

A DIFFERENTIAL-ALGEBRAIC RICCATI EQUATION FOR APPLICATIONS IN FLOW CONTROL*

JAN HEILAND†

Abstract. We investigate existence and structure of solutions to quadratic control problems with semiexplicit differential-algebraic constraints as they appear in the modeling of flows. We introduce a decoupled representation of the state and the formal adjoint equations to identify the conditions for the existence of solutions. We derive a differential-algebraic Riccati equation formulated in the original coefficients for the efficient numerical computation of the optimal solution.

Key words. optimal control, DAE, differential Riccati equation, formal adjoint equation

AMS subject classifications. 34H05, 49J15

DOI. 10.1137/151004963

1. Introduction. We investigate the solvability of the first order optimality conditions for the minimization of a cost functional of type

$$(1.1) \quad \mathcal{J}(v, u) = \frac{1}{2}[v - v^*]^\top V[v - v^*] \Big|_{t=T} + \frac{1}{2} \int_0^T [v - v^*]^\top W[v - v^*] + u^\top R u \, dt$$

subject to the structured linear time-varying differential-algebraic state equations

$$(1.2) \quad \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \frac{d}{dt} \begin{bmatrix} v \\ p \end{bmatrix} - \begin{bmatrix} A & J_1^\top \\ J_2 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} - \begin{bmatrix} B_1 \\ 0 \end{bmatrix} u = \begin{bmatrix} f \\ g \end{bmatrix}, \quad v(0) = \alpha.$$

Such setups occur in the modeling of optimal flow control problems where v and p stand for the velocity and the pressure of the fluid flow.

For the general case of linear-quadratic optimal control problems with constraints given through differential-algebraic equations (DAE) like

$$(1.3) \quad \mathcal{E} \dot{x} = \mathcal{A}x + \mathcal{B}u \quad \text{and} \quad l_0(x(0)) = 0,$$

it is known (cf. [6, 33, 37]) that the *formal* adjoint equation

$$(1.4) \quad -\mathcal{E}^\top \dot{\lambda} = (\mathcal{A}^\top + \dot{\mathcal{E}}^\top)\lambda + h(x) \quad \text{and} \quad l_T(x(T), \lambda(T)) = 0$$

may not be solvable because of inconsistencies or lack of regularity in the terminal conditions l_T or the inhomogeneity h which are defined through the cost functional of the optimization. In fact, it may happen that the considered optimal control problem has a solution while the *formal* optimality conditions, which (1.4) is a part of, are not solvable [33]. Thus, in order to employ the *formal* optimality conditions to solve the optimal control problem, one needs to establish solvability of (1.4) a priori; cf. [6, 8, 16, 36] and in particular [30] for the linear-quadratic case.

As a general alternative approach, one can reformulate the DAE (1.2) or (1.3) and formulate the optimality conditions for the resulting *decoupled* or *strangeness-free*

*Received by the editors January 21, 2015; accepted for publication (in revised form) January 15, 2016; published electronically March 24, 2016. This work was supervised by Volker Mehrmann and supported by a grant from *Studienstiftung des deutschen Volkes*.

<http://www.siam.org/journals/sicon/54-2/100496.html>

†Max Planck Institute for Dynamics of Complex Technical Systems, Magdeburg, Germany (heiland@mpi-magdeburg.mpg.de).

system; cf. [15, 32, 33]. In [33] it has been shown that for linear-quadratic systems, the first order optimality system for the *strangeness-free* formulation is solvable if and only if the optimal control problem is solvable. Therefore they are referred to as *true* optimality conditions, the more so as the solvability issues in the *formal* conditions are but an artifact of the chosen formulation.

However, in practice, the use of the *true* optimality conditions may not be feasible or desirable. First, the reformulations are based on implicit functions or projections that may be expensive to compute, whereas an approximate realization of the transformations can introduce a systematic error to the equations; cf. [3]. Second, a formulation in the original equations keeps the physical meaning of the computed quantities. Finally, for an efficient implementation, a possibly necessary reformulation can be tailored to the solution of the specific problem rather than to the derivation of theoretical results.

Except from the contributions in [5, 6], most specific approaches to necessary and sufficient optimality conditions for DAE constrained control problems are based on reformulations of the original equations. This concerns the works that exploit the special structure of semiexplicit equations (cf. [16, 20, 21, 43]) or consider linear DAEs with a *properly stated leading term*; cf. [5, 6, 7, 9, 36, 39]. The special case of Riccati-based feedback was investigated in [37].

We use the particular structure of the considered DAE (1.2) and the cost functional (1.1) to show when the solution of the untransformed or *formal* optimality system coincides with the solution of a *strangeness-free* reformulation and, thus, provide necessary and sufficient optimality conditions. The innovations we propose are of practical impact, as the theoretical sufficiency and necessity results are already contained in [5, 6]. We show that the gap between the *formal* and the *true* optimality conditions is given by a single condition on the cost functional that can be easily checked and enforced in practice. In view of numerical simulations, we provide a novel differential-algebraic Riccati decoupling of the *formal* optimality conditions that makes the computation of the optimal control in feedback form computable also for large-scale systems. We illustrate the capabilities of the approach via an example of output tracking with a time-varying target signal with linearized Navier–Stokes equations. To our knowledge the numerical approximation of these feedback control problems has not been considered before. However, there is extensive work on open loop control in this context; cf., e.g., [22, 29]. Concerning the algorithms employed here, our work relates to recent work on computing stabilizing controllers for laminar flows via large-scale Riccati equations [10].

The paper is structured as follows. In the beginning, we provide basic concepts from DAE theory in order to define the tractability index. In section 3, we introduce the considered class of linear *index-2* equations and decouple them using the projectors defined together with the tractability index. Section 4 is devoted to the linear-quadratic optimal control problem modeling a velocity tracking problem. We state the formal optimality system and formulate sufficient conditions for existence of solutions. Additionally, we provide a feedback representation of the optimal control that can be computed without resorting to projections or variable transformations. We illustrate the capabilities of the proposed methodology in an example for linearized Navier–Stokes equations in section 5.

2. Basic concepts from DAE calculus. There are several concepts of so-called indices [41] to classify a DAE. We will call on the *tractability index* since it will provide a decoupling of the state equations only by scaling the equations from

the left. This particular property will later enable us to derive a structure preserving decoupling of the formal optimality conditions.

To introduce the concepts, we consider a generic linear DAE of the form

$$(2.1) \quad \mathcal{E}_A \frac{d}{dt}(\mathcal{E}_D x) - \mathcal{A}x = q \quad \text{on } (0, T)$$

for a state $x(t) \in \mathbb{R}^{n_x}$ and with $\mathcal{E}_A^\top, \mathcal{E}_D \in \mathcal{C}(0, T; \mathbb{R}^{n_x, n_d})$, $\mathcal{A} \in \mathcal{C}(0, T; \mathbb{R}^{n_x, n_x})$, and $q \in \mathcal{C}(0, T; \mathbb{R}^{n_x})$, $n_x, n_d \in \mathbb{N}$, where, e.g., $\mathcal{C}(0, T; \mathbb{R}^{n_x, n_d})$ is short for $\mathcal{C}([0, T], \mathbb{R}^{n_x, n_d})$ and where, e.g., \mathbb{R}^{n_x, n_d} denotes the space of matrices with n_x rows, n_d columns, and real entries.

DEFINITION 2.1 (cf. [40, Def. 2.1]). *The DAE (2.1) has a properly stated leading term if*

$$(2.2) \quad \mathbb{R}^{n_x} = \text{im } \mathcal{E}_A(t) \oplus \ker \mathcal{E}_D(t)$$

for all $t \in (0, T)$.

DEFINITION 2.2 (cf. [40, eqn. (2.2)]). *Consider (2.1). Given*

$$\mathcal{E}_0 := \mathcal{E}_A \mathcal{E}_D \quad \text{and} \quad \mathcal{A}_0 := \mathcal{A}$$

for $i = 0, 1, 2, \dots$, define a sequence of subspaces and matrices via

$$\begin{aligned} N_i &:= \ker \mathcal{E}_i, \\ S_i &:= \{x \in \mathbb{R}^{n_x} : \mathcal{A}_i x \in \text{im } \mathcal{E}_i\}, \\ \mathcal{Q}_i &:= \mathcal{P}_{[N_i]} \text{ (projector onto } \ker \mathcal{E}_i\text{)}, \\ \mathcal{P}_i &:= I - \mathcal{Q}_i, \end{aligned}$$

and

$$\mathcal{E}_{i+1} := \mathcal{E}_i + \mathcal{A}_i \mathcal{Q}_i \quad \text{and} \quad \mathcal{A}_{i+1} := \mathcal{A}_i \mathcal{P}_i.$$

The definition of the spaces and matrices holds pointwise for $t \in [0, T]$.

We can now define the tractability index.

DEFINITION 2.3 (see [40, Def. 2.2]). *Consider (2.1) and let \mathcal{E}_A and \mathcal{E}_D fulfill Definition 2.1. Consider a sequence of operators and subspaces as defined in Definition 2.2. Then, the DAE (2.1) has tractability index $i_\mu = \mu$ if there is a $\mu \in \mathbb{N}$ such that $\dim(N_j \cap S_j) = d_j > 0$, for $j = 0, 1, \dots, \mu - 1$, and $\dim(N_\mu \cap S_\mu) = 0$.*

Remark 2.4. Note that the choice of the projectors in the chain in Definition 2.2 is not unique, since only the image but not the kernels is specified. This freedom does not affect the determination of the tractability index as in Definition 2.3, provided that the chosen projectors are *admissible* [38, Thm. 2.8].

For completeness, we mention the commonly used *differentiation index* [41, Def. 1] and its relation to the *tractability index*.

Remark 2.5. The definition of the *differentiation index* i_ν requires a certain smoothness of the coefficients. If both the *differentiation index* and the *tractability index* are defined, then they coincide; cf. [9, Rem. 2.3].

The definitions of the *tractability index* and the *differentiation index* apply for uncontrolled systems. Index concepts for systems that include inputs, like the considered system (1.2), have to take into account the relation of inputs u and variables (v, p) [15]. We will investigate the index only for the uncontrolled system but we will use the obtained projectors to decouple the controlled state equations and the formal adjoint equations.

3. Decoupling of the state equations. In this section, we introduce the state equations, namely, linear time-varying so-called index-2 systems [25]. We use an explicit realization of the projectors and spaces from Definition 2.2 to decouple the state equations and to state existence of solutions.

Problem 3.1. Consider $T > 0$ and $n_u, n_v,$ and $n_p \in \mathbb{N}$, where $n_v > n_p$. Consider $\alpha \in \mathbb{R}^{n_v}$ and right-hand sides f and g in $\mathcal{C}(0, T; \mathbb{R}^{n_v})$ and $\mathcal{C}(0, T; \mathbb{R}^{n_p})$, respectively, and for $t \in [0, T]$ let $M(t) \in \mathbb{R}^{n_v, n_v}$ invertible, $A(t), J_2(t), J_1^\top(t)$, and $B_1(t)$ be matrices of suitable sizes with entries in $\mathcal{C}(0, T; \mathbb{R})$. For $u \in \mathcal{C}(0, T; \mathbb{R}^{n_u})$, find $v \in \mathcal{C}^1(0, T; \mathbb{R}^{n_v})$ and $p \in \mathcal{C}(0, T; \mathbb{R}^{n_p})$ that fulfill

$$(3.1a) \quad \begin{bmatrix} M(t) & 0 \\ 0 & 0 \end{bmatrix} \frac{d}{dt} \begin{bmatrix} v \\ p \end{bmatrix} (t) - \begin{bmatrix} A(t) & J_1^\top(t) \\ J_2(t) & 0 \end{bmatrix} \begin{bmatrix} v(t) \\ p(t) \end{bmatrix} - \begin{bmatrix} B_1(t) \\ 0 \end{bmatrix} u(t) = \begin{bmatrix} f(t) \\ g(t) \end{bmatrix},$$

for time $t \in (0, T]$, and

$$(3.1b) \quad v(0) = \alpha.$$

In what follows, we will always omit the time dependency of the coefficient matrices.

We will refer to the first equation in (3.1a) as the *differential part* and to the second as the *algebraic part* or *algebraic constraint*. Note that there are other so-called hidden constraints, which constrain the motion of v, \dot{v} , and p but which are not apparent in (3.1a); see [31, p. 177] and [38, pp. 323, 426] for examples.

PROPOSITION 3.2. *Consider the setup of Problem 3.1 and assume that $u \equiv 0$. If $J_2 M^{-1} J_1^\top$ is pointwise invertible, then the DAE (3.1a) is of tractability index $i_\mu = 2$.*

Proof. Having factorized the leading matrix $\begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix}$ as $\mathcal{E}_A := \begin{bmatrix} M \\ 0 \end{bmatrix}$ and $\mathcal{E}_D := \begin{bmatrix} I & 0 \end{bmatrix}$, and using that M is invertible, we find that (3.1a) has a *properly stated leading term*; cf. Definition 2.3.

Using the invertibility of $J_2 M^{-1} J_1^\top$, we can give an explicit representation of a matrix sequence as defined in Definition 2.2 for the DAE (3.1a):

$$\begin{aligned} \mathcal{E}_0 &:= \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathcal{A}_0 := \begin{bmatrix} A & J_1^\top \\ J_2 & 0 \end{bmatrix}, \\ \mathcal{Q}_0 &:= \begin{bmatrix} 0 & 0 \\ 0 & I \end{bmatrix} \text{ (projector onto } \ker \mathcal{E}_0), \\ \mathcal{P}_0 &:= I - \mathcal{Q}_0 = \begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix}, \\ \mathcal{E}_1 &:= \mathcal{E}_0 + \mathcal{A}_0 \mathcal{Q}_0 = \begin{bmatrix} M & J_1^\top \\ 0 & 0 \end{bmatrix}, \quad \mathcal{A}_1 := \mathcal{A}_0 \mathcal{P}_0 = \begin{bmatrix} A & 0 \\ J_2 & 0 \end{bmatrix}, \end{aligned}$$

$$\begin{aligned}\mathcal{Q}_1 &:= \begin{bmatrix} M^{-1}J_1^\top(J_2M^{-1}J_1^\top)^{-1}J_2 & 0 \\ -(J_2M^{-1}J_1^\top)^{-1}J_2 & 0 \end{bmatrix} =: \begin{bmatrix} \mathcal{Q} & 0 \\ -\mathcal{Q}^- & 0 \end{bmatrix} \text{ (projector onto } \ker \mathcal{E}_1), \\ \mathcal{P}_1 &:= I - \mathcal{Q}_1 =: \begin{bmatrix} \mathcal{P} & 0 \\ \mathcal{Q}^- & I \end{bmatrix}, \\ \mathcal{E}_2 &:= \mathcal{E}_1 + \mathcal{A}_1\mathcal{Q}_1 = \begin{bmatrix} M + A\mathcal{Q} & J_1^\top \\ J_2 & 0 \end{bmatrix}, \quad \mathcal{A}_2 := \mathcal{A}_1\mathcal{P}_1.\end{aligned}$$

To check for $i_\mu = 2$ (see Definition 2.3), we now have to check the dimensions of the spaces $N_0 \cap S_0$ and $N_1 \cap S_1$, as defined in Definition 2.2. In what follows, we will arbitrarily switch between a space, e.g., N_0 , and a matrix, e.g., $N_0 = \begin{bmatrix} 0 \\ I \end{bmatrix}$, if the columns of $\begin{bmatrix} 0 \\ I \end{bmatrix}$ span N_0 .

From the assumption that $J_2M^{-1}J_1^\top$ is invertible, for all $t \in (0, T)$, we have that J_2 and J_1 have full rank n_p . With $N_0 = \begin{bmatrix} 0 \\ I \end{bmatrix}$ and $S_0 = [\ker J_2]$, we have that $N_0 \cap S_0$ has dimension n_p for all $t \in (0, T)$. With $N_1 = \begin{bmatrix} \mathcal{Q} \\ -\mathcal{Q}^- \end{bmatrix}$, with $S_1 = [\ker J_2]$, and with the observation that $\mathcal{Q} \cap \ker J_2 = \{0\}$, we find that $N_1 \cap S_1 = \begin{bmatrix} 0 \\ -\mathcal{Q}^- \end{bmatrix}$, which has a rank of $n_v - n_p > 0$ for all $t \in [0, T]$.

Finally, we find that

$$\mathcal{E}_2^{-1} = \begin{bmatrix} \mathcal{P}M^{-1} & [I - \mathcal{P}M^{-1}A]M^{-1}J_1^\top S^{-1} \\ \mathcal{Q}^-M^{-1} & -[I + \mathcal{Q}^-M^{-1}AM^{-1}J_1^\top]S^{-1} \end{bmatrix}$$

is an inverse to \mathcal{E}_2 , so that $\dim(\ker \mathcal{E}_2 \cap S_2)$ is zero for all $t \in [0, T]$. \square

Remark 3.3. The particular choice of the projectors \mathcal{Q}_0 and \mathcal{Q}_1 (cf. Remark 2.4) is taken out of two considerations. First, in the corresponding decoupling (see Theorem 3.6), certain coupling terms do not occur [38, Ex. 2.34]. Second, considering flow equations, these projections admit a physical interpretation; see Remark 3.8 below. In view of numerical realizations, other choices may inhibit other desirable properties like orthogonality. However, in this paper, the projectors are used only for theoretical considerations.

In order to guarantee existence of solutions $(v, p) \in \mathcal{C}^1(0, T; \mathbb{R}^{n_v}) \times \mathcal{C}(0, T; \mathbb{R}^{n_p})$ of Problem 3.1, we make the following assumption.

Assumption 3.4. Consider Problem 3.1. We assume

- (a) that $S := J_2(t)M^{-1}(t)J_1^\top(t)$ is pointwise invertible for all $t \in [0, T]$,
- (b) sufficient regularity of the data, i.e., g , $M^{-1}J_1^\top S^{-1}$, and J_2 are differentiable on $[0, T]$, and
- (c) consistency of the given initial condition, i.e., $-J_2\alpha = g(0)$.

We remark that the taken assumptions are fulfilled for standard finite element discretizations of incompressible Navier–Stokes equations.

Remark 3.5. Let us point out first that the proposed projector chain also applies to the nonlinear Navier–Stokes equations [27, Thm. 8.6]. For spatially discretized Navier–Stokes equations [24], the matrix M is the mass matrix and typically symmetric and strictly positive definite. Also, typically, one has $J_1 = J_2$. Thus, if one considers stable finite element discretization schemes, i.e., schemes that fulfill the so-called LBB condition, $J := J_1 = J_2$ has full rank, so that Assumption 3.4(a) is fulfilled. Also, typically, $g = 0$ and the coefficient matrices M , J_1 , and J_2 are independent of time such that also Assumption 3.4(b) is fulfilled. Part (c) of the assumption is fulfilled if the initial velocity is consistent, which is, typically, that it is discretely divergence-free.

Theorem 3.6 gives a decoupling of (3.1) into the *inherent ODE* and algebraic equations that will be used to ensure existence and uniqueness of solutions.

THEOREM 3.6. *Consider Problem 3.1 and define $\mathcal{P} := I - \mathcal{Q}$, $\mathcal{Q} := M^{-1}J_1^\top S^{-1}J_2$, and $\mathcal{Q}^- := S^{-1}J_2$. Let $u \in \mathcal{C}(0, T; \mathbb{R}^{n_u})$ be arbitrary.*

1. *If Assumption 3.4(a)–(b) holds, then any solution (v, p) to (3.1) can be decomposed as $(v_{\mathcal{P}} + \mathcal{Q}v, p)$ such that*

$$(3.2a) \quad \mathcal{Q}v = -M^{-1}J_1^\top S^{-1}g,$$

$$(3.2b) \quad p = -\mathcal{Q}^-[M^{-1}[A\mathcal{Q}v + v_{\mathcal{P}}] + B_1u + f] + \frac{d}{dt}(\mathcal{Q}v),$$

and $v_{\mathcal{P}} := \mathcal{P}v$ solves the ODE

$$(3.2c) \quad v_{\mathcal{P}} - \left[\frac{d}{dt}\mathcal{P} + \mathcal{P}M^{-1}A \right] [\mathcal{Q}v + v_{\mathcal{P}}] - \mathcal{P}M^{-1}[B_1u + f] = 0, \quad v_{\mathcal{P}}(0) = \mathcal{P}\alpha.$$

2. *If, in addition, Assumption 3.4(c) holds, then Problem 3.1 has a unique solution (v, p) given via $(v, p) = (v_{\mathcal{P}} + \mathcal{Q}v, p)$, where $(\mathcal{Q}v, p, v_{\mathcal{P}})$ is uniquely defined through (3.2).*

Proof. We consider the state equations (3.1a) given as

$$(3.3) \quad \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \frac{d}{dt} \begin{bmatrix} v \\ p \end{bmatrix} - \begin{bmatrix} A & J_1^\top \\ J_2 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} = \begin{bmatrix} B_1u + f \\ 0 \end{bmatrix}$$

and we compute the sequence of operators from Definition 2.2 as given in the proof of Proposition 3.2. In particular, with the projectors $\mathcal{Q} = M^{-1}J_1^\top S^{-1}J_2$, which satisfy

$$\mathcal{Q}^2 = \mathcal{Q}, \quad J_2\mathcal{Q} = J_2, \quad \mathcal{Q}M^{-1}J_1^\top = M^{-1}J_1^\top, \quad \text{and} \quad \mathcal{Q}^-\mathcal{Q} = \mathcal{Q}^-,$$

and $\mathcal{P} = I - \mathcal{Q}$, we have

$$(3.4) \quad \mathcal{E}_2^{-1} = \begin{bmatrix} \mathcal{P}M^{-1} & [I - \mathcal{P}M^{-1}A]M^{-1}J_1^\top S^{-1} \\ \mathcal{Q}^-M^{-1} & -[I + \mathcal{Q}^-M^{-1}AM^{-1}J_1^\top]S^{-1} \end{bmatrix}.$$

Scaling the state equations (3.3) by \mathcal{E}_2^{-1} we get

$$(3.5) \quad \begin{bmatrix} \mathcal{P} & 0 \\ \mathcal{Q}^- & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix} - \left[\begin{bmatrix} \mathcal{P}M^{-1}A\mathcal{P} & 0 \\ \mathcal{Q}^-M^{-1}A\mathcal{P} & 0 \end{bmatrix} + \begin{bmatrix} \mathcal{Q} & 0 \\ -\mathcal{Q}^- & I \end{bmatrix} \right] \begin{bmatrix} v \\ p \end{bmatrix} = \mathcal{E}_2^{-1} \begin{bmatrix} M^{-1}[B_1u + f] \\ g \end{bmatrix}.$$

Having applied the projectors $\mathcal{Q}_1, \mathcal{Q}_0\mathcal{P}_1$, and $\mathcal{P}_0\mathcal{P}_1$ (cf. Definition 2.2) to (3.5), we obtain the three subsystems

$$(3.6a) \quad \begin{aligned} - \begin{bmatrix} \mathcal{Q} & 0 \\ -\mathcal{Q}^- & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} &= \mathcal{Q}_1\mathcal{E}_2^{-1} \begin{bmatrix} B_1u + f \\ g \end{bmatrix} \\ &= \begin{bmatrix} M^{-1}J_1^\top S^{-1}g \\ -S^{-1}g \end{bmatrix}, \end{aligned}$$

$$(3.6b) \quad \begin{aligned} \begin{bmatrix} 0 & 0 \\ \mathcal{Q}^- & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix} - \begin{bmatrix} 0 & 0 \\ \mathcal{Q}^-M^{-1}A\mathcal{P} & I \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} &= \mathcal{Q}_0\mathcal{P}_1\mathcal{E}_2^{-1} \begin{bmatrix} B_1u + f \\ g \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ \mathcal{Q}^-M^{-1}[B_1u + f - AM^{-1}J_1^\top S^{-1}g] \end{bmatrix}, \end{aligned}$$

and

$$(3.6c) \quad \begin{aligned} \begin{bmatrix} \mathcal{P} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix} - \begin{bmatrix} \mathcal{P}M^{-1}A\mathcal{P} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} &= \mathcal{P}_0\mathcal{P}_1\mathcal{E}_2^{-1} \begin{bmatrix} B_1u + f \\ g \end{bmatrix} \\ &= \begin{bmatrix} \mathcal{P}M^{-1}[B_1u + f - AM^{-1}J_1^T S^{-1}g] \\ 0 \end{bmatrix}, \end{aligned}$$

respectively. Since $\mathcal{Q}_1 + \mathcal{P}_0\mathcal{P}_1 + \mathcal{Q}_0\mathcal{P}_1 = I$, equations (3.6) form a decomposition of (3.5). We decompose $v = v_{\mathcal{P}} + \mathcal{Q}v$, where $v_{\mathcal{P}} := \mathcal{P}v$, so that from (3.6a) we can deduce that

$$(3.7) \quad \mathcal{Q}v = -M^{-1}J_1^T S^{-1}g$$

and that $\mathcal{Q}v$ is differentiable by assumption. With $\dot{v} = \frac{d}{dt}(\mathcal{Q}v) + \dot{v}_{\mathcal{P}}$ and $\mathcal{Q}^{-}v_{\mathcal{P}} = 0$, (3.6b) gives

$$(3.8) \quad p = -\mathcal{Q}^{-}M^{-1}[A(\mathcal{Q}v + v_{\mathcal{P}})[\mathcal{Q}v + v_{\mathcal{P}}] + B_1u + f] + \mathcal{Q}^{-}\frac{d}{dt}(\mathcal{Q}v),$$

while $v_{\mathcal{P}} := \mathcal{P}v$ satisfies (3.6c), which is the *inherent* or *underlying* ODE:

$$(3.9) \quad \dot{v}_{\mathcal{P}} - \left[\frac{d}{dt}(\mathcal{P}) + \mathcal{P}M^{-1}A(\mathcal{Q}v + v_{\mathcal{P}}) \right] [\mathcal{Q}v + v_{\mathcal{P}}] = \mathcal{P}M^{-1}[B_1u + f], \quad v_{\mathcal{P}}(0) = \mathcal{P}\alpha.$$

The other way round, equation system (3.2) uniquely defines $(\mathcal{Q}v, p, v_{\mathcal{P}})$. Then, by construction, $(v_{\mathcal{P}} + \mathcal{Q}v, p)$ solves (3.3) and since by assumption the initial value is consistent, it also solves (3.1). \square

Remark 3.7. Note the necessity of the consistency condition given in Assumption 3.4(c), since by (3.7) the condition

$$J_2\alpha = J_2v(0) = J_2[\mathcal{Q}v(0) + \mathcal{P}v(0)] = J_2\mathcal{Q}v(0) = -g(0)$$

must hold, and note that an initial condition for p would have to fulfill (3.8) at $t = 0$ and an initial condition for \dot{v} would have to satisfy $J_2\dot{v}(0) = \dot{g}(0)$.

Remark 3.8. In the setting of the Navier–Stokes equation, the projector \mathcal{Q} realizes the discrete *Helmholtz-decomposition* that splits a vector field into a divergence-free part and a part that can be expressed as the gradient of a scalar potential; cf. [23, Cor. 3.4]. If J_2 is the discrete divergence operator, then the decomposition $v = \mathcal{Q}v + \mathcal{P}v =: \mathcal{Q}v + v_{\mathcal{P}}$ delivers that $J_2v_{\mathcal{P}} = 0$ and $\mathcal{Q}v$ is in the range of $M^{-1}J_1^T$, which is the discrete gradient operator in many discretization schemes. The matrix \mathcal{Q}^{-} is a generalized left inverse of $M^{-1}J_1^T$ and can be seen as the operator that maps the potential field $\mathcal{Q}v = M^{-1}J_1^T\rho$ onto its potential ρ . Accordingly, (3.2b) is equivalent to the discrete *Pressure Poisson equation* of a linearized Navier–Stokes equation; see [24, Chap. 3.16.1].

4. Linear-quadratic optimal control. We now define the linear-quadratic optimal control problem for the state equations introduced in section 3. We state the *formal* optimality conditions and show that they have a similar structure as the state equations. We reflect the decoupling of the state equations in the adjoint equations to decouple the formal optimality conditions while preserving their symmetry and to state existence of solutions of the optimality system. Once a solution is known to exist, we show that it can be obtained via a differential-algebraic Riccati equation that does not resort to projectors or variable transformations.

Problem 4.1. Consider Problem 3.1 and let Assumption 3.4 hold. For a given $v^* \in \mathcal{C}(0, T; \mathbb{R}^{n_v})$ and for given $V \in \mathbb{R}^{n_v, n_v}$ and $W \in \mathcal{C}(0, T; \mathbb{R}^{n_v, n_v})$ (pointwise) symmetric positive definite and $R \in \mathcal{C}(0, T; \mathbb{R}^{n_u, n_u})$ symmetric strictly positive definite, find $u \in \mathcal{C}(0, T; \mathbb{R}^{n_u})$ such that

$$(4.1) \quad \mathcal{J}(v, u) = \frac{1}{2}[v - v^*]^\top V [v - v^*] \Big|_{t=T} + \frac{1}{2} \int_0^T [v - v^*]^\top W [v - v^*] + u^\top R u \, dt,$$

becomes minimal, where v and u are subject to the state equations (3.1).

For Problem 4.1, the *formal* optimality conditions (cf. [32, 33]) are given by

$$(4.2a) \quad \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \frac{d}{dt} \begin{bmatrix} v \\ p \end{bmatrix} - \begin{bmatrix} A & J_1^\top \\ J_2 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} - \begin{bmatrix} B_1 u \\ 0 \end{bmatrix} = \begin{bmatrix} f \\ g \end{bmatrix},$$

$$(4.2b) \quad v(0) = \alpha,$$

$$(4.2c) \quad -\frac{d}{dt} \begin{bmatrix} M^\top \lambda \\ \mu \end{bmatrix} - \begin{bmatrix} A^\top & J_2^\top \\ J_1 & 0 \end{bmatrix} \begin{bmatrix} \lambda \\ \mu \end{bmatrix} + \begin{bmatrix} W v \\ 0 \end{bmatrix} = \begin{bmatrix} W v^* \\ 0 \end{bmatrix},$$

$$(4.2d) \quad M^\top \lambda(T) + V v(T) = V v^*(T),$$

$$(4.2e) \quad -B_1^\top \lambda + R u = 0.$$

Remark 4.2. The formal optimality system (4.2) is commonly referred to as *Euler–Lagrange equations*. In the considered DAE case, if it possesses a solution, then it provides necessary optimality conditions; cf. [5, 33, 39]. Recall, however, that the optimal control problem can be solvable also if (4.2) is not well posed [33]. Thus, the important intermediate step is to establish solvability of the *Euler–Lagrange equations* (4.2).

Remark 4.3. We will show that in the considered setup of Problem 4.1, the possible failure of the conditions is related only to the consistency condition that is imposed by (4.2d). For more general control setups that include the variable p in the cost functional or that comprise control action in the algebraic constraint, the structure of the formal optimality system can be fundamentally different; cf. [27, Chap. 8.6] and the concluding remarks at the end of this paper.

We use the invertibility of R to express u via

$$(4.3) \quad u = R^{-1} B_1^\top \lambda$$

and write (4.2) as

$$(4.4a) \quad \begin{bmatrix} M \dot{v} \\ 0 \\ -\frac{d}{dt}(M^\top \lambda) \\ 0 \end{bmatrix} - \begin{bmatrix} B_1 R^{-1} B_1^\top & 0 & A & J_1^\top \\ 0 & 0 & J_2 & 0 \\ A^\top & J_2^\top & -W & 0 \\ J_1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \lambda \\ \mu \\ v \\ p \end{bmatrix} = \begin{bmatrix} f \\ g \\ W v^* \\ 0 \end{bmatrix},$$

$$(4.4b) \quad M^\top \lambda(T) + V v(T) = V v^*(T), \quad \text{and} \quad v(0) = \alpha.$$

For completeness we state that for the particular choice of the cost functional the index of the *Euler–Lagrange equations* is the same as of the state equations (3.1) with zero inputs.

Remark 4.4. By inverting the mass matrices and permuting the rows and the columns, system (4.4) can be brought into the form of (3.1) with initial and terminal

conditions. Then, with Assumption 3.4(a), we have that system (4.4) is of tractability index 2. This follows from Proposition 3.2 and from

$$\begin{bmatrix} 0 & J_2 \\ J_1 & 0 \end{bmatrix} \begin{bmatrix} 0 & M \\ -M^T & 0 \end{bmatrix}^{-1} \begin{bmatrix} 0 & J_1^T \\ J_2^T & 0 \end{bmatrix} = \begin{bmatrix} 0 & J_2 M^{-1} J_1^T \\ -J_1 M^{-T} J_2^T & 0 \end{bmatrix}$$

being invertible by Assumption 3.4(a).

Remark 4.2 states that system (4.2) functions as a necessary optimality condition only if the existence of a solution is guaranteed. Recall that the considered equations are in the form of (3.1); cf. Remark 4.4. Therefore, one may apply Theorem 3.6 to identify the inherent ODE (3.9), so that one can use the theory for ODEs to state the existence of solutions to the obtained linear boundary value problem; cf. [4, Thm. 3.26]. However, the reformulation as used in Theorem 3.6 will not preserve the structure of (4.4) and, thus, makes it more difficult to investigate whether the boundary values admit the existence of a solution. We will rather use a reformulation that preserves the *Hamiltonian* structure such that the existence of a solution can be obtained via a differential-algebraic Riccati equation; cf. [17, 18, 19, 42] for the ODE case. Additionally, we show that the solution can also be obtained via a Riccati equation formulated directly for the formal optimality conditions.

LEMMA 4.5. *Consider Problem 4.1 with Assumption 3.4(a)–(b) holding.*

- (1) *Let (v, p, λ, μ) be a solution of the associated Euler–Lagrange equations as given by (4.2). Then (v, p, λ, μ) can be decomposed as $(v, p) = (v_{\mathcal{P}} + \mathcal{Q}v, p)$ and $(M^T \lambda, \mu) = (\lambda_{\mathcal{P}} + \mathcal{Q}^T M^T \lambda, \mu)$ such that*

$$(4.5a) \quad \mathcal{Q}v = -M^{-1} J_1^T S^{-1} g,$$

$$(4.5b) \quad \mathcal{Q}^T M^T \lambda = 0,$$

$$(4.5c) \quad \mu = -S^{-T} J_1 M^{-T} [A^T M^{-T} \lambda_{\mathcal{P}} - W[\mathcal{Q}v + v_{\mathcal{P}}]],$$

$$(4.5d) \quad \begin{aligned} p &= -\mathcal{Q}^- [M^{-1} [A[\mathcal{Q}v + v_{\mathcal{P}}] + f] - \frac{d}{dt}(\mathcal{Q}v)] \\ &\quad - \mathcal{Q}^- M^{-1} G M^{-T} [\lambda_{\mathcal{P}} + \mathcal{Q}^T A M^{-T} \lambda_{\mathcal{P}} + J_2^T \mu], \end{aligned}$$

and

$$(4.5e) \quad \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix} \begin{bmatrix} \dot{\lambda}_{\mathcal{P}} \\ \dot{v}_{\mathcal{P}} \end{bmatrix} - \begin{bmatrix} G_0 & A_0 \\ A_0^T & -W_0 \end{bmatrix} \begin{bmatrix} \lambda_{\mathcal{P}} \\ v_{\mathcal{P}} \end{bmatrix} = \begin{bmatrix} \mathcal{P} M^{-1} [f - A M^{-1} J_1^T S^{-1} g] \\ \mathcal{P}^T - W M^{-1} J_1^T S^{-1} g \end{bmatrix},$$

$$v_{\mathcal{P}}(0) = \mathcal{P}\alpha \quad \text{and} \quad \lambda_{\mathcal{P}}(T) = -\mathcal{P}^T V [v(T) - v^*(t)],$$

where $A_0 := \frac{d}{dt} \mathcal{P} + \mathcal{P} M^{-1} A \mathcal{P}$, $G_0 = G_0^T = \mathcal{P} M^{-1} B_1 R^{-1} B_1^T M^{-T} \mathcal{P}^T$, $W_0 = W_0^T := \mathcal{P}^T W \mathcal{P}$ and where $\mathcal{P} := I - \mathcal{Q}$, $\mathcal{Q} := M^{-1} J_1^T S^{-1} J_2$, $\mathcal{Q}^- := S^{-1} J_2$, and $S := J_2 M^{-1} J_1^T$ as defined in Theorem 3.6.

- (2) *If in addition*

$$(4.6) \quad -J_2 \alpha = g(0) \quad \text{and} \quad J_1 M^{-T} V = 0,$$

then the Euler–Lagrange equations (4.2) possess a unique solution.

- (3) *If in addition $f, g,$ and v^* are zero, then (4.2) can be decoupled via*

$$(4.7) \quad \begin{bmatrix} \lambda \\ \mu \end{bmatrix} = \begin{bmatrix} X_1 & X_2 \\ X_2 & 0 \end{bmatrix} \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix},$$

where $X_1 = X_1^\top$ and X_2 fulfill the differential-algebraic Riccati equation

$$(4.8a) \quad \frac{d}{dt}(M^\top X_1 M) + M^\top X_1 A + A^\top X_1 M + M^\top X_1 B_1 R^{-1} B_1^\top X_1 M - W + M^\top X_2^\top J_2 + J_2^\top X_2 M = 0$$

with the terminal condition

$$(4.8b) \quad M^\top X_1(T)M = -V$$

and the algebraic constraints

$$(4.8c) \quad M^\top J_1 X_1 = 0 \quad \text{and} \quad J_1 X_1 M = 0.$$

Equations (4.8) uniquely define a symmetric negative semidefinite X_1 .

- (4) If, however, f , g , and v^* are not identically zero, then the solution of (4.2) decouples via

$$(4.9) \quad \begin{bmatrix} \lambda \\ \mu \end{bmatrix} = \begin{bmatrix} X_1 & X_2^\top \\ X_2 & 0 \end{bmatrix} \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} + \begin{bmatrix} w_1 \\ w_2 \end{bmatrix},$$

where X_1 and X_2 are given by a solution of (4.8) and (w_1, w_2) solve

$$-\frac{d}{dt} \begin{bmatrix} M^\top w_1 \\ 0 \end{bmatrix} - \begin{bmatrix} M^\top X_1 G + A^\top & J_2^\top \\ J_1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} W v^* + M^\top [X_1 f + X_2 g] \\ 0 \end{bmatrix},$$

$$M^\top w_1(T) = V v^*(T),$$

where $G := B_1 R^{-1} B_1^\top$. The solution (w_1, w_2) exists and w_1 is unique.

Proof. ad (1) We define $G := B_1 R^{-1} B_1^\top$ and write the Euler–Lagrange system (4.2) as

$$\begin{bmatrix} 0 & 0 & I & 0 \\ 0 & 0 & 0 & 0 \\ -I & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \frac{d}{dt} \begin{bmatrix} M^\top \lambda \\ \mu \\ v \\ p \end{bmatrix} - \begin{bmatrix} M^{-1} G M^{-\top} & 0 & M^{-1} A & M^{-1} J_1^\top \\ 0 & 0 & J_2 & 0 \\ A^\top M^{-\top} & J_2^\top & -W & 0 \\ J_1 M^{-\top} & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} M^\top \lambda \\ \mu \\ v \\ p \end{bmatrix} = \begin{bmatrix} f^\top M^{-\top} & g^\top & v^{*\top} W & 0 \end{bmatrix}^\top,$$

$$v(0) = \alpha \quad \text{and} \quad M^\top \lambda(T) = -V[v(T) - v^*(T)].$$

In order to preserve the self-adjoint structure (cf. [35]), only congruence transformations should be applied, i.e., a scaling of the equations must be accompanied by the transpose inverse scaling of the variables. In accordance to (3.5) we congruently transform the system by

$$S_2 := \begin{bmatrix} \mathcal{E}_2^{-1} & 0 & 0 \\ 0 & I & 0 \\ 0 & 0 & I \end{bmatrix} = \begin{bmatrix} \mathcal{P} & [I - \mathcal{P} M^{-1} A] M^{-1} J_1^\top S^{-1} & 0 & 0 \\ \mathcal{Q}^- & -[I + \mathcal{Q}^- M^{-1} A M^{-1} J_1^\top] S^{-1} & 0 & 0 \\ 0 & 0 & I & 0 \\ 0 & 0 & 0 & I \end{bmatrix},$$

where $\mathcal{E}_2 = \begin{bmatrix} I + M^{-1} A \mathcal{Q} & M^{-1} J_1^\top \\ J_2 & 0 \end{bmatrix}$ as defined in Definition 2.2 with the inverse as given in (3.4). The summand that comes from the time dependency in the variable trans-

formation S_2^\top is given by

$$S_2 \begin{bmatrix} 0 & 0 & I & 0 \\ 0 & 0 & 0 & 0 \\ -I & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \dot{S}_2^\top = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -\frac{d}{dt}\mathcal{P}^\top & -\dot{\mathcal{Q}}^\top & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

With this we get the scaled and transformed system

$$\begin{aligned} \tilde{f} = & \begin{bmatrix} 0 & 0 & \mathcal{P} & 0 \\ 0 & 0 & \mathcal{Q}^- & 0 \\ -\mathcal{P}^\top & -\dot{\mathcal{Q}}^\top & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{\tilde{\lambda}} \\ \dot{\tilde{\mu}} \\ \dot{v} \\ \dot{p} \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -\frac{d}{dt}\mathcal{P}^\top & -\dot{\mathcal{Q}}^\top & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \tilde{\lambda} \\ \tilde{\mu} \\ v \\ p \end{bmatrix} \\ & - \begin{bmatrix} \mathcal{P}M^{-1}GM^{-\top}\mathcal{P}^\top & M^{-1}GM^{-\top}\mathcal{Q}^\top & \mathcal{P}M^{-1}A\mathcal{P} + \mathcal{Q} & 0 \\ \mathcal{Q}^-M^{-1}GM^{-\top} & 0 & \mathcal{Q}^-M^{-1}A\mathcal{P} - \mathcal{Q}^- & I \\ \mathcal{P}^\top A^\top M^{-\top}\mathcal{P}^\top + \mathcal{Q}^\top & \mathcal{P}^\top A^\top M^{-\top}\mathcal{Q}^\top - \mathcal{Q}^\top & -W & 0 \\ 0 & I & 0 & 0 \end{bmatrix} \begin{bmatrix} \tilde{\lambda} \\ \tilde{\mu} \\ v \\ p \end{bmatrix} \end{aligned} \quad (4.10)$$

with the transformed state and scaled right-hand side

$$\begin{bmatrix} \tilde{\lambda} \\ \tilde{\mu} \\ v \\ p \end{bmatrix} := S_2^{-\top} \begin{bmatrix} M^\top \lambda \\ \mu \\ v \\ p \end{bmatrix} = \begin{bmatrix} [I + \mathcal{Q}^\top A^\top M^{-\top}]M^\top \lambda + J_2^\top \mu \\ J_1 \lambda \\ v \\ p \end{bmatrix} \quad (4.11)$$

and $\tilde{f} := S_2 [f^\top M^{-\top} \quad g^\top \quad v^{*\top} W \quad 0]^\top$, respectively. From the last line in (4.10) we find that $\tilde{\mu} = 0$ so that we can rewrite the equations for $(\tilde{\lambda}, v, p)$ as

$$\begin{aligned} \begin{bmatrix} \mathcal{P} & 0 \\ \mathcal{Q}^- & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix} - \begin{bmatrix} \mathcal{P}M^{-1}GM^{-\top}\mathcal{P}^\top & 0 \\ \mathcal{Q}^-M^{-1}GM^{-\top} & 0 \end{bmatrix} \begin{bmatrix} \tilde{\lambda} \\ \tilde{\mu} \end{bmatrix} - \begin{bmatrix} \mathcal{P}M^{-1}A\mathcal{P} + \mathcal{Q} & 0 \\ \mathcal{Q}^-M^{-1}A\mathcal{P} - \mathcal{Q}^- & I \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} \\ = \mathcal{E}_2^{-1} \begin{bmatrix} M^{-1}f \\ g \end{bmatrix} \end{aligned} \quad (4.12a)$$

and

$$-\frac{d}{dt}(\mathcal{P}^\top \tilde{\lambda}) - [\mathcal{P}^\top A^\top M^{-\top}\mathcal{P}^\top + \mathcal{Q}^\top] \tilde{\lambda} + Wv = Wv^*. \quad (4.12b)$$

We apply the projectors

$$\mathcal{Q}_1 = \begin{bmatrix} \mathcal{Q} & 0 \\ -\mathcal{Q}^- & 0 \end{bmatrix}, \quad \mathcal{Q}_0 \mathcal{P}_1 = \begin{bmatrix} 0 & 0 \\ \mathcal{Q}^- & I \end{bmatrix} \quad \text{and} \quad \mathcal{P}_0 \mathcal{P}_1 = \begin{bmatrix} \mathcal{P} & 0 \\ 0 & 0 \end{bmatrix}$$

to (4.12a)(cf. the treatment of (3.6)) to obtain the three subsystems

$$-\begin{bmatrix} \mathcal{Q} & 0 \\ -\mathcal{Q}^- & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} = \begin{bmatrix} M^{-1}J_1^\top S^{-1}g \\ -S^{-1}g \end{bmatrix}, \quad (4.13a)$$

$$\begin{aligned} \begin{bmatrix} 0 & 0 \\ \mathcal{Q}^- & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix} - \begin{bmatrix} 0 & 0 \\ \mathcal{Q}^-M^{-1}GM^{-\top} & 0 \end{bmatrix} \begin{bmatrix} \tilde{\lambda} \\ \tilde{\mu} \end{bmatrix} \\ - \begin{bmatrix} 0 & 0 \\ \mathcal{Q}^-M^{-1}A\mathcal{P} & I \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} = \begin{bmatrix} 0 \\ \mathcal{Q}^-M^{-1}[f - AM^{-1}J_1^\top S^{-1}g] \end{bmatrix}, \end{aligned} \quad (4.13b)$$

and

$$(4.13c) \quad \begin{bmatrix} \mathcal{P} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix} - \begin{bmatrix} \mathcal{P}M^{-1}GM^{-T}\mathcal{P}^T & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \tilde{\lambda} \\ \tilde{\mu} \end{bmatrix} - \begin{bmatrix} \mathcal{P}M^{-1}A\mathcal{P} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} = \begin{bmatrix} \mathcal{P}M^{-1}[f - AM^{-1}J_1^T S^{-1}g] \\ 0 \end{bmatrix},$$

respectively. Using the projector property $\mathcal{P}^T = \mathcal{P}^T\mathcal{P}^T$ to obtain the relation

$$\frac{d}{dt}(\mathcal{P}^T\tilde{\lambda}) = \mathcal{P}^T\frac{d}{dt}(\mathcal{P}^T\tilde{\lambda}) + \dot{\mathcal{P}}^T\mathcal{P}^T\tilde{\lambda} = \frac{d}{dt}(\mathcal{P}^T\tilde{\lambda}) - \mathcal{Q}^T\frac{d}{dt}(\mathcal{P}^T\tilde{\lambda}) + \dot{\mathcal{P}}^T\mathcal{P}^T\tilde{\lambda},$$

we split (4.12b) into the two subsystems

$$(4.14a) \quad \mathcal{Q}^T\frac{d}{dt}(\mathcal{P}^T\tilde{\lambda}) - \mathcal{Q}^T\tilde{\lambda} + \mathcal{Q}^TWv = \mathcal{Q}^TWv^*$$

and

$$(4.14b) \quad -\frac{d}{dt}(\mathcal{P}^T\tilde{\lambda}) - \frac{d}{dt}(\mathcal{P}^T)\mathcal{P}^T\tilde{\lambda} - \mathcal{P}^TA^TM^{-T}\mathcal{P}^T\tilde{\lambda}_1 + \mathcal{P}^TWv = \mathcal{P}^TWv^*.$$

If we then define $v_{\mathcal{P}} := \mathcal{P}v$ and $\tilde{\lambda}_{\mathcal{P}} := \mathcal{P}^T\tilde{\lambda}$ and decompose $\tilde{\lambda} = \tilde{\lambda}_{\mathcal{P}} + \mathcal{Q}^T\tilde{\lambda}$ and $v = v_{\mathcal{P}} + \mathcal{Q}v$ we find that (4.13a)–(4.13b) and (4.14a) define algebraic relations for

$$(4.15a) \quad \mathcal{Q}v = -M^{-1}J_1^T S^{-1}g,$$

and, with $\mathcal{Q}^{-}\dot{v}_{\mathcal{P}} = 0$ and $\mathcal{Q}^T\dot{\tilde{\lambda}}_{\mathcal{P}} = 0$,

$$(4.15b) \quad \mathcal{Q}^T\tilde{\lambda} = \mathcal{Q}^TW[\mathcal{Q}v + v_{\mathcal{P}} - v^*],$$

$$(4.15c) \quad p = -\mathcal{Q}^-[M^{-1}A[\mathcal{Q}v + v_{\mathcal{P}}] + M^{-1}f + M^{-1}GM^{-T}\tilde{\lambda} - \frac{d}{dt}(\mathcal{Q}v)],$$

while $\tilde{\lambda}_{\mathcal{P}}$ and $v_{\mathcal{P}}$ solve the coupled ODEs given by (4.14b) and (4.13c):

$$(4.16a) \quad -\dot{\tilde{\lambda}}_{\mathcal{P}} - \left[\frac{d}{dt}\mathcal{P}^T + \mathcal{P}^TA^TM^{-T}\mathcal{P}^T\right]\tilde{\lambda}_{\mathcal{P}} + \mathcal{P}^TW\mathcal{P}v_{\mathcal{P}} = \mathcal{P}^TWM^{-1}J_1^T S^{-1}g + \mathcal{P}^TWv^*$$

and

$$(4.16b) \quad \dot{v}_{\mathcal{P}} - \mathcal{P}M^{-1}GM^{-T}\mathcal{P}^T\tilde{\lambda}_{\mathcal{P}} - \left[\frac{d}{dt}\mathcal{P} + \mathcal{P}M^{-1}A\mathcal{P}\right]v_{\mathcal{P}} = \mathcal{P}M^{-1}[f - AM^{-1}J_1^T S^{-1}g],$$

with the projected initial and terminal conditions

$$(4.16c) \quad v_{\mathcal{P}}(0) = \mathcal{P}\alpha \quad \text{and} \quad \lambda_{\mathcal{P}}(T) = \mathcal{P}^TM^T\lambda(T) = -\mathcal{P}^TV[v(T) - v^*(T)].$$

Note that we have used the projector property $\mathcal{P} = \mathcal{P}^2$ to make the symmetry in (4.16) obvious.

In view of expressing the obtained relations in terms of the original variables (λ, μ) we use the relations defined in (4.11) and observe that

$$\tilde{\lambda}_{\mathcal{P}} = \mathcal{P}^T\tilde{\lambda} = \mathcal{P}^T[M^T\lambda + \mathcal{Q}^TA^T\lambda + J_2^T\mu] = \mathcal{P}^TM^T\lambda =: \lambda_{\mathcal{P}}.$$

For μ we use $\mathcal{Q}^T \tilde{\lambda} = \mathcal{Q}^T [I + \mathcal{Q}^T A M^{-T}] M^T \lambda + \mathcal{Q}^T J_2^T \mu = \mathcal{Q}^T A M^{-T} \lambda_{\mathcal{P}} + J_2^T \mu$, relation (4.15b), and the invertibility of $S^T = J_1 M^{-1} J_2^T$ to get

$$\mu = S^{-T} J_1 M^{-T} [W[\mathcal{Q}v + v_{\mathcal{P}} - v^*] + A M^{-T} \lambda_{\mathcal{P}}].$$

Similarly, one can express the equation for p in terms of (λ, μ) , which completes the derivation of equations (4.5).

ad (2) We proceed as follows. We show that for any α and $\mathcal{P}^T V$ symmetric positive semidefinite the decoupled system (4.5) has a unique solution $(v_{\mathcal{P}}, \mathcal{Q}v, p, \lambda_{\mathcal{P}}, \mathcal{Q}^T M^T \lambda, \mu)$. Then, we confer that under the consistency conditions (4.6), the solution of (4.5) provides a solution of the Euler–Lagrange equations (4.2). Finally, by (1) every solution of (4.2) has a representation in (4.5), such that—in summary—the Euler–Lagrange equations must possess a unique solution.

In (4.5e), consider the case with a zero right-hand side. With the ansatz $\lambda_{\mathcal{P}} = X_0(t)v_{\mathcal{P}}(t)$, we derive the differential Riccati equation

$$(4.17) \quad \dot{X}_0 = -X_0 G_0 X_0 - X_0 A_0 - A_0^T X_0 + W_0, \quad X_0(T) = -\mathcal{P}^T V,$$

which has a unique solution (cf. [1, Thm. 4.1.6]), since $\mathcal{P}^T V$, G_0 , and W_0 are symmetric positive semidefinite. With this X_0 we get $v_{\mathcal{P}}$, and $\lambda_{\mathcal{P}}$ as the solution of $\dot{v}_{\mathcal{P}} - [G_0 X_0 + A_0]v_{\mathcal{P}} = 0$, $v_{\mathcal{P}}(0) = \mathcal{P}\alpha$, and $\lambda_{\mathcal{P}} = X_0 v_{\mathcal{P}}$, respectively.

One can show that if there exists a solution to (4.5e) with a zero right-hand side, then it is unique. This is equivalent to the fact that the linear part of the affine boundary conditions are stated such that (4.5e) with $\mathcal{P}^T V$ symmetric positive semidefinite has a unique solution (cf. [4, Thm. 3.26]) for any continuous right-hand side.

A solution of (4.2) uniquely defines a solution to (4.5). The converse is true if and only if the given initial and terminal conditions are consistent, namely,

$$(4.18a) \quad \mathcal{Q}v(0) = \mathcal{Q}\alpha = M^{-1} J_1^T S^{-1} J_2 \alpha$$

and

$$(4.18b) \quad \mathcal{Q}^T M^T \lambda(T) = -\mathcal{Q}^T V[v(T) - v^*(T)] = -J_2^T S^{-T} J_1 M^{-T} V[v(T) - v^*(T)];$$

cf. Theorem 3.6. By (4.5a) we have that $\mathcal{Q}v(0) = -M^{-1} J_1^T S^{-1} g(0)$ such that $-J_2 \alpha = g(0)$ is necessary and sufficient for (4.18a). By (4.5b) we have that $\mathcal{Q}^T M^T \lambda = 0$ such that $J_1 M^{-T} V = 0$ is sufficient for (4.18b). Note that in this case we can infer that $J_1^T M^{-T} V = 0$, so that $V M^{-1} J_1 = 0$ or $V \mathcal{Q} = 0$, which means that $Vv = V \mathcal{P}v$ and that, in (4.5f), $\mathcal{P}^T V$ can be replaced by $\mathcal{P}^T V \mathcal{P}$. Thus, condition (4.6) implies the symmetry in the terminal condition that was sufficient for the existence of X_0 in (4.17).

ad (3) With the ansatz

$$(4.19) \quad \begin{bmatrix} \lambda \\ \mu \end{bmatrix} = \begin{bmatrix} X_1 & X_2^T \\ X_2 & 0 \end{bmatrix} \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix}$$

we obtain that

$$(4.20) \quad \frac{d}{dt} \left(\begin{bmatrix} M^T & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \lambda \\ \mu \end{bmatrix} \right) = \begin{bmatrix} \frac{d}{dt} M^T X_1 M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} + \begin{bmatrix} M^T X_1 & M^T X_2^T \\ 0 & 0 \end{bmatrix} \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix}.$$

In (4.20) we eliminate $\frac{d}{dt} \begin{bmatrix} M^\top & 0 \\ 0 & \lambda \end{bmatrix}$ and $\begin{bmatrix} M & 0 \\ 0 & \dot{p} \end{bmatrix}$ via the relations given in the state and the adjoint equations in (4.2) and we replace every occurrence of $\begin{bmatrix} \lambda \\ \mu \end{bmatrix}$ by the ansatz (4.19). Thus, we obtain that $\mathcal{X} \begin{bmatrix} v \\ p \end{bmatrix} = 0$, where

$$(4.21) \quad \mathcal{X} := \begin{bmatrix} \frac{d}{dt}(M^\top X_1 M) + A^\top X_1 M + M^\top X_1 A + M^\top X_1 G X_1 M - W + J_2^\top X_2 M + M^\top X_2^\top J_2 & M^\top X_1 J_1^\top \\ J_1 X_1 M & 0 \end{bmatrix}.$$

Since $\mathcal{X} \begin{bmatrix} v \\ p \end{bmatrix} = 0$ must hold for every state trajectory, one requires $\mathcal{X} = 0$, which gives the equations for X_1 and X_2 :

$$(4.22a) \quad \frac{d}{dt}(M^\top X_1 M) + M^\top X_1 A + A^\top X_1 M + M^\top X_1 G X_1 M - W + M^\top X_2^\top J_2 + J_2^\top X_2 M = 0,$$

$$(4.22b) \quad M^\top X_1(T)M = -V,$$

$$(4.22c) \quad M^\top J_1 X_1 = 0 \quad \text{and} \quad J_1 X_1 M = 0.$$

The terminal condition (4.22b) is defined via (4.2d), and (4.7), namely,

$$M^\top \lambda(T) = M^\top X_1(T)M v(T) = -V v(T) \quad \text{or} \quad M^\top X_1(T)M = -V.$$

To show that (4.22) has a solution, we consider (4.22a)–(4.22b) in the transformed variables $X := -M^\top X_1 M$ and $Y := X_2 M$:

$$(4.23) \quad \begin{aligned} -\dot{X} - A^\top M^{-\top} X - X M^{-1} A + X M^{-1} G^{-\top} M^{-\top} X - W + J_2^\top Y + Y^\top J_2 &= 0, \\ X(T) &= V. \end{aligned}$$

By means of the projector $Q := M^{-1} J_1^\top [J_2 M^{-1} J_1^\top]^{-1} J_2$ we write $X = [Q^\top + P^\top] X [P + Q]$. From (4.22c) we obtain that $Q^\top X = X Q = 0$ and thus X is completely defined via $X_0 := P^\top X P$. Applying P^\top and P to (4.23) from the left and the right, respectively, we get a standard differential Riccati equation

$$\begin{aligned} -\dot{X}_0 - F_0^\top X_0 - X_0 A_0 + X_0 M^{-1} G M^{-\top} X_0 - P^\top W P &= 0, \\ X_0(T) &= P^\top V P, \end{aligned}$$

which has a unique and symmetric positive semidefinite solution (cf. [1, Thm. 4.1.6]), since V , G , and W are symmetric positive semidefinite. Again, the consistency condition (4.6) ensures that $X_0(T)$ also satisfies the initial condition and the algebraic constraints (4.22c). Since $Q^\top X = 0$ and $X Q = 0$, we have $X_1 = -M^{-\top} X M^{-1}$ is unique and symmetric negative semidefinite.

Application of P^\top from the left and Q from the right to (4.23) gives

$$-X_0 \dot{Q} - X_0 M^{-1} A Q - P^\top W Q = -P^\top Y^\top J_2 Q = -P^\top Y^\top J_2,$$

which is uniquely solvable for $P^\top Y^\top$. The projected equation obtained via Q^\top and P is the transpose of the above equation and bears no additional information.

Finally, one can determine $Q^\top Y^\top$ from the projection of (4.23) onto the range of Q^\top and Q which reads

$$(4.24) \quad -Q^\top W Q + Q^\top Y^\top J_2 Q + Q^\top J_2^\top Y Q = 0.$$

With $J_2 Q = J_2$, we find that (4.24) is of the form $[YQ]^\top J_2 + J_2^\top [YQ] = Q^\top WQ$ that was investigated in [14]. With $Q^- := M^{-1} J_1^\top [J_2 M^{-1} J_1^\top]^{-1}$ being a generalized inverse to J_2 , we obtain the projectors $P_1 := Q^- J_2 = Q$ and $P_2 := J_2 Q^- = I$ and the existence of solutions to (4.24) follows by [14, Thm. 1], since $-Q^\top WQ$ is symmetric and $-[I - P_1]^\top Q^\top WQ [I - P_1] = 0$.

The general solution to (4.24) is given by

$$YQ = -\frac{1}{2}[J_1 M^{-\top} J_2^\top]^{-1} J_1 M^{-\top} WQ + ZJ_2,$$

where Z is arbitrary with $Z^\top = -Z$. Thus existence of $M^\top X_1 M$ and $M^\top X_2^\top = Y^\top = P^\top Y^\top + Q^\top Y^\top$ and therefore X_1 and X_2 is proved.

By construction, with X_1 and X_2 as determined above, the solution of

$$\begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix} - \left(\begin{bmatrix} G & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} X_1 & X_2^\top \\ X_2 & 0 \end{bmatrix} \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} A & J_1^\top \\ J_2 & 0 \end{bmatrix} \right) \begin{bmatrix} v \\ p \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix},$$

$$v(0) = \alpha,$$

and

$$\begin{bmatrix} \lambda \\ \mu \end{bmatrix} = \begin{bmatrix} X_1 & X_2^\top \\ X_2 & 0 \end{bmatrix} \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix}$$

gives the solution of (4.2) with a zero right-hand side.

ad (4) The result for the affine-linear case is obtained via the affine-linear ansatz (4.9) using similar arguments as in the proof of (3). Proceeding analogously to the first steps for part (3), but with the affine linear ansatz (4.9) instead of the linear (4.7), we come to the expression

$$\begin{aligned} \mathcal{X} \begin{bmatrix} v \\ p \end{bmatrix} + \frac{d}{dt} \left(\begin{bmatrix} M^\top & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \right) - \begin{bmatrix} M^\top X_1 G + A^\top & J_2^\top \\ J_1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \\ (4.26) \qquad \qquad \qquad = \begin{bmatrix} Wv^* + M^\top X_1 f + M^\top X_2 g \\ 0 \end{bmatrix}, \end{aligned}$$

where \mathcal{X} is as in (4.21). Again, the requirement $\mathcal{X} = 0$ uniquely defines X_1 and $X_2 := \tilde{X}_2 + ZJ_2 M^{-1}$ up to an arbitrary skew-symmetric matrix Z . We write $M^\top X_2^\top g = M^\top \tilde{X}_2^\top g - J_2^\top Zg$ and define $\tilde{w}_2 := w_2 - Zg$. With this, (4.26), and an according decomposition of the terminal condition give a system for (w_1, \tilde{w}_2) ,

$$\begin{aligned} (4.27a) \qquad \qquad \qquad \begin{bmatrix} -\frac{d}{dt}(M^\top w_1) \\ 0 \end{bmatrix} - \begin{bmatrix} M^\top X_1 G + A^\top & J_2^\top \\ J_1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ \tilde{w}_2 \end{bmatrix} = \begin{bmatrix} Wv^* + M^\top X_1 f + M^\top \tilde{X}_2 g \\ 0 \end{bmatrix}, \\ (4.27b) \qquad \qquad \qquad M^\top w_1(T) = -Vv^*(T), \end{aligned}$$

which is of type (3.1). Since by (4.6) the terminal condition is consistent, system (4.27) has a unique solution (w_1, \tilde{w}_2) ; cf. Theorem 3.6. In particular, w_1 is independent of Zg ; cf. (3.7) and (3.9). For the solution w_2 of (4.26) we have $w_2 = \tilde{w}_2 + Zg$. Thus the existence of the functions used for the ansatz (4.9) is shown.

By construction, we have that the ansatz (4.9) leads to the solution of the Euler-Lagrange equations (4.2) via the decoupled system

$$\begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix} - \left(\begin{bmatrix} G & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} X_1 & X_2^\top \\ X_2 & 0 \end{bmatrix} \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} A & J_1^\top \\ J_2 & 0 \end{bmatrix} \right) \begin{bmatrix} v \\ p \end{bmatrix} = \begin{bmatrix} f + Gw_1 \\ g \end{bmatrix},$$

$$v(0) = \alpha,$$

and

$$\begin{bmatrix} \lambda \\ \mu \end{bmatrix} = \begin{bmatrix} X_1 & X_2^T \\ X_2 & 0 \end{bmatrix} \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} + \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}. \quad \square$$

Remark 4.6. For the considered setup, the sufficient conditions for solvability (4.6) coincide with the conditions derived in [5, Thm. 4].

Remark 4.7. The solution of (4.8) is unique up to an additive term ZJ_2M^{-1} in X_2 , with an arbitrary matrix Z , that fulfills $Z^T = -Z$. However, this does not contradict the unique solvability of the Euler–Lagrange equations, since λ and μ as defined via (4.7) are independent of Z .

In view of optimal control, the above results can be summarized as follows. To obtain an optimal input u for (3.1) with respect to a cost functional of type (4.1) it is sufficient to have a solution of the associated Euler–Lagrange equations (4.2); cf. [39]. By Lemma 4.5 it follows that for the considered state equations and cost functionals this solution exists, that it is unique, that it can be obtained via the separation ansatz (4.7), and that an optimal u is obtained via expression (4.3). For the inhomogeneous and for the trajectory tracking case, one can use an affine linear Riccati ansatz; cf. [34]. Thus, we can state the following theorem.

THEOREM 4.8. *Let $T > 0$ and consider the time interval $[0, T]$, let $n_u, n_v, n_p \in \mathbb{N}$, $n_v > n_p$, $M \in \mathcal{C}(0, T; \mathbb{R}^{n_v, n_v})$ be pointwise invertible, $A \in \mathcal{C}(0, T; \mathbb{R}^{n_v, n_v})$, and let $J_1, J_2 \in \mathcal{C}(0, T; \mathbb{R}^{n_p, n_v})$, such that $J_2M^{-1}J_1^T$ is invertible and such that $M^{-1}J_1^T S^{-1}J_2$ is differentiable. Let $B_1 \in \mathcal{C}(0, T; \mathbb{R}^{n_v, \nu})$, let $W, V \in \mathbb{R}^{n_v, n_v}$ be symmetric positive semidefinite, and let $R \in \mathbb{R}^{n_u, n_u}$ be symmetric positive definite. For a given $v^* \in \mathcal{C}^1(0, T; \mathbb{R}^{n_v})$ consider the linear-quadratic optimal control problem of finding $u \in \mathcal{C}(0, T; \mathbb{R}^{n_u})$ such that*

$$\frac{1}{2}[v - v^*]^T V [v - v^*] \Big|_{t=T} + \frac{1}{2} \int_0^T [v - v^*]^T W [v - v^*] + u^T R u \, dt$$

is minimal, where v and u on $(0, T]$ satisfy the state equations

$$\begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{v} \\ \dot{p} \end{bmatrix} - \begin{bmatrix} A & J_1^T \\ J_2 & 0 \end{bmatrix} \begin{bmatrix} v \\ p \end{bmatrix} - \begin{bmatrix} B_1 \\ 0 \end{bmatrix} u = \begin{bmatrix} f \\ g \end{bmatrix}, \quad v(0) = \alpha.$$

If $f \in \mathcal{C}(0, T; \mathbb{R}^{n_v})$, $g \in \mathcal{C}^1(0, T; \mathbb{R}^{n_p})$, and if

$$-J_2\alpha = g(0) \quad \text{and} \quad J_1M^{-T}V = 0,$$

then the optimal control problem is solvable and an optimal input u is given via

$$u = R^{-1}B_1^T[X_1Mv + w_1],$$

where X_1 and w_1 are the unique solutions of

$$\begin{aligned} \frac{d}{dt}(M^T X_1 M) + A^T X_1 M + M^T X_1 A + M^T X_1 B_1 R^{-1} B_1^T X_1 M - W \\ + J_2^T X_2 M + M^T X_2^T J_2 = 0, \\ M^T X_1(T)M = -V, \\ J_1 X_1 M = 0 \quad \text{and} \quad M^T X_1 J_1^T = 0, \end{aligned}$$

and

$$\begin{aligned} -\frac{d}{dt}(M^\top w_1) - [M^\top X_1 G + A^\top]w_1 - J_2^\top w_2 &= Wv^* + M^\top X_1 f + M^\top X_2 g, \\ M^\top w_1(T) &= Vv^*(T), \\ J_1 w_1 &= 0. \end{aligned}$$

Theorem 4.8 gives—in particular—sufficient optimality conditions in terms of the original variables and coefficients. As laid out in Remark 4.2, the necessity of these conditions is not guaranteed in general, since an inconsistent V makes them ill-posed. For practical applications, the following modification may be considered.

Remark 4.9. By Theorem 3.6 one has that if v solves (3.1), then it can be expressed as $v = \mathcal{P}v - c$, with $c := M^{-1}J_1^\top S^{-1}g$ independent of u and v . Therefore, the terminal point penalization in the cost functional (4.1) can be replaced like

$$\frac{1}{2}[v - v^*]^\top(T)V[v(T) - v^*] \leftarrow \frac{1}{2}[\mathcal{P}v(T) - c(T) - v^*]^\top V[\mathcal{P}v(T) - c(T) - v^*].$$

With this equivalent formulation, the terminal condition on $M^\top \lambda$ in (4.2d) coming from the variation of the cost functional with respect to v reads $M^\top \lambda(T) = -\mathcal{P}^\top V[\mathcal{P}v(T) - c(T) - v^*]$. Then the end condition for the gain matrix X_1 is given via $\mathcal{P}^\top V\mathcal{P}$ and for the affine part w_1 via $M^\top w_1(T) = \mathcal{P}^\top V[v^*(T) + M^{-1}J_1^\top S^{-1}g(T)]$. Both conditions are consistent since $J_1 M^{-\top} \mathcal{P}^\top = 0$. With this modification, in Theorem 4.8, the restriction $J_1 M^{-\top} V = 0$ is obsolete.

5. Numerical example. We illustrate the practical use of the results of Theorem 4.8 for an output tracking problem constrained by linear Navier–Stokes equations.

Therefore, we consider a lid driven cavity on the unit square $\Omega = [0, 1]^2$, with the boundary Γ , for time $t \in (0, 0.2]$, modeled by the incompressible Navier–Stokes equations for the velocity \mathbf{v} and the pressure \mathbf{p} ,

$$(5.1a) \quad \dot{\mathbf{v}} + (\mathbf{v} \cdot \nabla)\mathbf{v} - \frac{1}{Re}\Delta\mathbf{v} + \nabla\mathbf{p} = 0,$$

$$(5.1b) \quad \nabla \cdot \mathbf{v} = 0,$$

for a given *Reynolds number* $Re = 200$, completed by boundary conditions for $\mathbf{v}|_\Gamma$ and an initial condition for $\mathbf{v}(0)$.

We introduce distributed control and observation on the unit square as described in [27, Chap. 9.3]; see Figure 1 for an illustration. Therefore, we define input and output state spaces U and Y as subspaces of $L^2([0, 1])$ of piecewise linear functions on three equally sized segments. We define the input operator $B: U \times U \rightarrow [L^2(\Omega)]^2$ such that the control acts distributed in the domain of control Ω_c with two spatial components. Also, we consider an output operator that is defined such that for $v(t) \in [L^2(\Omega)]^2$ the output $Cv(t) = y(t) \in Y \times Y$ corresponds to the velocity in the domain of observation Ω_o averaged in the x_2 direction and projected onto $Y \times Y$ in the x_1 direction.

A spatial discretization of (5.1) with the input operator using mixed finite elements leads to

$$(5.2a) \quad \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \frac{d}{dt} \begin{bmatrix} v \\ p \end{bmatrix} (t) - \begin{bmatrix} N(v(t)) & J^\top \\ J & 0 \end{bmatrix} \begin{bmatrix} v(t) \\ p(t) \end{bmatrix} - \begin{bmatrix} B \\ 0 \end{bmatrix} u(t) = \begin{bmatrix} \hat{f} \\ 0 \end{bmatrix},$$

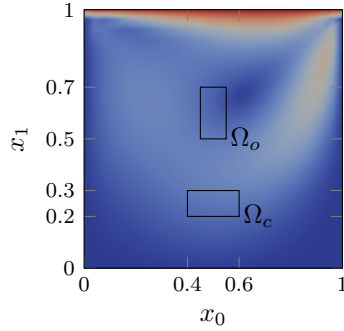


FIG. 1. Illustration of the driven cavity problem showing the velocity magnitude at the steady state for Reynolds number $Re = 200$ and the domain of control and observation $\Omega_c = [0.4, 0.6] \times [0.2, 0.3]$ and $\Omega_o = [0.45, 0.55] \times [0.5, 0.7]$.

on $(0, T]$, and

$$(5.2b) \quad v(0) = \alpha,$$

where v and p are finite dimensional approximations to \mathbf{v} and \mathbf{p} and where, among others, $N(v)v$ is a discrete approximation to $(\mathbf{v} \cdot \nabla)\mathbf{v} - \frac{1}{Re}\Delta\mathbf{v}$ and J and J^T are discrete approximations of the divergence and the gradient; cf. [24]. Finally, taking a velocity trajectory v_0 and defining $A(t)v = N(v_0(t))v + N(v)v_0(t)$ and $f := N(v_0)v_0 + \hat{f}$, we obtain the linearized state equations:

$$(5.3a) \quad \begin{bmatrix} M & 0 \\ 0 & 0 \end{bmatrix} \frac{d}{dt} \begin{bmatrix} v \\ p \end{bmatrix} (t) - \begin{bmatrix} A(t) & J^T \\ J & 0 \end{bmatrix} \begin{bmatrix} v(t) \\ p(t) \end{bmatrix} - \begin{bmatrix} B \\ 0 \end{bmatrix} u(t) = \begin{bmatrix} f \\ 0 \end{bmatrix},$$

on $(0, T]$, and

$$(5.3b) \quad v(0) = \alpha.$$

The linearized state equations are taken as the constraints when minimizing the cost functional

$$(5.4) \quad \mathcal{J}(v, u) = \frac{\gamma}{2} \|Cv(T) - y^*(T)\|_{Y \times Y}^2 + \frac{1}{2} \int_0^T \|Cv(t) - y^*(t)\|_{Y \times Y}^2 + \beta \|u(t)\|_{U \times U}^2 dt,$$

where $T = 0.2$, where $y^* \in \mathcal{C}(0, T; Y \times Y)$ is a target signal, and $\gamma, \beta > 0$ are parameters.

For the spatial discretization, we used *Taylor-Hood* finite elements on a regular triangulation with 2500 triangles. For the time discretization, we employed the implicit trapezoidal rule on a discretization of the time interval $[0, 0.2]$ with 128 time instances that are clustered toward the marginal points. For the linearization, we took the velocity solution to (5.2) without control input that starts with the steady-state Stokes solution α , which is also taken as the initial value for the linearized system (5.3).

With the *LBB* stability of the *Taylor-Hood* elements and with the modification of the terminal costs in (5.4) as proposed in Remark 4.9 we have that the optimization problem (5.4) constrained by (5.3) fulfills all assumptions of Theorem 4.8; cf. Remark 3.5.

Thus, the optimal control can be obtained via the feedback gain matrix X_1 and the affine correction w_1 solving a constrained differential Riccati equation and a backward in time linear DAE; see Theorem 4.8. To solve the constrained differential Riccati equation, we applied an implicit Euler scheme for time-stepping and a *Newton-ADI* iteration to solve the resulting constrained algebraic Riccati equation. See [27] for details on the algorithm and the author's *github* account [26] for the code and parameters used.

In the reported setup, we considered 5400 spatial degrees of freedom in the state space and 8 degrees of freedom in the input and outputs space. In the cost functional (5.4), we fixed $\gamma = 10^{-1}$ and solved the problem for various $\beta \in \{10^{-7}, 10^{-9}, 10^{-11}\}$. As shown in Figure 2, for a given target trajectory, the computed feedback matrices and affine corrections force the system into a trajectory that approximates the target.

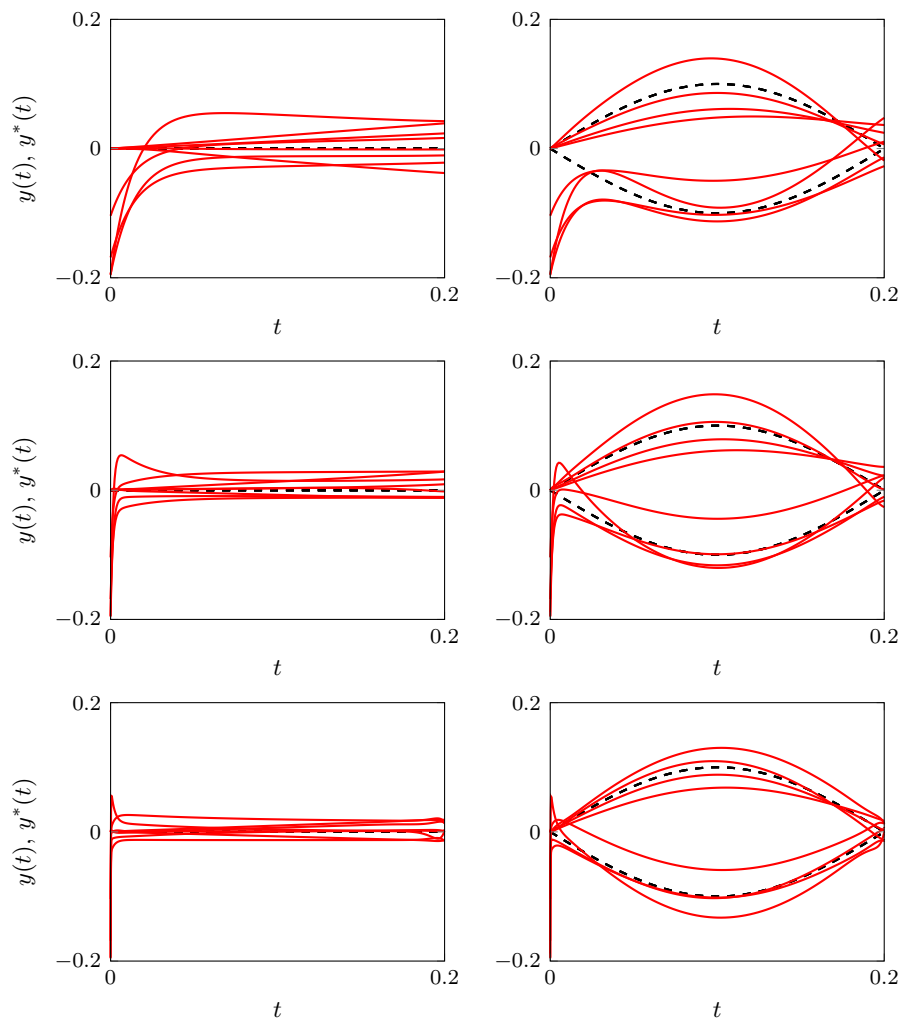


FIG. 2. Illustration of the eight components of the output signals (solid red) corresponding to the optimized input for a constant and a time dependent target signal (dashed) and for varying input penalization parameter $\beta = 10^{-7}$ (top), $\beta = 10^{-9}$ (middle), and $\beta = 10^{-11}$ (bottom row).

For smaller values of β , which increases the weight on the tracking as opposed to the penalization of the input, the approximation improves. However, for ever smaller β , the underlying *Newton-ADI* iteration converges more slowly. For $\beta = 10^{-13}$, with the current implementation, no results could be obtained.

6. Discussion and outlook. As the numerical example indicates, the proposed methodology is well suited for optimal control problems for linearized Navier–Stokes equations with time-varying control targets. The feedback representation with a negative semidefinite gain matrix seems beneficial for the task of stabilizing a trajectory stemming, e.g., from an *open-loop* optimization. The numerical framework allows for the incorporation of recent developments for large-scale Riccati equations associated with flow control problems addressing observer based control [12] or the efficient solution of the arising linear equation systems [10, 13] so that more complex setups will not pose severe problems in terms of dimensionality.

In view of applications, it is worth investigating how models for boundary control and for pressure measurements can be adjusted to fit the presented framework. By now, the exclusion of control action in the constraint equation, i.e., no B_2 in (3.1), makes the approach inapplicable to standard models for boundary control of flows. A possible remedy is proposed in [28], namely, to partially resolve the constrained equation for the part of the velocity that is affected by B_2 . Alternatively, one may consider the weak imposition of the boundary control (see, e.g., [11]) to avoid the occurrence of B_2 after spatial discretizations. If one includes the pressure in the cost functional, the differential-algebraic structure of the formal adjoint equation may completely change [27, Chap. 8.6]. A solution might be the use of (3.8) to replace p in the cost functional by terms of v , \dot{v} , and u .

Another possible extension of the results is directed to systems of higher index but with a similar structure. The general structure of the considered state equations (3.1) can be seen as a dynamical equation for v with a constrained coupled to it through the Lagrange multiplier p . The same structure can be found in mechanical or multibody systems [44] with the difference that the dynamical differential equation is of second order in time and that the constraints are often nonlinear. Since the proposed projections and splittings act in the state space independently of the order of the time derivatives, in the case of linear constraints, Theorem 3.6 can be extended to decouple these types of *index-3* state equations. Those parts of Lemma 4.8 that concern the decoupling of the formal optimality conditions also apply to second order systems. In the case of nonlinear constraints, the formulation of the results with time dependent projectors may be applied to linearized versions. For theoretical considerations, it may be worth extending the proposed decoupling for nonlinear constraints using similar *implicit function* arguments as provided in [2].

REFERENCES

- [1] H. ABOU-KANDIL, G. FREILING, V. IONESCU, AND G. JANK, *Matrix Riccati Equations in Control and Systems Theory*, Birkhäuser, Basel, Switzerland, 2003.
- [2] R. ALTMANN AND J. HEILAND, *Regularization of Constrained PDEs of Semi-Explicit Structure*, Preprint 2014–05, Institut für Mathematik, Technische Universität Berlin, 2014.
- [3] R. ALTMANN AND J. HEILAND, *Finite element decomposition and minimal extension for flow equations*, ESAIM Math. Model. Numer. Anal., 49 (2015), pp. 1489–1509.
- [4] U. M. ASCHER, R. M. M. MATTHEIJ, AND R. D. RUSSELL, *Numerical Solution of Boundary Value Problems for Ordinary Differential Equations*, SIAM, Philadelphia, 1995.
- [5] A. BACKES, *Optimale Steuerung der linearen DAE im Fall Index 2*, Preprint 2003–04, Institut für Mathematik, Humboldt-Universität zu Berlin, 2003.

- [6] A. BACKES, *Extremalbedingungen für Optimierungs-Probleme mit Algebro-Differentialgleichungen*, Ph.D. thesis, Institut für Mathematik, Humboldt-Universität zu Berlin, 2006.
- [7] K. BALLA, G. A. KURINA, AND R. MÄRZ, *Index criteria for differential algebraic equations arising from linear-quadratic optimal control problems*, J. Dyn. Control Syst., 12 (2006), pp. 289–311.
- [8] K. BALLA AND V. H. LINH, *Adjoint pairs of differential-algebraic equations and Hamiltonian systems*, Appl. Numer. Math., 53 (2005), pp. 131–148.
- [9] K. BALLA AND R. MÄRZ, *A unified approach to linear differential algebraic equations and their adjoints*, Z. Anal. Anwend., 21 (2002), pp. 783–802.
- [10] E. BÄNSCH, P. BENNER, J. SAAK, AND H. K. WEICHELT, *Riccati-based boundary feedback stabilization of incompressible Navier-Stokes flows*, SIAM J. Sci. Comput., 37 (2015), pp. A832–A858.
- [11] F. B. BELGACEM, H. EL FEKIH, AND J.-P. RAYMOND, *A penalized Robin approach for solving a parabolic equation with nonsmooth Dirichlet boundary conditions*, Asymptot. Anal., 34 (2003), pp. 121–136.
- [12] P. BENNER AND J. HEILAND, *LQG-Balanced truncation low-order controller for stabilization of laminar flows*, in Active Flow and Combustion Control 2014, Notes Numer. Fluid Mech. Multidiscip. Des. 127, R. King, ed., Springer, Berlin, 2015, pp. 365–379.
- [13] P. BENNER, J. SAAK, M. STOLL, AND H. WEICHELT, *Efficient solution of large-scale saddle point systems arising in Riccati-based boundary feedback stabilization of incompressible Stokes flow*, SIAM J. Sci. Comput., 35 (2013), pp. S150–S170.
- [14] H. W. BRADEN, *The equations $A^T X \pm X^T A = B$* , SIAM J. Matrix Anal. Appl., 20 (1998), pp. 295–302.
- [15] S. CAMPBELL, P. KUNKEL, AND V. MEHRMANN, *Regularization of linear and nonlinear descriptor systems*, in Control and Optimization with Differential-Algebraic Constraints, Adv. Des. Control 23, SIAM, Philadelphia, 2012, pp. 17–36.
- [16] M. DO R DE PINHO AND R. B. VINTER, *Necessary conditions for optimal control problems involving nonlinear differential algebraic equations*, J. Math. Anal. Appl., 212 (1997), pp. 493–516.
- [17] L. DIECI, *Numerical integration of the differential Riccati equation and some related issues*, SIAM J. Numer. Anal., 29 (1992), pp. 781–815.
- [18] L. DIECI, *On the decoupling of dichotomic linear Hamiltonians. Considerations on integrating symmetric differential Riccati equations*, J. Comput. Appl. Math., 45 (1993), pp. 47–63.
- [19] L. DIECI, M. R. OSBORNE, AND R. D. RUSSELL, *A Riccati transformation method for solving linear BVPs. I: Theoretical aspects*, SIAM J. Numer. Anal., 25 (1988), pp. 1055–1073.
- [20] M. GERDTS, *Optimal control and real-time optimization of mechanical multi-body systems*, Z. Angew. Math. Mech., 83 (2003), pp. 705–719.
- [21] M. GERDTS, *Local minimum principle for optimal control problems subject to differential-algebraic equations of index two*, J. Optim. Theory Appl., 130 (2006), pp. 443–462.
- [22] M. GERDTS AND M. KUNKEL, *A globally convergent semi-smooth Newton method for control-state constrained DAE optimal control problems*, Comput. Optim. Appl., 48 (2011), pp. 601–633.
- [23] V. GIRAULT AND P.-A. RAVIART, *Finite Element Methods for Navier–Stokes Equations. Theory and Algorithms*, Springer, Berlin, 1986.
- [24] P. M. GRESHO AND R. L. SANI, *Incompressible Flow and the Finite Element Method. Vol. 2: Isothermal Laminar Flow*, Wiley, Chichester, UK, 2000.
- [25] E. HAIRER AND G. WANNER, *Solving Ordinary Differential Equations II; Stiff and Differential-Algebraic Problems*, 2nd ed., Springer Ser. Comput. Math., Springer, Berlin, 2002.
- [26] J. HEILAND, *Optconpy—A Python Module for the Solution of DAE-Riccati Equations via Newton-ADI and Application Example*, Release V 3.0, <https://github.com/highlando/optconpy> (2013).
- [27] J. HEILAND, *Decoupling and Optimization of Differential-Algebraic Equations with Application in Flow Control*, Ph.D. thesis, TU Berlin, 2014; available online from <http://opus4.kobv.de/opus4-tuberlin/frontdoor/index/index/docId/5243>.
- [28] M. HEINKENSCHLOSS, D. C. SORENSEN, AND K. SUN, *Balanced truncation model reduction for a class of descriptor systems with application to the Oseen equations*, SIAM J. Sci. Comput., 30 (2008), pp. 1038–1063.
- [29] M. HINZE, *Optimal and Instantaneous Control of the Instationary NAVier-STokes Equations*, Habilitationsschrift, Institut für Mathematik, Technische Universität Berlin, 2000.
- [30] P. KUNKEL AND V. MEHRMANN, *The linear quadratic optimal control problem for linear descriptor systems with variable coefficients*, Math. Control Signals Syst., 10 (1997), pp. 247–264.

- [31] P. KUNKEL AND V. MEHRMANN, *Differential-Algebraic Equations. Analysis and Numerical Solution*, European Mathematical Society Publishing House, Zürich, Switzerland, 2006.
- [32] P. KUNKEL AND V. MEHRMANN, *Optimal control for unstructured nonlinear differential-algebraic equations of arbitrary index*, *Math. Control Signals Syst.*, 20 (2008), pp. 227–269.
- [33] P. KUNKEL AND V. MEHRMANN, *Formal adjoints of linear DAE operators and their role in optimal control*, *Electron. J. Linear Algebra*, 22 (2011), pp. 672–693.
- [34] P. KUNKEL AND V. MEHRMANN, *Optimal control for linear descriptor systems with variable coefficients*, in *Numerical Linear Algebra in Signals, Systems and Control*, P. Van Dooren, S. P. Bhattacharyya, R. H. Chan, V. Olshevsky, and A. Routray, eds., Springer, Berlin, 2011, pp. 313–339.
- [35] P. KUNKEL, V. MEHRMANN, AND L. SCHOLZ, *Self-adjoint differential-algebraic equations*, *Math. Control Signals Systems*, 26 (2014), pp. 47–76.
- [36] G. A. KURINA AND R. MÄRZ, *On linear-quadratic optimal control problems for time-varying descriptor systems*, *SIAM J. Control Optim.*, 42 (2004), pp. 2062–2077.
- [37] G. A. KURINA AND R. MÄRZ, *Feedback solutions of optimal control problems with DAE constraints*, *SIAM J. Control Optim.*, 46 (2007), pp. 1277–1298.
- [38] R. LAMOUR, R. MÄRZ, AND C. TISCHENDORF, *Differential-algebraic equations: A projector based analysis*, in *Differential-Algebraic Equations Forum*, Springer, Heidelberg, 2013.
- [39] R. MÄRZ, *Adjoint equations of differential-algebraic systems and optimal control problems*, *Proc. Inst. Math. NAS Belarus*, 7 (2001), pp. 88–97.
- [40] R. MÄRZ, *Differential algebraic systems anew.*, *Appl. Numer. Math.*, 42 (2002), pp. 315–335.
- [41] V. MEHRMANN, *Index Concepts for Differential-Algebraic Equations*, Tech. report 3–2012, Institut für Mathematik, Technische Universität Berlin, 2012.
- [42] W. T. REID, *Riccati Differential Equations*, Academic Press, New York, 1972.
- [43] T. ROUBIČEK AND M. VALÁŠEK, *Optimal control of causal differential-algebraic systems*, *J. Math. Anal. Appl.*, 269 (2002), pp. 616–641.
- [44] B. SIMEON, *DAEs and PDEs in elastic multibody systems*, *Numer. Algorithms*, 19 (1998), pp. 235–246.