,T]}(\mathbb{R})} \\
&+ \gamma_{s, LB} \sum_{k \in \bar{s}} \sum_{d=1}^D \left\| (-s_{d,k})^+ \right\|_{L^1_{[0,T]}(\mathbb{R})} + \gamma_{s, UB} \sum_{k \in \bar{s}} \sum_{d=1}^D \left\| (s_{d,k} - S^\dagger)^+ \right\|_{L^1_{[0,T]}(\mathbb{R})}.
\end{aligned} \tag{24}$$

For a given stepwidth $t > 0$ and search direction $d \in V$ we define

$$M_{t,d}(y, u, \alpha, s) := M(y + td_y, u + td_u, \alpha + td_\alpha, s + td_s).$$

Hereby $d_\#$ denotes the vector, which consists only out of the entries of d associated with the corresponding symbol $\# \in \{y, u, \alpha, s\}$.

Algorithm 2. Fix $\beta \in (0, 1)$ and $0 < \text{Max}_{Iter} < \infty$. Initialise $k = 0$ and $t = 1$.

```

while  $M_{t,d}(y, u, \alpha, s) \geq M(y, u, \alpha, s)$  and  $k < \text{Max}_{Iter}$  do
  | Update  $t = t \cdot \beta$ ;
  | Update  $k = k + 1$ ;
end
if  $M_{t,d}(y, u, \alpha, s) < M(y, u, \alpha, s)$  then
  | Accept stepwidth  $t = t$ ;
else
  | Reset stepwidth  $t = 1$ ;
end

```

We substitute the update step in Algorithm 1 by $x_{k+1} = x_k + td_k$, where the stepwidth t results from Algorithm 2. For alternative globalisation strategies and rigorous convergence results of the semismooth Newton algorithm we refer to, e.g., [16]. The convergence theory for the path following technique applied to pure state constraints is discussed in [15], while a differentiable penalty approach towards MPECs is topic of [7]. However, a rigorous combined convergence study is out of scope for this paper. Instead we report on our numerical experiment.

5. Numerical results

In this section we present numerical results for Algorithm 1, where the stepsize is selected according to Algorithm 2, on a selected instance. As a benchmark test, we attempt to numerically recover distinguished input parameters to the problem proposed in Example 2.1.

We select a one-dimensional space domain $\Omega := (0, L)$ with an $L > 0$. Hence the domain under consideration for our problem is denoted by $Q := (0, L) \times (0, T)$. Continuing

the indirect approach, we solve the state equation (14a) and corresponding adjoint system (12) by the implicit Euler method. As proposed in Example 2.1 we investigate $D = 2$ modes, whose areas of effect are declared by $\Omega_1 = (\frac{1}{6}, \frac{2}{6})$ and $\Omega_2 = (\frac{4}{6}, \frac{5}{6})$. The switching rule is formulated via

$$C : L^2(\Omega) \rightarrow [2], y \mapsto \begin{cases} 1, & \text{if } \int_{\Omega} y(x) dx \leq 0.1, \\ 2, & \text{if else } \int_{\Omega} y(x) dx > 0.1. \end{cases}$$

Hence the switching threshold is selected as $\sigma = 0.1$. We aim to reconstruct the (constant) desired distributed control $u_{des}(x, t) = 7.5$ by tracking the distance towards the associated state y_{des} obtained by evaluation of (1). Carefully note that the acting control takes the area of effect associated with each mode into account and is therefore given by $u_{des,act}(x, t) = u_{des}(x, t)\chi_{\Omega_{d(t)}}(x)$. We denote the associated binary control by $\omega_{des}(t)$.

Note that we will not insist to verify the Assumptions 2 and 3 for this example. Instead, we consider generalised solutions according to Remark 1 and understand the tracking task to be successful if the optimal value approaches zero for one of all possible branches. Since the branches are implicitly selected by the reformulation, we will perform a post processing step of the switching threshold σ in order to identify the solution branch y possessing the desired tracking properties towards y_{des} when we pass back from the reformulation to the original problem formulation. Accordingly, the tracking goal will be considered successful, if the objective value is sufficiently small for a modified switching threshold σ^* sufficiently close to σ .

We initialise the system with $u_{start}(x, t) = 1$, $\alpha_{start}(t) = (0.5, 0.5)^T$ and $s_{start}(t) = (10^{-9}, 10^{-9})^T$ for all $(x, t) \in Q$ respectively $t \in [0, T]$. The initial input is completed by the multiplier $\lambda(t) = 0$, $\rho_{\alpha}(t) = (0, 0)^T$, $\rho_s(t) = (0, 0)^T$, $\zeta_s(t) = (0, 0)^T$. We initialise the respective penalty parameters as $\gamma_u = 10^{-9}$, $\gamma_{\alpha} = 10^{-7}$, $\gamma_s = 10^{-7}$, $\gamma_{EC} = 10^{-8}$, $\gamma_{MY} = 10^{-8}$. We apply a spatial and time discretisation with gridwidth $\frac{1}{72}$ for each dimension on an equidistant grid.

The evaluation of the merit function (24) includes, besides the evaluation of the cost function, the primal admissibility of the involved variables. This covers the satisfaction of the discretised state equation, the sum constraint on the relaxed control together with the fulfilment of the lower and upper bounds on the entries of the relaxed control respectively slack control. We set all parameters involved in M to 1.

We abort the current Newton iteration of the penalty homotopy once the condition $\|\nabla \mathcal{L}(y, u, \alpha, s, p, \lambda, \rho_{\alpha}, \rho_s, \zeta_s)\| < \text{Tol}_N = 10^{-4}$ is satisfied. Furthermore we initialise $\text{Tol}_{EC} = \text{Tol}_{MY} = 10^{-5}$ with $\delta_{EC} = 10$ and $\delta_{MY} = 10$. The value of S^{\dagger} is 10^1 .

In our computational experiments the proposed algorithm terminates with the results presented in Figures 1a – 1f. In Figures 1a – 1b the desired state, created by the mentioned input control u_{des} , for the initial problem formulation together with the computed state via the evaluation of (1) is presented. The algorithmically reconstructed state fits well to the desired state, if the switching threshold $\sigma = 0.1$ is adapted to $\sigma^* = 0.95$.

In Figure 1c the evolution of the Euclidean norm of the gradient of the Lagrangian in combination with the cost and merit function is displayed. The black vertical lines indicate an adaption of the involved penalty parameters γ_{EC} and γ_{MY} according to the rule formulated in Algorithm 1. In the first section of the penalty homotopy the cost and merit function decay until the associated curves nearly intersect. This indicates sufficient primal admissibility of the obtained point. Likewise the norm of the gradient

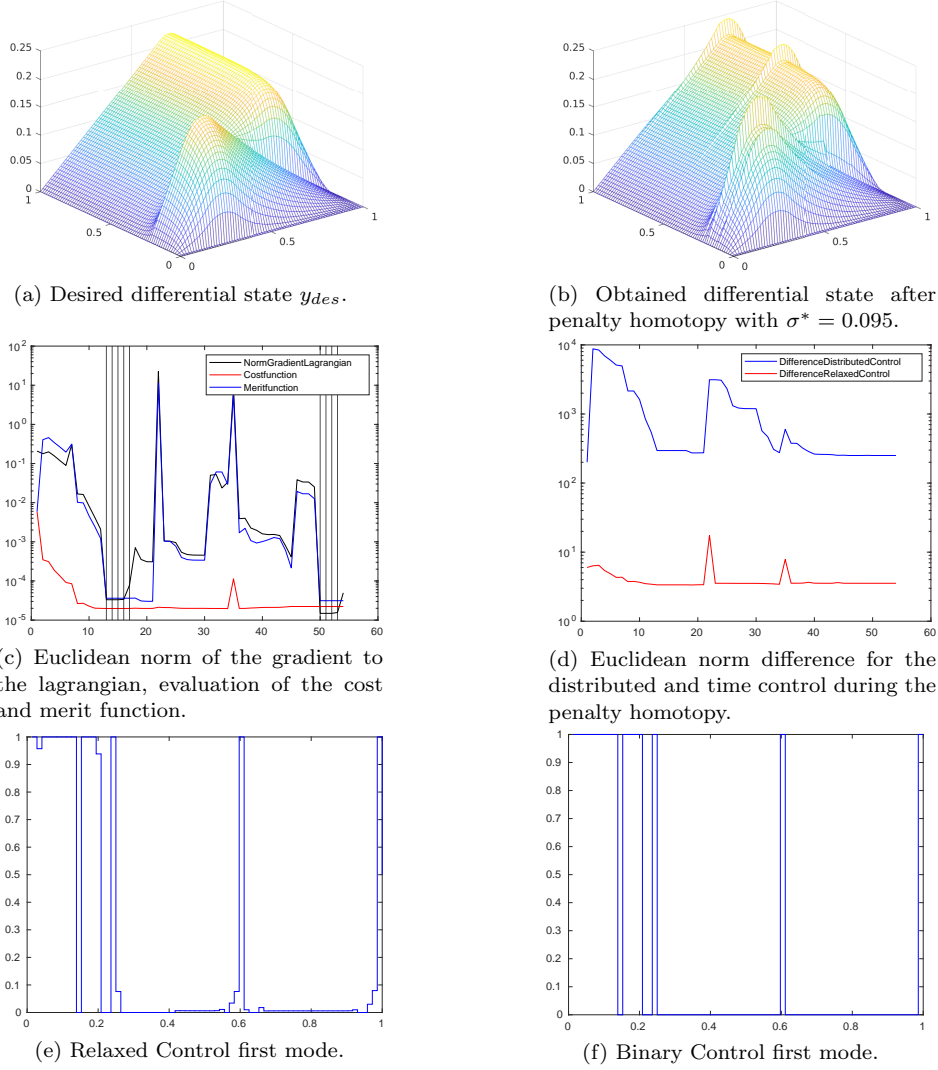


Figure 1. Algorithmic results

of the Lagrangian gradually decreases towards Tol_N . Then the algorithm proceeds with several adaptations of the penalty parameters γ_{EC} and γ_{MY} . This process is visualised by regions with few down to one iteration for the underlying semismooth Newton algorithm. The final iteration of the penalty homotopy is coined by a spike of the merit function at the beginning, which indicates primal infeasibility of the visited points due to the visible gap between the curves of the cost and merit function. Until termination of the algorithm the merit function and norm of the gradient of the Lagrangian decrease, whereas the cost remains more or less stable.

In Figure 1d the Euclidean norm difference for the obtained distributed and relaxed control towards $u_{des,act}$ and ω_{des} is displayed. During the iterations of the algorithm both distances decrease, which indicates a successful reconstruction of the aforementioned input parameters.

The Figures 1e – 1f contain the obtained relaxed control after termination of the algorithm and the associated binary control after application of the rounding scheme (SUR-SOS-VC). The corresponding results for the $\alpha_2(t)$ and $\omega_2(t)$ can be deduced via $\alpha_2(t) = 1 - \alpha_1(t)$.

6. Conclusion

We presented a promising algorithmic approach to solve linear systems with implicit switching formulated in a semigroup setting. We reformulated the original problem as a MIOCP instance with vanishing constraints by methods of disjunctive programming. Based on the theory of the associated rounding schemes, we proceeded with the solution of the corresponding relaxed formulation. During this procedure we embedded the appearing vanishing constraints into the framework of equilibrium constraints. We performed a final penalization step to elaborate a system of necessary optimality conditions, which formed the core of our indirect approach. In the process we also provided theoretical justification for the selected method in the form of an approximation result. A numerical experiment was conducted to underline the promising nature of the demonstrated method. In a succeeding step the generalization of the investigated dynamics to cover more applications seems possible. This is an important detail to keep in mind, especially in comparison to the approach proposed in [13], which is so far restricted to the application on parabolic PDEs. Also the implementation of the numerical routine suggested in the paper at hand appears to require less effort than the time transformation method utilised in [13]. Our numerical approach motivates a detailed convergence study of the surrogate penalty formulations. Also, a closer numerical comparison to the approach in [13] is desirable.

7. Acknowledgements

This research was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under projects A03 of the Sonderforschungsbereich/ Transregio 154 “Mathematical Modelling, Simulation and Optimisation using the Example of Gas Networks” (project ID: 239904186) and Germany’s Excellence Strategy - The Berlin Mathematics Research Center MATH+ (EXC-2046/1, project ID: 390685689).

8. Data availability statement

The data that support the findings of this study are available from the corresponding author, C.K., upon reasonable request.

References

- [1] W. Achtziger, C. Kanzow, and T. Hoheisel. On a relaxation method for mathematical programs with vanishing constraints. *GAMM-Mitt.*, 35(2):110–130, 2012.
- [2] A. Bensoussan, G. Da Prato, M. C. Delfour, and S. Mitter. *Representation and control of infinite dimensional systems*. Systems & Control: Foundations & Applications. Birkhäuser Boston, Inc., Boston, MA, second edition, 2007.
- [3] M. Boccadoro, Y. Wardi, M. Egerstedt, and E. Verriest. Optimal control of switching surfaces in hybrid dynamical systems. *Discrete Event Dyn. Syst.*, 15(4):433–448, 2005.
- [4] H. G. Bock, C. Kirches, A. Meyer, and A. Potschka. Numerical solution of optimal control problems with explicit and implicit switches. *Optim. Methods Softw.*, 33(3):450–474, 2018.
- [5] L. Cesari. *Optimization—theory and applications*, volume 17 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, 1983. Problems with ordinary differential equations.

- [6] F. H. Clarke. Generalized gradients and applications. *Trans. Amer. Math. Soc.*, 205:247–262, 1975.
- [7] C. Clason, Y. Deng, P. Mehrlitz, and U. Prüfert. Optimal control problems with control complementarity constraints: existence results, optimality conditions, and a penalty method. *Optim. Methods Softw.*, 35(1):142–170, 2020.
- [8] Aleksej F. Filippov. Differential equations with discontinuous righthand sides. In *Mathematics and Its Applications*, 1988.
- [9] M. L. Flegel and C. Kanzow. Abadie-type constraint qualification for mathematical programs with equilibrium constraints. *J. Optim. Theory Appl.*, 124(3):595–614, 2005.
- [10] M. L. Flegel and C. Kanzow. On the Guignard constraint qualification for mathematical programs with equilibrium constraints. *Optimization*, 54(6):517–534, 2005.
- [11] M. Gerdts. A variable time transformation method for mixed-integer optimal control problems. *Optimal Control Appl. Methods*, 27(3):169–182, 2006.
- [12] F. M. Hante. *Hybrid Dynamics Comprising Modes Governed by Partial Differential Equations: Modeling, Analysis and Control for Semilinear Hyperbolic Systems in One Space Dimension*. PhD thesis, Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU), 2010.
- [13] F. M. Hante and C. Kuchler. An algorithmic framework for optimal control of hybrid dynamical systems with parabolic PDEs. *TRR154 Preprint 519*, 2023.
- [14] F. M. Hante, G. Leugering, A. Martin, L. Schewe, and M. Schmidt. Challenges in optimal control problems for gas and fluid flow in networks of pipes and canals: from modeling to industrial application. In *Industrial mathematics and complex systems*, Ind. Appl. Math., pages 77–122. Springer, Singapore, 2017.
- [15] M. Hintermüller and K. Kunisch. Feasible and noninterior path-following in constrained minimization with low multiplier regularity. *SIAM J. Control Optim.*, 45(4):1198–1221, 2006.
- [16] M. Hinze, R. Pinnau, M. Ulbrich, and S. Ulbrich. *Optimization with PDE constraints*, volume 23 of *Mathematical Modelling: Theory and Applications*. Springer, New York, 2009.
- [17] T. Hoheisel. *Mathematical Programs with Vanishing Constraints*. doctoralthesis, Universität Würzburg, 2009.
- [18] K. Ito and K. Kunisch. Semi-smooth Newton methods for state-constrained optimal control problems. *Systems Control Lett.*, 50(3):221–228, 2003.
- [19] J. Jahn. *Introduction to the theory of nonlinear optimization*. Springer, Berlin, third edition, 2007.
- [20] C. Kirches, F. Lenders, and P. Manns. Approximation properties and tight bounds for constrained mixed-integer optimal control. *SIAM J. Control Optim.*, 58(3):1371–1402, 2020.
- [21] V. Laha, V. Singh, Y. Pandey, and S. K. Mishra. Nonsmooth mathematical programs with vanishing constraints in Banach spaces. In *High-dimensional optimization and probability—with a view towards data science*, volume 191 of *Springer Optim. Appl.*, pages 395–417. Springer, Cham, [2022] ©2022.
- [22] H. W. J. Lee, K. L. Teo, V. Rehbock, and L. S. Jennings. Control parametrization enhancing technique for optimal discrete-valued control problems. *Automatica J. IFAC*, 35(8):1401–1407, 1999.
- [23] Z. Luo, J. Pang, and D. Ralph. *Mathematical programs with equilibrium constraints*. Cambridge University Press, Cambridge, 1996.
- [24] P. Manns and C. Kirches. Improved regularity assumptions for partial outer convexification of mixed-integer PDE-constrained optimization problems. *ESAIM Control Optim. Calc. Var.*, 26:Paper No. 32, 16, 2020.
- [25] P. Manns, C. Kirches, and F. Lenders. Approximation properties of sum-up rounding in the presence of vanishing constraints. *Math. Comp.*, 90(329):1263–1296, 2021.
- [26] A. Pazy. Semigroups of operators in Banach spaces. In *Equadiff 82 (Würzburg, 1982)*, volume 1017 of *Lecture Notes in Math.*, pages 508–524. Springer, Berlin, 1983.
- [27] F. Rüffler. *Control and Optimization for Switched Systems of Evolution Equations*. PhD

- thesis, Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU), 2019.
- [28] S. Sager, H. Bock, and M. Diehl. The integer approximation error in mixed-integer optimal control. *Math. Program.*, 133(1-2):1–23, 2012.
 - [29] T. I. Seidman. The residue of model reduction. In Rajeev Alur, Thomas A. Henzinger, and Eduardo D. Sontag, editors, *Hybrid Systems III: Verification and Control, Proceedings of the DIMACS/SYCON Workshop on Verification and Control of Hybrid Systems, October 22-25, 1995, Rutgers University, New Brunswick, NJ, USA*, volume 1066 of *Lecture Notes in Computer Science*, pages 201–208. Springer, 1995.
 - [30] T. I. Seidman. Optimal control of a diffusion/reaction/switching system. *Evol. Equ. Control Theory*, 2(4):723–731, 2013.
 - [31] Thomas I. Seidman. Feedback modal control of partial differential equations. In *Optimal control of coupled systems of partial differential equations*, volume 158 of *Internat. Ser. Numer. Math.*, pages 239–253. Birkhäuser Verlag, Basel, 2009.
 - [32] G. Wachsmuth. Mathematical programs with complementarity constraints in Banach spaces. *J. Optim. Theory Appl.*, 166(2):480–507, 2015.