

# CONVERGENCE RATE FOR A RADAU COLLOCATION METHOD APPLIED TO UNCONSTRAINED OPTIMAL CONTROL \*

WILLIAM W. HAGER<sup>†</sup>, HONGYAN HOU<sup>‡</sup>, AND ANIL V. RAO<sup>§</sup>

**Abstract.** A local convergence rate is established for an orthogonal collocation method based on Radau quadrature applied to an unconstrained optimal control problem. If the continuous problem has a sufficiently smooth solution and the Hamiltonian satisfies a strong convexity condition, then the discrete problem possesses a local minimizer in a neighborhood of the continuous solution, and as the number of collocation points increases, the discrete solution convergences exponentially fast in the sup-norm to the continuous solution. An earlier paper analyzes an orthogonal collocation method based on Gauss quadrature, where neither end point of the problem domain is a collocation point. For the Radau quadrature scheme, one end point is a collocation point.

**Key words.** Radau collocation method, convergence rate, optimal control, orthogonal collocation

**1. Introduction.** A convergence rate is established for an orthogonal collocation method applied to an unconstrained control problem of the form

$$\begin{aligned} & \text{minimize} && C(\mathbf{x}(1)) \\ & \text{subject to} && \dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)), \quad t \in [-1, 1], \\ & && \mathbf{x}(-1) = \mathbf{x}_0, \quad (\mathbf{x}, \mathbf{u}) \in \mathcal{C}^1(\mathbb{R}^n) \times \mathcal{C}^0(\mathbb{R}^m) \end{aligned} \quad (1.1)$$

where the state  $\mathbf{x}(t) \in \mathbb{R}^n$ ,  $\dot{\mathbf{x}} \equiv \frac{d}{dt}\mathbf{x}$ , the control  $\mathbf{u}(t) \in \mathbb{R}^m$ ,  $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ ,  $C : \mathbb{R}^n \rightarrow \mathbb{R}$ , and  $\mathbf{x}_0$  is the initial condition, which we assume is given;  $\mathcal{C}^k$  denotes the space of  $k$  times continuously differentiable functions. It is assumed that  $\mathbf{f}$  and  $C$  are at least continuous. Assuming the dynamics  $\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t))$  can be solved for the state  $\mathbf{x}$  as a function of the control  $\mathbf{u}$ , the control problem reduces to an unconstrained minimization over  $\mathbf{u}$ .

Let  $\mathcal{P}_N$  denote the space of polynomials of degree at most  $N$  defined on the interval  $[-1, +1]$ , and let  $\mathcal{P}_N^n$  denote the  $n$ -fold Cartesian product  $\mathcal{P}_N \times \dots \times \mathcal{P}_N$ . We analyze a discrete approximation to (1.1), introduced in [11, 12], of the form

$$\begin{aligned} & \text{minimize} && C(\mathbf{x}(1)) \\ & \text{subject to} && \dot{\mathbf{x}}(\tau_i) = \mathbf{f}(\mathbf{x}(\tau_i), \mathbf{u}_i), \quad 1 \leq i \leq N, \\ & && \mathbf{x}(-1) = \mathbf{x}_0, \quad \mathbf{x} \in \mathcal{P}_N^n. \end{aligned} \quad (1.2)$$

At the collocation points  $\tau_i$ ,  $1 \leq i \leq N$ , the equation should be satisfied. The control approximation at time  $\tau_i$  is  $\mathbf{u}_i$ . We focus on the Radau quadrature points satisfying

$$-1 < \tau_1 < \tau_2 < \dots < \tau_N = +1.$$

The analysis we give also applies to the flipped Radau scheme obtained by reversing the sign for each collocation point. In (1.2) the dimension of  $\mathcal{P}_N$  is  $N + 1$ , while there

---

\* August 17, 2015. Revised September 12, 2015. The authors gratefully acknowledge support by the Office of Naval Research under grants N00014-11-1-0068 and N00014-15-1-2048, and by the National Science Foundation under grants DMS-1522629 and CBET-1404767.

<sup>†</sup>[hager@ufl.edu](mailto:hager@ufl.edu), <http://people.clas.ufl.edu/hager/>, PO Box 118105, Department of Mathematics, University of Florida, Gainesville, FL 32611-8105. Phone (352) 294-2308. Fax (352) 392-8357.

<sup>‡</sup>[hongyan388@gmail.com](mailto:hongyan388@gmail.com), Chemical Engineering, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213.

<sup>§</sup>[anilvrao@ufl.edu](mailto:anilvrao@ufl.edu), <http://www.mae.ufl.edu/rao> Department of Mechanical and Aerospace Engineering, P.O. Box 116250, Gainesville, FL 32611-6250. Phone:(352) 392-0961. Fax:(352) 392-7303

are  $N + 1$  equations in (1.2) corresponding to the collocated dynamics at  $N$  points and the initial condition. When the discrete dynamics is nice enough, we can solve for the discrete state  $\mathbf{x} \in \mathcal{P}_N^n$  as a function of the discrete controls  $\mathbf{u}_i$ ,  $1 \leq i \leq N$ , and the discrete approximation reduces to an unconstrained minimization over the discrete controls.

In an earlier paper [16] we analyzed a scheme based on Gauss quadrature, where the collocation points lie in the interior of the interval  $[-1, +1]$ . The Gauss scheme is easier to analyze than the Radau scheme of this paper due to the symmetry of the Gauss collocation points, and the fact that the none of the collocation points lies at an end of the problem domain. For the Radau scheme, symmetry is broken and the polynomials in the discrete adjoint equation have degree  $N - 1$  compared to the degree  $N$  polynomials used for the state approximation. As will be seen, the presence of the Radau collocation point at the end of the interval  $[-1, +1]$  leads to the embedding of the terminal adjoint condition of the discrete problem into the discrete adjoint dynamics. Despite these differences in the analysis and despite the fact that Gauss quadrature has a higher degree of accuracy than Radau quadrature, the convergence rate obtained for the Radau scheme is exactly the same as that of the Gauss scheme. Moreover, in numerical experiments with test problems where an exact solution is known, the observed error in the state, control, and adjoint for the Radau scheme is very similar to the observed error for the Gauss scheme. The fact that a Radau quadrature point can be placed at the end point of the interval leads to a simpler implementation of terminal constraints and terminal cost. And the Radau scheme is easier than the Gauss scheme to extend to an *hp*-framework [1, 2, 20, 22] where the interval  $[-1, +1]$  is partitioned into a mesh and a different polynomial is employed in each mesh interval.

The analysis in this and the earlier paper [16] needs further extensions in order to handle Lobatto collocation schemes such as those in [7, 9] where  $\tau_1 = -1$  and  $\tau_N = +1$ . Although the Gauss and Radau scheme have similar errors, the Lobatto scheme can converge at a slower rate, as observed in [12]. An advantage of the Lobatto scheme is that the value of the optimal control is estimated at the initial point  $t = -1$  and the terminal point  $t = +1$ . Moreover, both initial and terminal constraints are easier to implement in the Lobatto framework. A consistency result for a Lobatto collocation scheme applied to optimal control is given in [13]. Other quadrature points that have been exploited in the optimal control literature include the Chebyshev quadrature points [8, 10], and the extrema of Jacobi polynomials [26].

Our goal is to show that if  $(\mathbf{x}^*, \mathbf{u}^*)$  is a local minimizer for (1.1), then the discrete problem (1.2) has a local minimizer  $(\mathbf{x}^N, \mathbf{u}^N)$  that converges exponentially fast in  $N$  to  $(\mathbf{x}^*, \mathbf{u}^*)$  at the collocation points. Convergence rates have been obtained previously when the approximating space consists of piecewise polynomials as in [3, 4, 6, 5, 14, 19, 23]. In these earlier results, convergence is achieved by letting the mesh spacing tend to zero. In our results, on the other hand, convergence is achieved by letting  $N$ , the degree of the approximating polynomials, tend to infinity.

Let  $\mathcal{C}^k(\mathbb{R}^n)$  denote the space of  $k$  times continuously differentiable functions  $\mathbf{x} : [-1, +1] \rightarrow \mathbb{R}^n$  with the sup-norm  $\|\cdot\|_\infty$  given by

$$\|\mathbf{x}\|_\infty = \sup\{|\mathbf{x}(t)| : t \in [-1, +1]\}, \quad (1.3)$$

where  $|\cdot|$  is the Euclidean norm. It is assumed that (1.1) has a local minimizer  $(\mathbf{x}^*, \mathbf{u}^*)$  in  $\mathcal{C}^1(\mathbb{R}^n) \times \mathcal{C}^0(\mathbb{R}^m)$ . Given  $\mathbf{y} \in \mathbb{R}^n$ , the ball with center  $\mathbf{y}$  and radius  $\rho$  is denoted

$$\mathcal{B}_\rho(\mathbf{y}) = \{\mathbf{x} \in \mathbb{R}^n : |\mathbf{x} - \mathbf{y}| \leq \rho\}.$$

It is assumed that there exists an open set  $\Omega \subset \mathbb{R}^{m+n}$  and  $\rho > 0$  such that

$$\mathcal{B}_\rho(\mathbf{x}^*(t), \mathbf{u}^*(t)) \subset \Omega \text{ for all } t \in [-1, +1].$$

Moreover, the first two derivative of  $f$  and  $C$  are continuous on the closure of  $\Omega$  and on  $\mathcal{B}_\rho(\mathbf{x}^*(1))$  respectively.

Let  $\boldsymbol{\lambda}^*$  denote the solution of the linear costate equation

$$\dot{\boldsymbol{\lambda}}^*(t) = -\nabla_x H(\mathbf{x}^*(t), \mathbf{u}^*(t), \boldsymbol{\lambda}^*(t)), \quad \boldsymbol{\lambda}^*(1) = \nabla C(\mathbf{x}^*(1)), \quad (1.4)$$

where  $H$  is the Hamiltonian defined by  $H(\mathbf{x}, \mathbf{u}, \boldsymbol{\lambda}) = \boldsymbol{\lambda}^\top \mathbf{f}(\mathbf{x}, \mathbf{u})$ . Here  $\nabla C$  denotes the gradient of  $C$ . By the first-order optimality conditions (Pontryagin's minimum principle), we have

$$\nabla_u H(\mathbf{x}^*(t), \mathbf{u}^*(t), \boldsymbol{\lambda}^*(t)) = \mathbf{0} \quad (1.5)$$

for all  $t \in [-1, +1]$ .

Since the discrete collocation problem (1.2) is finite dimensional, the first-order optimality conditions (Karush-Kuhn-Tucker conditions) imply that when a constraint qualification holds [21], the gradient of the Lagrangian vanishes. By the analysis in [11, 12], the gradient of the Lagrangian vanishes if and only if there exists  $\boldsymbol{\lambda} \in \mathcal{P}_{N-1}^n$  such that

$$\boldsymbol{\lambda}_0 = \boldsymbol{\lambda}(-1) \quad (1.6)$$

$$\dot{\boldsymbol{\lambda}}(\tau_i) = -\nabla_x H(\mathbf{x}(\tau_i), \mathbf{u}_i, \boldsymbol{\lambda}(\tau_i)), \quad 1 \leq i < N, \quad (1.7)$$

$$\dot{\boldsymbol{\lambda}}(1) = -\nabla_x H(\mathbf{x}(1), \mathbf{u}_N, \boldsymbol{\lambda}(1)) + (\boldsymbol{\lambda}(1) - \nabla C(\mathbf{x}(1)))/\omega_N, \quad (1.8)$$

$$\mathbf{0} = \nabla_u H(\mathbf{x}(\tau_i), \mathbf{u}_i, \boldsymbol{\lambda}(\tau_i)), \quad 1 \leq i \leq N, \quad (1.9)$$

where  $\omega_i$  is the Radau quadrature weight associated with  $\tau_i$ , and  $\boldsymbol{\lambda}_0$  is the multiplier associated with the initial condition  $\mathbf{x}(-1) = \mathbf{x}_0$  in (1.2). Note that in [12], (1.6) is written in the form

$$\nabla C(\mathbf{x}(1)) = \boldsymbol{\lambda}_0 - \sum_{i=1}^N \omega_i \nabla_x H(\mathbf{x}(\tau_i), \mathbf{u}_i, \boldsymbol{\lambda}(\tau_i)).$$

However, utilizing (1.7), (1.8), and the fundamental theorem of calculus, this reduces to the more compact form (1.6).

In comparing the first-order conditions for Radau collocation to the first-order conditions for Gauss collocation [16], the differences are that in Gauss collocation,  $\boldsymbol{\lambda} \in \mathcal{P}_N^n$  not  $\mathcal{P}_{N-1}^n$ . Moreover, in Gauss collocation, the terminal condition for the discrete adjoint is simply  $\boldsymbol{\lambda}(1) = \nabla C(\mathbf{x}(1))$ , while in Radau collocation, the discrete costate dynamics and the terminal condition are mixed together as in (1.8).

The assumptions utilized in the convergence analysis are the following:

(A1)  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^* \in \mathcal{C}^{\eta+1}$  for some  $\eta \geq 3$ .

(A2) For some  $\alpha > 0$ , the smallest eigenvalue of the Hessian matrices

$$\nabla^2 C(\mathbf{x}^*(1)) \quad \text{and} \quad \nabla_{(x,u)}^2 H(\mathbf{x}^*(t), \mathbf{u}^*(t), \boldsymbol{\lambda}^*(t))$$

is greater than  $\alpha$ , uniformly for  $t \in [-1, +1]$ .

(A3) The Jacobian of the dynamics satisfies

$$\|\nabla_x \mathbf{f}(\mathbf{x}^*(t), \mathbf{u}^*(t))\|_\infty \leq 1/4 \quad \text{and} \quad \|\nabla_x \mathbf{f}(\mathbf{x}^*(t), \mathbf{u}^*(t))^\top\|_\infty \leq 1/4$$

for all  $t \in [-1, +1]$  where  $\|\cdot\|_\infty$  is the matrix sup-norm (largest absolute row sum), and the Jacobian  $\nabla_x \mathbf{f}$  is an  $n$  by  $n$  matrix whose  $i$ -th row is  $(\nabla_x f_i)^\top$ .

The smoothness assumption (A1) is used to obtain a bound for the accuracy with which the interpolant of the continuous state  $\mathbf{x}^*$  satisfies the discrete dynamics. The coercivity assumption (A2) ensures that the solution of the discrete problem is a local minimizer. The condition (A3) enters into the analysis of stability for the perturbed dynamics; this condition can be eliminated when the problem domain  $[-1, +1]$  is partitioned into subintervals with a different polynomial on each subinterval [3, 4, 5, 6, 14, 19, 23]. For the global polynomials analyzed in this paper, (A3) could be replaced by any condition that ensures stability of the linearized dynamics.

In addition to the 3 assumptions, the analysis utilizes 4 properties of the Radau collocation scheme. Let  $\tau_0 = -1$ , a noncollocated point, and define

$$D_{ij} = \dot{L}_j(\tau_i), \quad 1 \leq i \leq N, \quad 0 \leq j \leq N, \quad \text{where } L_j(\tau) := \prod_{\substack{i=0 \\ i \neq j}}^N \frac{\tau - \tau_i}{\tau_j - \tau_i}. \quad (1.10)$$

Here the dot denotes differentiation, and  $\mathbf{D}$  is a differentiation matrix in the sense that  $(\mathbf{D}\mathbf{p})_i = \dot{p}(\tau_i)$ ,  $1 \leq i \leq N$ , whenever  $p \in \mathcal{P}_N$  is the polynomial that satisfies  $p(\tau_j) = p_j$  for  $0 \leq j \leq N$ . The submatrix  $\mathbf{D}_{1:N}$  consisting of the tailing  $N$  columns of  $\mathbf{D}$  has the following properties:

- (P1)  $\mathbf{D}_{1:N}$  is invertible and  $\|\mathbf{D}_{1:N}^{-1}\|_\infty = 2$ .
- (P2) If  $\mathbf{W}$  is the diagonal matrix containing the Radau quadrature weights  $\omega$  on the diagonal, then the rows of the matrix  $[\mathbf{W}^{1/2}\mathbf{D}_{1:N}]^{-1}$  have Euclidean norm bounded by  $\sqrt{2}$ .

The fact that  $\mathbf{D}_{1:N}$  is invertible is established in [11, Prop. 1], and a formula for the elements of  $\mathbf{D}_{1:N}^{-1}$  is given in [12, equation (53)]. From the formula, the elements in the last row of  $\mathbf{D}_{1:N}^{-1}$  are the Radau quadrature weights, which are positive and sum to 2. Although elements in the earlier rows of  $\mathbf{D}_{1:N}^{-1}$  can be either positive or negative, we find numerically that their absolute sum is always less than 2. Similarly, the elements in the last row of  $[\mathbf{W}^{1/2}\mathbf{D}_{1:N}]^{-1}$  are the square roots of the Radau quadrature weights. Hence, the Euclidean norm of the last row of  $[\mathbf{W}^{1/2}\mathbf{D}_{1:N}]^{-1}$  is  $\sqrt{2}$ . Numerically, we find that Euclidean norm of the earlier rows is always less than  $\sqrt{2}$ .

Let  $\mathbf{D}^\ddagger$  by the  $N$  by  $N$  matrix defined by

$$D_{ij}^\ddagger = - \left( \frac{\omega_j}{\omega_i} \right) D_{ji}, \quad 1 \leq i \leq N, \quad 1 \leq j \leq N.$$

The matrix  $\mathbf{D}^\ddagger$  arises in the analysis of the costate equation. In Section 4.2.1 of [12], we introduce a matrix  $\mathbf{D}^\dagger$  which is a differentiation matrix for the collocation points  $\tau_i$ ,  $1 \leq i \leq N$ . That is, if  $p$  is a polynomial of degree at most  $N - 1$  and  $\mathbf{p}$  is the vector with components  $p(\tau_i)$ ,  $1 \leq i \leq N$ , then  $(\mathbf{D}^\dagger \mathbf{p})_i = \dot{p}(\tau_i)$ . The matrix  $\mathbf{D}^\ddagger$  only differs from  $\mathbf{D}^\dagger$  in a single entry:  $D_{NN}^\ddagger = D_{NN}^\dagger - 1/\omega_N$ . As a result,

$$(\mathbf{D}^\ddagger \mathbf{p})_i = \dot{p}(\tau_i), \quad 1 \leq i < N, \quad (\mathbf{D}^\ddagger \mathbf{p})_N = \dot{p}(\tau_N) - p(1)/\omega_N. \quad (1.11)$$

If  $\mathbf{D}^\ddagger \mathbf{p} = \mathbf{0}$ , then  $\dot{p}(\tau_i) = 0$  for  $i < N$  by (1.11). Since  $\dot{p}$  has degree  $N - 2$  and it vanishes at  $N - 1$  points,  $\dot{p}$  is identically zero and  $p$  is constant. By the final equation

in (1.11),  $p(1) = 0$  when  $\mathbf{D}^\dagger \mathbf{p} = \mathbf{0}$ , which implies that  $p$  is identically zero. This shows that  $\mathbf{D}^\dagger$  is invertible. We find that  $\mathbf{D}^\dagger$  has the following properties:

- (P3)  $\mathbf{D}^\dagger$  is invertible and  $\|(\mathbf{D}^\dagger)^{-1}\|_\infty \leq 2$ .
- (P4) The rows of the matrix  $[\mathbf{W}^{1/2} \mathbf{D}^\dagger]^{-1}$  have Euclidean norm bounded by  $\sqrt{2}$ .

In Proposition 9.1 at the end of the paper, an explicit formula is given for the inverse of  $\mathbf{D}^\dagger$ . However, it is not clear from the formula that  $\|(\mathbf{D}^\dagger)^{-1}\|_\infty$  is bounded by 2. Numerically, we find that the norms in (P3) and (P4) achieve their maximum in the first row of the matrix, and these norms increase monotonically towards the given bounds. Properties (P1)–(P4) differ from the assumptions (A1)–(A3) in the sense that the properties seem to hold for any choice of  $N$ , although a proof is missing, while (A1)–(A3) only hold for certain control problems. In the analysis of the Gauss scheme [16], properties (P3) and (P4) followed immediately from (P1) and (P2) since it could be shown that the discrete costate matrix was related to the state differentiation matrix through an exchange operation. However, due to the asymmetry of the Radau collocation points and the lower degree of the polynomials in the discrete adjoint system (1.6)–(1.9), the relation between the state and costate matrices for the Radau scheme is not clear. Nonetheless, the bounds in (P3) and (P4) are observed to be the same as the bounds in (P1) and (P2).

If  $\mathbf{x}^N$  is a solution of (1.2) associated with the discrete controls  $\mathbf{u}_i$ ,  $1 \leq i \leq N$ , and if  $\boldsymbol{\lambda}^N \in \mathcal{P}_{N-1}^n$  satisfies (1.6)–(1.9), then we define

$$\begin{aligned} \mathbf{X}^N &= [ \mathbf{x}^N(\tau_0), \mathbf{x}^N(\tau_1), \dots, \mathbf{x}^N(\tau_N) ], \\ \mathbf{X}^* &= [ \mathbf{x}^*(\tau_0), \mathbf{x}^*(\tau_1), \dots, \mathbf{x}^*(\tau_N) ], \\ \mathbf{U}^N &= [ \mathbf{u}_1, \dots, \mathbf{u}_N ], \\ \mathbf{U}^* &= [ \mathbf{u}^*(\tau_1), \dots, \mathbf{u}^*(\tau_N) ], \\ \boldsymbol{\Lambda}^N &= [ \boldsymbol{\lambda}^N(\tau_0), \boldsymbol{\lambda}^N(\tau_1), \dots, \boldsymbol{\lambda}^N(\tau_N) ], \\ \boldsymbol{\Lambda}^* &= [ \boldsymbol{\lambda}^*(\tau_0), \boldsymbol{\lambda}^*(\tau_1), \dots, \boldsymbol{\lambda}^*(\tau_N) ]. \end{aligned}$$

For any of the discrete variables, we define a discrete sup-norm analogous to the continuous sup-norm in (1.3). For example, if  $\mathbf{U}^N \in \mathbb{R}^{mN}$  with  $\mathbf{U}_i \in \mathbb{R}^m$ , then

$$\|\mathbf{U}^N\|_\infty = \sup\{|\mathbf{U}_i| : 1 \leq i \leq N\}.$$

The following convergence result is established:

**THEOREM 1.1.** *If  $(\mathbf{x}^*, \mathbf{u}^*)$  is a local minimizer for the continuous problem (1.1) and both (A1)–(A3) and (P1)–(P4) hold, then for  $N$  sufficiently large with  $N > \eta + 1$ , the discrete problem (1.2) has a local minimizer  $(\mathbf{X}^N, \mathbf{U}^N)$  for which*

$$\max \{ \|\mathbf{X}^N - \mathbf{X}^*\|_\infty, \|\mathbf{U}^N - \mathbf{U}^*\|_\infty, \|\boldsymbol{\Lambda}^N - \boldsymbol{\Lambda}^*\|_\infty \} \leq cN^{2-\eta}, \quad (1.12)$$

where  $c$  is independent of  $N$ .

The discrete problem provides an estimate for optimal control at  $\tau_N = +1$ , however, there is no discrete control at  $\tau_0 = -1$ . Due to the strong convexity assumption (A2), an estimate for the discrete control at  $t = -1$  can be obtained from the minimum principle (1.5) since we have estimates for the discrete state and costate at  $\tau_0 = -1$ . Alternatively, polynomial interpolation could be used to obtain estimates for the optimal control at  $t = -1$ .

The paper is organized as follows. In Section 2 the discrete optimization problem (1.2) is reformulated as a nonlinear system of equations obtained from the first-order optimality conditions, and a general approach to convergence analysis is presented.

Section 3 obtains an estimate for how closely the solution to the continuous problem satisfies the first-order optimality conditions for the discrete problem. Section 4 proves that the linearization of the discrete control problem around a solution of the continuous problem is invertible. Section 5 establishes an  $L^2$  stability property for the linearization, while Section 6 strengthens the norm to  $L^\infty$ . This stability property is the basis for the proof of Theorem 1.1. A numerical example illustrating the exponential convergence result is given in Section 7.

**Notation.** The meaning of the norm  $\|\cdot\|_\infty$  is based on context. If  $\mathbf{x} \in C^0(\mathbb{R}^n)$ , then  $\|\mathbf{x}\|_\infty$  denotes the maximum of  $|\mathbf{x}(t)|$  over  $t \in [-1, +1]$ , where  $|\cdot|$  is the Euclidean norm. If  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , then  $\|\mathbf{A}\|_\infty$  is the largest absolute row sum (the matrix norm induces by the  $\ell_\infty$  vector norm). If  $\mathbf{U} \in \mathbb{R}^{mN}$  is the discrete control with  $\mathbf{U}_i \in \mathbb{R}^m$ , then  $\|\mathbf{U}\|_\infty$  is the maximum of  $|\mathbf{U}_i|$ ,  $1 \leq i \leq N$ . The dimension of the identity matrix  $\mathbf{I}$  is often clear from context; when necessary, the dimension of  $\mathbf{I}$  is specified by a subscript. For example,  $\mathbf{I}_n$  is the  $n$  by  $n$  identity matrix.  $\nabla C$  denotes the gradient, a column vector, while  $\nabla^2 C$  denotes the Hessian matrix. Throughout the paper,  $c$  denotes a generic constant which has different values in different equations. The value of this constant is always independent of  $N$ , the degree of the polynomials used in the discrete approximation of the state.  $\mathbf{1}$  denotes a vector whose entries are all equal to one, while  $\mathbf{0}$  is a vector whose entries are all equal to zero, their dimension should be clear from context. If  $\mathbf{D}$  is the differentiation matrix introduced in (1.10), the  $\mathbf{D}_j$  is the  $j$ -th column of  $\mathbf{D}$  and  $\mathbf{D}_{i:j}$  is the submatrix formed by columns  $i$  through  $j$ .

**2. Abstract setting.** Given  $\mathbf{x} \in \mathcal{P}_N^n$  and  $\mathbf{u} \in \mathbb{R}^{mN}$  that are feasible in (1.2), define  $\mathbf{X}_i = \mathbf{x}(\tau_i)$  and  $\mathbf{U}_i = \mathbf{u}_i$ . As shown in [12], the discrete problem (1.2) can be reformulated as the nonlinear programming problem

$$\begin{aligned} & \text{minimize} && C(\mathbf{X}_N) \\ & \text{subject to} && \sum_{j=0}^N D_{ij} \mathbf{X}_j = \mathbf{f}(\mathbf{X}_i, \mathbf{U}_i), \quad 1 \leq i \leq N, \\ & && \mathbf{X}_0 = \mathbf{x}_0. \end{aligned} \quad (2.1)$$

Also, [12] shows that the equations obtained by setting the gradient of the Lagrangian to zero are equivalent to the system of equations

$$\mathbf{\Lambda}_0 = \nabla C(\mathbf{X}_N) + \sum_{i=1}^N \omega_i \nabla_x H(\mathbf{X}_i, \mathbf{U}_i, \mathbf{\Lambda}_i), \quad (2.2)$$

$$\sum_{j=1}^N D_{ij}^\dagger \mathbf{\Lambda}_j = -\nabla_x H(\mathbf{X}_i, \mathbf{U}_i, \mathbf{\Lambda}_i), \quad 1 \leq i < N, \quad (2.3)$$

$$\sum_{j=1}^N D_{Nj}^\dagger \mathbf{\Lambda}_j = -\nabla_x H(\mathbf{X}_N, \mathbf{U}_N, \mathbf{\Lambda}_N) - \nabla C(\mathbf{X}_N)/\omega_N, \quad (2.4)$$

$$\mathbf{0} = \nabla_u H(\mathbf{X}_i, \mathbf{U}_i, \mathbf{\Lambda}_i), \quad 1 \leq i \leq N, \quad (2.5)$$

where  $\mathbf{\Lambda}_0$  is the multiplier associated with the equation  $\mathbf{X}_0 = \mathbf{x}_0$  and  $\mathbf{\Lambda}_i$  for  $i > 0$  is related to the Lagrange multiplier  $\boldsymbol{\lambda}_i$  associated with the  $i$ -th equation in the discrete dynamics by

$$\mathbf{\Lambda}_i = \boldsymbol{\lambda}_i / \omega_i. \quad (2.6)$$

The first-order optimality conditions for the nonlinear program (2.1) consist of the equations (2.2)–(2.5), and the constraints in (2.1). This system can be written as

$\mathcal{T}(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) = \mathbf{0}$  where

$$(\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_6)(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) \in \mathbb{R}^{nN} \times \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^{n(N-1)} \times \mathbb{R}^n \times \mathbb{R}^{mN}.$$

The 6 components of  $\mathcal{T}$  are defined as follows:

$$\begin{aligned} \mathcal{T}_{1i}(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) &= \left( \sum_{j=0}^N D_{ij} \mathbf{X}_j \right) - \mathbf{f}(\mathbf{X}_i, \mathbf{U}_i), \quad 1 \leq i \leq N, \\ \mathcal{T}_2(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) &= \mathbf{X}_0 - \mathbf{x}_0, \\ \mathcal{T}_3(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) &= \mathbf{\Lambda}_0 - \nabla C(\mathbf{X}_N) - \sum_{i=1}^N \omega_i \nabla_x H(\mathbf{X}_i, \mathbf{U}_i, \mathbf{\Lambda}_i), \\ \mathcal{T}_{4i}(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) &= \left( \sum_{j=1}^N D_{ij}^\dagger \mathbf{\Lambda}_j \right) + \nabla_x H(\mathbf{X}_i, \mathbf{U}_i, \mathbf{\Lambda}_i), \quad 1 \leq i < N, \\ \mathcal{T}_5(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) &= \sum_{j=1}^N D_{Nj}^\dagger \mathbf{\Lambda}_j + \nabla_x H(\mathbf{X}_N, \mathbf{U}_N, \mathbf{\Lambda}_N) + \nabla C(\mathbf{X}_N) / \omega_N, \\ \mathcal{T}_{6i}(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) &= \nabla_u H(\mathbf{X}_i, \mathbf{U}_i, \mathbf{\Lambda}_i), \quad 1 \leq i \leq N. \end{aligned}$$

The proof of Theorem 1.1 reduces to a study of solutions to  $\mathcal{T}(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) = \mathbf{0}$  in a neighborhood of  $(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)$ . Our analysis is based on [6, Proposition 3.1], which we simplify below to take into account the structure of our  $\mathcal{T}$ . Other results like this are contained in Theorem 3.1 of [4], in Proposition 5.1 of [14], and in Theorem 2.1 of [15].

**PROPOSITION 2.1.** *Let  $\mathcal{X}$  be a Banach space and  $\mathcal{Y}$  be a linear normed space with the norms in both spaces denoted  $\|\cdot\|$ . Let  $\mathcal{T}: \mathcal{X} \mapsto \mathcal{Y}$  with  $\mathcal{T}$  continuously Fréchet differentiable in  $\mathcal{B}_r(\boldsymbol{\theta}^*)$  for some  $\boldsymbol{\theta}^* \in \mathcal{X}$  and  $r > 0$ . Suppose that*

$$\|\nabla \mathcal{T}(\boldsymbol{\theta}) - \nabla \mathcal{T}(\boldsymbol{\theta}^*)\| \leq \varepsilon \text{ for all } \boldsymbol{\theta} \in \mathcal{B}_r(\boldsymbol{\theta}^*)$$

*and  $\nabla \mathcal{T}(\boldsymbol{\theta}^*)$  is invertible; and define  $\mu := \|\nabla \mathcal{T}(\boldsymbol{\theta}^*)^{-1}\|$ . If  $\varepsilon\mu < 1$  and  $\|\mathcal{T}(\boldsymbol{\theta}^*)\| \leq (1 - \mu\varepsilon)r/\mu$ , then there exists a unique  $\boldsymbol{\theta} \in \mathcal{B}_r(\boldsymbol{\theta}^*)$  such that  $\mathcal{T}(\boldsymbol{\theta}) = \mathbf{0}$ . Moreover, we have the estimate*

$$\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq \frac{\mu}{1 - \mu\varepsilon} \|\mathcal{T}(\boldsymbol{\theta}^*)\|. \quad (2.7)$$

We apply Proposition 2.1 with  $\boldsymbol{\theta}^* = (\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)$  and  $\boldsymbol{\theta} = (\mathbf{X}^N, \mathbf{U}^N, \mathbf{\Lambda}^N)$ . In our context,  $\boldsymbol{\theta}^*$  and  $\boldsymbol{\theta} \in \mathcal{X} = \mathbb{R}^{n(N+1)} \times \mathbb{R}^{mN} \times \mathbb{R}^{n(N+1)}$ , where the norm on  $\mathcal{X}$  is

$$\|\boldsymbol{\theta}\| = \|(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda})\|_\infty = \max\{\|\mathbf{X}\|_\infty, \|\mathbf{U}\|_\infty, \|\mathbf{\Lambda}\|_\infty\}. \quad (2.8)$$

For this norm, the left side of (1.12) and the left side of (2.7) are the same. The key steps in the analysis are the estimation of the residual  $\|\mathcal{T}(\boldsymbol{\theta}^*)\|$ , the proof that  $\nabla \mathcal{T}(\boldsymbol{\theta}^*)$  is invertible, and the derivation of a bound for  $\|\nabla \mathcal{T}(\boldsymbol{\theta}^*)^{-1}\|$  that is independent of  $N$ . The norm on  $\mathcal{Y}$  enters into the estimation of both the residual  $\|\mathcal{T}(\boldsymbol{\theta}^*)\|$  in (2.7) and the parameter  $\mu := \|\nabla \mathcal{T}(\boldsymbol{\theta}^*)^{-1}\|$ . In our context, we think of an element of  $\mathcal{Y}$  as a vector with components  $\mathbf{y}_i$ ,  $1 \leq i \leq 3N + 2$ , where  $\mathbf{y}_i \in \mathbb{R}^n$  for  $1 \leq i \leq 2N + 2$  and  $\mathbf{y}_i \in \mathbb{R}^m$  for  $i > 2N + 2$ . For example,  $\mathcal{T}_1(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) \in \mathbb{R}^{nN}$  corresponds to the components  $\mathbf{y}_i \in \mathbb{R}^n$ ,  $1 \leq i \leq N$ . For the norm in  $\mathcal{Y}$ , we take

$$\|\mathbf{y}\|_\infty = \sup\{|\mathbf{y}_i| : 1 \leq i \leq 3N + 2\}. \quad (2.9)$$

**3. Analysis of the residual.** We now establish a bound for  $\mathcal{T}(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)$ , the residual which appears on the right side of the error bound (2.7). This bound for the residual ultimately appears in the right side of the error estimate (1.12).

LEMMA 3.1. *If (A1) holds, then there exists a constant  $c$ , independent of  $N$ , such that*

$$\|\mathcal{T}(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)\|_\infty \leq cN^{2-\eta} \quad (3.1)$$

for all  $N > \eta + 1$ .

*Proof.* By the definition of  $\mathcal{T}$ , we have

$$\mathcal{T}_2(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*) = \mathbf{X}_0^* - \mathbf{x}_0 = \mathbf{x}^*(\tau_0) - \mathbf{x}_0 = \mathbf{x}^*(-1) - \mathbf{x}_0 = \mathbf{0}$$

since  $\mathbf{x}^*$  satisfies the initial condition in (1.1). Likewise,  $\mathcal{T}_6(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*) = \mathbf{0}$  since (1.5) holds for all  $t \in [-1, +1]$ , which implies that (1.5) holds at the collocation points.

Now consider  $\mathcal{T}_1$ . Since  $\mathbf{D}$  is a differentiation matrix associated with the collocation points, we have

$$\sum_{j=0}^N D_{ij} \mathbf{X}_j^* = \dot{\mathbf{x}}^I(\tau_i), \quad 1 \leq i \leq N, \quad (3.2)$$

where  $\mathbf{x}^I \in \mathcal{P}_N^n$  is the (interpolating) polynomial that passes through  $\mathbf{x}^*(\tau_j)$  for  $0 \leq j \leq N$ . Since  $\mathbf{x}^*$  satisfies the dynamics of (1.1),

$$\mathbf{f}(\mathbf{X}_i^*, \mathbf{U}_i^*) = \mathbf{f}(\mathbf{x}^*(\tau_i), \mathbf{u}^*(\tau_i)) = \dot{\mathbf{x}}^*(\tau_i). \quad (3.3)$$

Combine (3.2) and (3.3) to obtain

$$\mathcal{T}_{1i}(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*) = \dot{\mathbf{x}}^I(\tau_i) - \dot{\mathbf{x}}^*(\tau_i). \quad (3.4)$$

Proposition 2.1 and Lemma 2.2 in [17] yield

$$\|\dot{\mathbf{x}}^I - \dot{\mathbf{x}}^*\|_\infty \leq \left( \frac{6e}{N-1} \right)^\eta [(1 + 2N^2) + 6eN(1 + c_1 \log N)] \left( \frac{12\|x^{(\eta+1)}\|}{\eta + 1} \right)$$

for all  $N > \eta + 1$ , where  $\mathbf{x}^{(\eta+1)}$  is the  $(\eta+1)$ -st derivative of  $\mathbf{x}$  and  $c_1 \log N$  is a bound for the Lebesgue constant of the point set  $\tau_j$ ,  $0 \leq j \leq N$ , given in Theorem 2.1 of [25]. Hence, there exists a constant  $c_2$ , independent of  $N$  but dependent on  $\eta$ , such that

$$|\mathcal{T}_{1i}(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)| = |\dot{\mathbf{x}}^I(\tau_i) - \dot{\mathbf{x}}^*(\tau_i)| \leq \|\dot{\mathbf{x}}^I - \dot{\mathbf{x}}^*\|_\infty \leq c_2 N^{2-\eta}, \quad (3.5)$$

which complies with the bound (3.1).

Next, let us consider  $\mathcal{T}_4$ . By (1.11) if  $\boldsymbol{\lambda}^I \in \mathcal{P}_{N-1}^n$  is the (interpolating) polynomial that passes through  $\boldsymbol{\lambda}^*(\tau_j)$  for  $1 \leq j \leq N$ , we have

$$\sum_{j=1}^N D_{ij}^\dagger \boldsymbol{\Lambda}_j^* = \dot{\boldsymbol{\lambda}}^I(\tau_i), \quad 1 \leq i < N, \quad \sum_{j=1}^N D_{Nj}^\dagger \boldsymbol{\Lambda}_j^* = \dot{\boldsymbol{\lambda}}^I(1) - \boldsymbol{\lambda}^I(1)/\omega_N. \quad (3.6)$$

Since  $\boldsymbol{\lambda}^*$  satisfies (1.4), it follows that

$$\nabla_x H(\mathbf{X}_i^*, \mathbf{U}_i^*, \boldsymbol{\Lambda}_i^*) = \nabla_x H(\mathbf{x}^*(\tau_i), \mathbf{u}^*(\tau_i), \boldsymbol{\lambda}^*(\tau_i)) = -\dot{\boldsymbol{\lambda}}^*(\tau_i). \quad (3.7)$$



Hence, for  $i < N$ , we have

$$\mathcal{T}_{4i}(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*) = \dot{\lambda}^I(\tau_i) - \dot{\lambda}^*(\tau_i). \quad (3.8)$$

By a similar analysis as that used for  $\mathcal{T}_1$  in (3.5), we conclude that

$$|\mathcal{T}_{4i}(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)| \leq cN^{2-\eta} \quad \text{for } i < N. \quad (3.9)$$

The difference between the analysis of the state in (3.5) and the analysis of the costate in (3.9) is that the  $O(\log N)$  bound for the Lebesgue constant of  $\tau_0, \dots, \tau_N$  must be replaced by an  $O(\sqrt{N})$  bound, derived in Theorem 5.1 of [17], for the Lebesgue constant of  $\tau_1, \dots, \tau_N$ . This difference between the state and the costate arises since the state interpolant  $\mathbf{x}^I \in \mathcal{P}_N^n$  interpolates  $\mathbf{x}^*$  at  $\tau_0, \dots, \tau_N$  while the costate interpolant  $\boldsymbol{\lambda}^I \in \mathcal{P}_{N-1}^n$  interpolates  $\boldsymbol{\lambda}^*$  at  $\tau_1, \dots, \tau_N$ . Since the Lebesgue constant term in the bound for  $\|\mathcal{T}_{4i}(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)\|$  is dominated by the other terms, the right sides of (3.5) and (3.9) are the same.

Similarly, using (3.7) and the second equation in (3.6), it follows that

$$\mathcal{T}_5(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*) = \dot{\lambda}^I(1) - \dot{\lambda}^*(1) + (\nabla C(\mathbf{x}^*(1)) - \boldsymbol{\lambda}^*(1))/\omega_N = \dot{\lambda}^I(1) - \dot{\lambda}^*(1)$$

since  $\boldsymbol{\lambda}^I(1) = \boldsymbol{\lambda}^*(1) = \nabla C(\mathbf{x}^*(1))$  by (1.4). Hence, just like the bound for  $\mathcal{T}_{4i}$  in (3.9), we have

$$|\mathcal{T}_5(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)| \leq cN^{2-\eta}.$$

Now consider  $\mathcal{T}_3$ . By (3.7), the definition  $\boldsymbol{\Lambda}_0^* = \boldsymbol{\lambda}^*(-1)$ , and the terminal condition  $\boldsymbol{\lambda}^*(1) = \nabla C(\mathbf{x}^*(1))$  from (1.4), we have

$$\mathcal{T}_3(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*) = \boldsymbol{\lambda}^*(-1) - \boldsymbol{\lambda}^*(1) + \sum_{i=1}^N \omega_i \dot{\lambda}^*(\tau_i). \quad (3.10)$$

By the fundamental theorem of calculus and the fact that  $N$ -point Radau quadrature is exact for polynomials of degree up to  $2N - 2$ , we have

$$\mathbf{0} = \boldsymbol{\lambda}^I(-1) - \boldsymbol{\lambda}^I(1) + \int_{-1}^1 \dot{\lambda}^I(t) dt = \boldsymbol{\lambda}^I(-1) - \boldsymbol{\lambda}^I(1) + \sum_{j=1}^N \omega_j \dot{\lambda}^I(\tau_j). \quad (3.11)$$

Subtract (3.11) from (3.10) to obtain

$$\mathcal{T}_5(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*) = \boldsymbol{\lambda}^*(-1) - \boldsymbol{\lambda}^I(-1) + \sum_{j=1}^N \omega_j \left( \dot{\lambda}^*(\tau_j) - \dot{\lambda}^I(\tau_j) \right) \quad (3.12)$$

since  $\boldsymbol{\lambda}^I(1) = \boldsymbol{\lambda}^*(1)$ . Since  $\omega_i > 0$  and their sum is 2, it follows from (3.8) and (3.9) that

$$\sum_{j=1}^N \omega_j \left| \dot{\lambda}^I(\tau_j) - \dot{\lambda}^*(\tau_j) \right| \leq cN^{2-\eta}. \quad (3.13)$$

By Theorem 15.1 in [24] and Lemma 2.2 and Theorem 5.1 in [17], we have

$$\begin{aligned} |\boldsymbol{\lambda}^*(-1) - \boldsymbol{\lambda}^I(-1)| &\leq (1 + c_1\sqrt{N}) \left( \frac{12}{\eta + 2} \right) \left( \frac{6e}{N} \right)^{\eta+1} \|\boldsymbol{\lambda}^{*(\eta+1)}\|_{\infty} \\ &\leq cN^{-(0.5+\eta)}. \end{aligned} \quad (3.14)$$

We combine (3.12)–(3.14) to see that  $\mathcal{T}_5(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)$  also complies with the bound (3.1). This completes the proof.  $\square$

**4. Invertibility.** In this section, we show that the derivative  $\nabla\mathcal{T}(\boldsymbol{\theta}^*)$  is invertible. This is equivalent to showing that for each  $\mathbf{y} \in \mathcal{Y}$ , there is a unique  $\boldsymbol{\theta} \in \mathcal{X}$  such that  $\nabla\mathcal{T}(\boldsymbol{\theta}^*)[\boldsymbol{\theta}] = \mathbf{y}$ . In our case,  $\boldsymbol{\theta}^* = (\mathbf{X}^*, \mathbf{U}^*, \boldsymbol{\Lambda}^*)$  and  $\boldsymbol{\theta} = (\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda})$ . To simplify the notation, we let  $\nabla\mathcal{T}^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}]$  denote the derivative of  $\mathcal{T}$  evaluated at  $(\mathbf{X}^*, \mathbf{U}^*, \boldsymbol{\Lambda}^*)$  operating on  $(\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda})$ . This derivative involves the following 6 matrices:

$$\begin{aligned} \mathbf{A}_i &= \nabla_x \mathbf{f}(\mathbf{x}^*(\tau_i), \mathbf{u}^*(\tau_i)), & \mathbf{B}_i &= \nabla_u \mathbf{f}(\mathbf{x}^*(\tau_i), \mathbf{u}^*(\tau_i)), \\ \mathbf{Q}_i &= \nabla_{xx} H(\mathbf{x}^*(\tau_i), \mathbf{u}^*(\tau_i), \boldsymbol{\lambda}^*(\tau_i)), & \mathbf{S}_i &= \nabla_{xu} H(\mathbf{x}^*(\tau_i), \mathbf{u}^*(\tau_i), \boldsymbol{\lambda}^*(\tau_i)), \\ \mathbf{R}_i &= \nabla_{uu} H(\mathbf{x}^*(\tau_i), \mathbf{u}^*(\tau_i), \boldsymbol{\lambda}^*(\tau_i)), & \mathbf{T} &= \nabla^2 C(\mathbf{x}^*(1)). \end{aligned}$$

With this notation, the 6 components of  $\nabla\mathcal{T}^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}]$  are as follows:

$$\begin{aligned} \nabla\mathcal{T}_{1i}^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}] &= \left( \sum_{j=0}^N D_{ij} \mathbf{X}_j \right) - \mathbf{A}_i \mathbf{X}_i - \mathbf{B}_i \mathbf{U}_i, \quad 1 \leq i \leq N, \\ \nabla\mathcal{T}_2^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}] &= \mathbf{X}_0, \\ \nabla\mathcal{T}_3^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}] &= \boldsymbol{\Lambda}_0 - \mathbf{T} \mathbf{X}_N - \sum_{i=1}^N \omega_i (\mathbf{A}_i^\top \boldsymbol{\Lambda}_i + \mathbf{Q}_i \mathbf{X}_i + \mathbf{S}_i \mathbf{U}_i), \\ \nabla\mathcal{T}_{4i}^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}] &= \left( \sum_{j=1}^N D_{ij}^\dagger \boldsymbol{\Lambda}_j \right) + \mathbf{A}_i^\top \boldsymbol{\Lambda}_i + \mathbf{Q}_i \mathbf{X}_i + \mathbf{S}_i \mathbf{U}_i, \quad 1 \leq i < N, \\ \nabla\mathcal{T}_5^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}] &= \left( \sum_{j=1}^N D_{Nj}^\dagger \boldsymbol{\Lambda}_j \right) + \mathbf{A}_N^\top \boldsymbol{\Lambda}_N + \mathbf{Q}_N \mathbf{X}_N + \mathbf{S}_N \mathbf{U}_N + \mathbf{T} \mathbf{X}_N / \omega_N, \\ \nabla\mathcal{T}_{6i}^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}] &= \mathbf{S}_i^\top \mathbf{X}_i + \mathbf{R}_i \mathbf{U}_i + \mathbf{B}_i^\top \boldsymbol{\Lambda}_i, \quad 1 \leq i \leq N. \end{aligned}$$

The analysis of invertibility starts with results concerning the invertibility of the linearized discrete state dynamics.

**LEMMA 4.1.** *If (P1) and (A3) hold, then for each  $\mathbf{q} \in \mathbb{R}^n$  and  $\mathbf{p} \in \mathbb{R}^{nN}$  with  $\mathbf{p}_i \in \mathbb{R}^n$ , the linear system*

$$\left( \sum_{j=0}^N D_{ij} \mathbf{X}_j \right) - \mathbf{A}_i \mathbf{X}_i = \mathbf{p}_i \quad 1 \leq i \leq N, \quad (4.1)$$

$$\mathbf{X}_0 = \mathbf{q}, \quad (4.2)$$

has a unique solution  $\mathbf{X}_j \in \mathbb{R}^n$ ,  $0 \leq j \leq N$ . This solution has the bound

$$\|\mathbf{X}_j\|_\infty \leq 4\|\mathbf{p}\|_\infty + 2\|\mathbf{q}\|_\infty, \quad 0 \leq j \leq N. \quad (4.3)$$

*Proof.* If  $\mathbf{X}$  is a solution of (4.1)–(4.2), then  $\mathbf{X}_0 = \mathbf{q}$  and  $\mathbf{X}_0$  trivially satisfies (4.3). Next, focus on the remaining components of  $\mathbf{X}$ . Let  $\bar{\mathbf{X}}$  be the vector obtained by vertically stacking  $\mathbf{X}_1$  through  $\mathbf{X}_N$ , let  $\mathbf{A}$  be the block diagonal matrix with  $i$ -th diagonal block  $\mathbf{A}_i$ ,  $1 \leq i \leq N$ , and define  $\bar{\mathbf{D}} = \mathbf{D}_{1:N} \otimes \mathbf{I}_n$  and  $\bar{\mathbf{D}}_0 = \mathbf{D}_0 \otimes \mathbf{I}_n$  where  $\otimes$  is the Kronecker product. With this notation, the linear system (4.1)–(4.2) reduces to

$$(\bar{\mathbf{D}} - \mathbf{A})\bar{\mathbf{X}} = \mathbf{p} - \bar{\mathbf{D}}_0 \mathbf{q}. \quad (4.4)$$

By (P1),  $\mathbf{D}_{1:N}$  is invertible which implies that  $\bar{\mathbf{D}}$  is invertible and  $\bar{\mathbf{D}}^{-1} = \mathbf{D}_{1:N}^{-1} \otimes \mathbf{I}_n$ . Since the polynomial that is identically equal to 1 has derivative 0 and since  $\mathbf{D}$  is a differentiation matrix, we have  $\mathbf{D}\mathbf{1} = \mathbf{0}$ , which implies that  $\mathbf{D}_{1:N}^{-1}\mathbf{D}_0 = -\mathbf{1}$ . It follows that

$$\bar{\mathbf{D}}^{-1}\bar{\mathbf{D}}_0 = [\mathbf{D}_{1:N}^{-1} \otimes \mathbf{I}_n][\mathbf{D}_0 \otimes \mathbf{I}_n] = -\mathbf{1} \otimes \mathbf{I}_n.$$

Multiply (4.4) by  $\bar{\mathbf{D}}^{-1}$  to obtain

$$(\mathbf{I} - \bar{\mathbf{D}}^{-1}\mathbf{A})\bar{\mathbf{X}} = \bar{\mathbf{D}}^{-1}\mathbf{p} + (\mathbf{1} \otimes \mathbf{I}_n)\mathbf{q}. \quad (4.5)$$

By (P1)  $\|\mathbf{D}_{1:N}^{-1}\|_\infty \leq 2$ , which implies that  $\|\bar{\mathbf{D}}^{-1}\|_\infty \leq 2$ . By (A3)  $\|\mathbf{A}\|_\infty \leq 1/4$ . By [18, p. 351],  $\mathbf{I} - \bar{\mathbf{D}}^{-1}\mathbf{A}$  is invertible and  $\|(\mathbf{I} - \bar{\mathbf{D}}^{-1}\mathbf{A})^{-1}\|_\infty \leq 2$ . Multiply (4.5) by  $(\mathbf{I} - \bar{\mathbf{D}}^{-1}\mathbf{A})^{-1}$  and take the norm of each side to obtain  $\|\bar{\mathbf{X}}\|_\infty \leq 4\|\mathbf{p}\|_\infty + 2\|\mathbf{q}\|_\infty$ . This complete the proof of (4.3).  $\square$

Next, we establish the invertibility of  $\nabla\mathcal{T}^*$ .

PROPOSITION 4.2. *If (P1), (A2), and (A3) hold, then  $\nabla\mathcal{T}^*$  is invertible.*

*Proof.* Our approach is to formulate a strongly convex quadratic programming problem which has a unique solution  $(\mathbf{X}, \mathbf{U})$  by (A2), and which has the property that the associated first-order optimality condition is  $\nabla\mathcal{T}^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$ . Since  $\nabla\mathcal{T}^*$  is square and  $\nabla\mathcal{T}^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$  has a solution for each choice of  $\mathbf{y}$ , we conclude the  $\nabla\mathcal{T}^*$  is invertible.

The quadratic program is

$$\left. \begin{aligned} & \text{minimize} && \frac{1}{2}\mathcal{Q}(\mathbf{X}, \mathbf{U}) + \mathcal{L}(\mathbf{X}, \mathbf{U}) \\ & \text{subject to} && \sum_{j=0}^N D_{ij}\mathbf{X}_j = \mathbf{A}_i\mathbf{X}_i + \mathbf{B}_i\mathbf{U}_i + \mathbf{y}_{1i}, \quad 1 \leq i \leq N, \\ & && \mathbf{X}_0 = \mathbf{y}_2, \end{aligned} \right\} \quad (4.6)$$

where the quadratic and linear terms in the objective are

$$\mathcal{Q}(\mathbf{X}, \mathbf{U}) = \mathbf{X}_N^\top \mathbf{T} \mathbf{X}_N + \sum_{i=1}^N \omega_i (\mathbf{X}_i^\top \mathbf{Q}_i \mathbf{X}_i + 2\mathbf{X}_i^\top \mathbf{S}_i \mathbf{U}_i + \mathbf{U}_i^\top \mathbf{R}_i \mathbf{U}_i), \quad (4.7)$$

$$\mathcal{L}(\mathbf{X}, \mathbf{U}) = \mathbf{X}_0^\top \left( \mathbf{y}_3 + \sum_{i=1}^N \omega_i \mathbf{y}_{4i} \right) - \sum_{i=1}^N \omega_i (\mathbf{y}_{4i}^\top \mathbf{X}_i + \mathbf{y}_{6i}^\top \mathbf{U}_i). \quad (4.8)$$

In (4.8) we simplified the formula for  $\mathcal{L}$  by introducing  $\mathbf{y}_{4N} = \mathbf{y}_5$ . By Lemma 4.1, the quadratic programming problem (4.6) is feasible. Since the Radau quadrature weights  $\omega_i$  are strictly positive, it follows from (A2) that  $\mathcal{Q}$  is strongly convex. Hence, there exists a unique optimal solution to (4.6) for any choice of  $\mathbf{y}$ . Since the constraints are linear, the first-order optimality conditions hold. The linear term was chosen so that the first-order optimality conditions for (4.6) reduce to  $\nabla\mathcal{T}^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$ . Hence, the existence of a solution to (4.6) for each  $\mathbf{y}$  would imply that  $\nabla\mathcal{T}^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$  has a solution for each choice of  $\mathbf{y}$ . To complete the proof, we need to show that the first-order optimality conditions for (4.6) are equivalent to  $\nabla\mathcal{T}^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$ .

The Lagrangian of (4.6) is given by

$$\frac{1}{2}\mathcal{Q}(\mathbf{X}, \mathbf{U}) + \mathcal{L}(\mathbf{X}, \mathbf{U}) + \sum_{i=1}^N \lambda_i^\top \left( \mathbf{A}_i\mathbf{X}_i + \mathbf{B}_i\mathbf{U}_i + \mathbf{y}_{1i} - \sum_{j=0}^N D_{ij}\mathbf{X}_j \right) + \mathbf{\Lambda}_0^\top (\mathbf{y}_2 - \mathbf{X}_0).$$

The first-order optimality conditions are obtained by setting to zero the derivative of the Lagrangian with respect to each of the components of  $\mathbf{X}$  and  $\mathbf{U}$ . We give the derivation of the 3rd, 4th, and 5th components of  $\nabla T^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$ , while the 6th component follows in a similar fashion, and the 1st and 2nd components are simply the constraints in (4.6).

Setting to zero the partial derivative of the Lagrangian with respect to  $\mathbf{X}_i$ ,  $1 \leq i < N$ , yields the equation

$$\mathbf{A}_i^\top \boldsymbol{\lambda}_i + \omega_i \mathbf{Q}_i \mathbf{X}_i + \omega_i \mathbf{S}_i \mathbf{U}_i - \sum_{j=1}^N D_{ji} \boldsymbol{\lambda}_j = \omega_i \mathbf{y}_{4i}.$$

Substituting  $D_{ji} = -(\omega_i/\omega_j)D_{ij}^\ddagger$  and  $\boldsymbol{\lambda}_j = \omega_j \mathbf{\Lambda}_j$ , we obtain

$$\left( \sum_{j=1}^N D_{ij}^\ddagger \mathbf{\Lambda}_j \right) + \mathbf{A}_i^\top \mathbf{\Lambda}_i + \mathbf{Q}_i \mathbf{X}_i + \mathbf{S}_i \mathbf{U}_i = \mathbf{y}_{4i}, \quad (4.9)$$

which gives the 4th component of  $\nabla T^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$ . In a similar fashion, setting to zero the partial derivative of the Lagrangian with respect to  $\mathbf{X}_N$  yields the 5th component of  $\nabla T^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$ :

$$\left( \sum_{j=1}^N D_{Nj}^\ddagger \mathbf{\Lambda}_j \right) + \mathbf{A}_N^\top \mathbf{\Lambda}_N + \mathbf{Q}_N \mathbf{X}_N + \mathbf{S}_N \mathbf{U}_N + \mathbf{T} \mathbf{X}_N / \omega_N = \mathbf{y}_{4N}. \quad (4.10)$$

Setting to zero the partial derivative of the Lagrangian with respect to  $\mathbf{X}_0$  gives the equation

$$\mathbf{\Lambda}_0 + \sum_{i=1}^N (\boldsymbol{\lambda}_i D_{i0} - \omega_i \mathbf{y}_{4i}) = \mathbf{y}_3. \quad (4.11)$$

Since  $\mathbf{D}$  is a differentiation matrix and  $\mathbf{D}\mathbf{1} = \mathbf{0}$ , it follows that

$$D_{i0} = - \sum_{j=1}^N D_{ij}.$$

Consequently, we have

$$\sum_{i=1}^N D_{i0} \boldsymbol{\lambda}_i = - \sum_{i=1}^N \sum_{j=1}^N D_{ij} \boldsymbol{\lambda}_i = \sum_{i=1}^N \sum_{j=1}^N \omega_j D_{ji}^\ddagger \mathbf{\Lambda}_i = \sum_{i=1}^N \sum_{j=1}^N \omega_i D_{ij}^\ddagger \mathbf{\Lambda}_j.$$

We make this substitution as well as (4.9) and (4.10) into (4.11). The  $\mathbf{D}^\ddagger$  terms cancel to give

$$\mathbf{\Lambda}_0 - \mathbf{T} \mathbf{X}_N - \sum_{i=1}^N \omega_i (\mathbf{A}_i^\top \mathbf{\Lambda}_i + \mathbf{Q}_i \mathbf{X}_i + \mathbf{S}_i \mathbf{U}_i) = \mathbf{y}_3,$$

which is the 3rd component of  $\nabla T^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$ . This completes the proof.  $\square$

**5.  $\omega$ -norm bounds for the state and control.** In this section we obtain a bound for the  $(\mathbf{X}, \mathbf{U})$  component of the solution to  $\nabla \mathcal{T}^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}] = \mathbf{y}$  in terms of  $\mathbf{y}$ . Since  $\mathbf{X}_0$  must satisfy the constraint  $\mathbf{X}_0 = \mathbf{y}_2$ , it is trivially bounded in terms of  $\|\mathbf{y}\|_\infty$ . Hence, we focus on  $(\mathbf{X}_i, \mathbf{U}_i)$ ,  $1 \leq i \leq N$ . The bound we derive in this section is in terms of the  $\omega$ -norms defined by

$$\|\mathbf{X}\|_\omega^2 = |\mathbf{X}_N|^2 + \sum_{i=1}^N \omega_i |\mathbf{X}_i|^2 \quad \text{and} \quad \|\mathbf{U}\|_\omega^2 = \sum_{i=1}^N \omega_i |\mathbf{U}_i|^2. \quad (5.1)$$

This defines a norm since the Radau quadrature weight  $\omega_i > 0$  for each  $i$ . Since the  $(\mathbf{X}, \mathbf{U})$  component of the solution to  $\nabla \mathcal{T}^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}] = \mathbf{y}$  is a solution of the quadratic program (4.6), we will bound the solution to the quadratic program.

First, let us think more abstractly. Let  $\pi$  be a symmetric, continuous bilinear functional defined on a Hilbert space  $\mathcal{H}$ , let  $\ell$  be a continuous linear functional, let  $\phi \in \mathcal{H}$ , and consider the quadratic program

$$\min \left\{ \frac{1}{2} \pi(v + \phi, v + \phi) + \ell(v + \phi) : v \in \mathcal{V} \right\},$$

where  $\mathcal{V}$  is a subspace of  $\mathcal{H}$ . If  $w$  is a minimizer, then by the first-order optimality conditions, we have

$$\pi(w, v) + \pi(\phi, v) + \ell(v) = 0 \quad \text{for all } v \in \mathcal{V}.$$

Inserting  $v = w$  yields

$$\pi(w, w) = -(\pi(w, \phi) + \ell(w)). \quad (5.2)$$

We apply this observation to the quadratic program (4.6) where we treat  $\mathbf{X}_0 = \mathbf{y}_2$  as fixed, so the minimization is over  $(\mathbf{X}_i, \mathbf{U}_i)$ ,  $1 \leq i \leq N$ . We identify  $\ell$  with the linear functional  $\mathcal{L}$  in (4.8) but with the  $\mathbf{X}_0$  term dropped since it is fixed, and  $\pi$  with the bilinear form associated with the quadratic term (4.7). The subspace  $\mathcal{V}$  is the null space of the linear operator in (4.6) and  $\phi$  is a particular solution of the linear system. The complete solution of (4.6) is the particular solution plus the minimizer over the null space.

In more detail, let  $\boldsymbol{\chi}$  denote the solution to (4.1)–(4.2) given by Lemma 4.1 for  $\mathbf{p} = \mathbf{y}_1$  and  $\mathbf{q} = \mathbf{y}_2$ . We consider the particular solution  $(\mathbf{X}, \mathbf{U})$  of the linear system in (4.6) given by  $(\boldsymbol{\chi}, \mathbf{0})$ . The relation (5.2) describing the null space component  $(\mathbf{X}, \mathbf{U})$  of the solution is

$$\mathcal{Q}(\mathbf{X}, \mathbf{U}) = - \left( \boldsymbol{\chi}_N^\top \mathbf{T} \mathbf{X}_N + \sum_{i=1}^N \omega_i [(\mathbf{Q}_i \boldsymbol{\chi}_i - \mathbf{y}_{4i})^\top \mathbf{X}_i - \mathbf{y}_{6i}^\top \mathbf{U}_i] \right). \quad (5.3)$$

Here the terms containing  $\boldsymbol{\chi}$  are associated with  $\pi(w, \phi)$ , while the remaining terms are associated with  $\ell$ , or equivalently with  $\mathcal{L}$ . By (A2) we have the lower bound

$$\mathcal{Q}(\mathbf{X}, \mathbf{U}) \geq \alpha (\|\mathbf{X}\|_\omega^2 + \|\mathbf{U}\|_\omega^2). \quad (5.4)$$

All the terms on the right side of (5.3) can be bounded with the Schwarz inequality; for example,

$$\begin{aligned} \sum_{i=1}^N \omega_i \mathbf{y}_{4i}^\top \mathbf{X}_i &\leq \left( \sum_{i=1}^N \omega_i |\mathbf{y}_{4i}|^2 \right)^{1/2} \left( \sum_{i=1}^N \omega_i |\mathbf{X}_i|^2 \right)^{1/2} \\ &\leq \sqrt{2} \|\mathbf{y}_4\|_\infty (\|\mathbf{X}\|_\omega^2 + \|\mathbf{U}\|_\omega^2)^{1/2}. \end{aligned} \quad (5.5)$$

The last inequality exploits the fact that the  $\omega_i$  sum to 2 and  $|\mathbf{y}_{4i}| \leq \|\mathbf{y}_4\|_\infty$ . To handle the terms involving  $\boldsymbol{\chi}$  in (5.3), we utilize the upper bound  $\|\boldsymbol{\chi}_j\|_\infty \leq 6\|\mathbf{y}\|_\infty$  based on Lemma 4.1 with  $\mathbf{p} = \mathbf{y}_1$  and  $\mathbf{q} = \mathbf{y}_2$ . Combining upper bounds of the form (5.5) with the lower bound (5.4), we conclude from (5.3) that both  $\|\mathbf{X}\|_\omega$  and  $\|\mathbf{U}\|_\omega$  are bounded by a constant times  $\|\mathbf{y}\|_\infty$ . The complete solution of (4.6) is the null space component that we just estimated plus the particular solution  $(\boldsymbol{\chi}, \mathbf{0})$ . Again, since  $\|\boldsymbol{\chi}_j\|_\infty \leq 6\|\mathbf{y}\|_\infty$ , we obtain the following result.

LEMMA 5.1. *If (A2)–(A3) and (P1) hold, then there exists a constant  $c$ , independent of  $N$ , such that the solution  $(\mathbf{X}, \mathbf{U})$  of (4.6) satisfies  $\|\mathbf{X}\|_\omega \leq c\|\mathbf{y}\|_\infty$  and  $\|\mathbf{U}\|_\omega \leq c\|\mathbf{y}\|_\infty$ .*

**6.  $\infty$ -norm bounds.** We now need to convert these  $\omega$ -norm bounds for  $\mathbf{X}$  and  $\mathbf{U}$  into  $\infty$ -norm bounds and at the same time, obtain an  $\infty$ -norm bound for  $\boldsymbol{\Lambda}$ . As in Lemma 4.1, the solution to the dynamics in (4.6) can be expressed

$$\bar{\mathbf{X}} = (\mathbf{I} - \bar{\mathbf{D}}^{-1}\mathbf{A})^{-1} [\bar{\mathbf{D}}^{-1}\mathbf{B}\mathbf{U} + \bar{\mathbf{D}}^{-1}\mathbf{y}_1 + (\mathbf{1} \otimes \mathbf{I}_n)\mathbf{y}_2], \quad (6.1)$$

where  $\mathbf{B}$  is the block diagonal matrix with  $i$ -th diagonal block  $\mathbf{B}_i$  and  $\mathbf{U}$  is obtained by vertically stacking  $\mathbf{U}_1$  through  $\mathbf{U}_N$ . By Lemma 4.1,

$$\|(\mathbf{I} - \bar{\mathbf{D}}^{-1}\mathbf{A})^{-1} [\bar{\mathbf{D}}^{-1}\mathbf{y}_1 + (\mathbf{1} \otimes \mathbf{I}_n)\mathbf{y}_2]\| \leq 4\|\mathbf{y}_1\|_\infty + 2\|\mathbf{y}_2\|_\infty. \quad (6.2)$$

The term  $\bar{\mathbf{D}}^{-1}\mathbf{B}\mathbf{U}$  can be bounded using (P2) and the strategy given in Section 6 of [16]. That is, we first observe that

$$\bar{\mathbf{D}}^{-1}\mathbf{B}\mathbf{U} = [\mathbf{D}_{1:N}^{-1} \otimes \mathbf{I}_n]\mathbf{B}\mathbf{U} = [(\mathbf{W}^{1/2}\mathbf{D}_{1:N})^{-1} \otimes \mathbf{I}_n]\mathbf{B}\mathbf{U}_\omega, \quad (6.3)$$

where  $\mathbf{W}$  is the diagonal matrix with the quadrature weights on the diagonal and  $\mathbf{U}_\omega$  is the vector whose  $i$ -th element is  $\sqrt{\omega_i}\mathbf{U}_i$ ; the  $\sqrt{\omega_i}$  factors in (6.3) cancel each other. An element of the vector  $\bar{\mathbf{D}}^{-1}\mathbf{B}\mathbf{U}$  is the dot product between a row of  $(\mathbf{W}^{1/2}\mathbf{D}_{1:N})^{-1} \otimes \mathbf{I}_n$  and the column vector  $\mathbf{B}\mathbf{U}_\omega$ . By (P2) the rows of  $(\mathbf{W}^{1/2}\mathbf{D}_{1:N})^{-1} \otimes \mathbf{I}_n$  have Euclidean length bounded by  $\sqrt{2}$ . By the properties of matrix norms induced by vector norms, we have

$$\|\mathbf{B}\mathbf{U}_\omega\|_2 \leq \|\mathbf{B}\|_2\|\mathbf{U}_\omega\|_2 = \|\mathbf{B}\|_2\|\mathbf{U}\|_\omega.$$

It follows that

$$\|\bar{\mathbf{D}}^{-1}\mathbf{B}\mathbf{U}\|_\infty \leq \sqrt{2}\|\mathbf{B}\|_2\|\mathbf{U}\|_\omega \leq c\|\mathbf{y}\|_\infty,$$

where the generic constant  $c$  is independent of  $N$  by Lemma 5.1. Combining this with (6.2), we conclude that for some constant  $c$ , the  $\bar{\mathbf{X}}$  component of the solution to (4.6) satisfies  $\|\bar{\mathbf{X}}\|_\infty \leq c\|\mathbf{y}\|_\infty$ . Since  $|\mathbf{X}_0| = |\mathbf{y}_2| \leq \|\mathbf{y}\|_\infty$ , the entire  $\mathbf{X}$ -component of the solution to (4.6) satisfies  $\|\mathbf{X}\|_\infty \leq c\|\mathbf{y}\|_\infty$  for some  $c$ .

Next, we focus on the 4th and 5th components of  $\nabla\mathcal{T}^*[\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}] = \mathbf{y}$  which can be written

$$\bar{\mathbf{D}}^\dagger\bar{\boldsymbol{\Lambda}} + \mathbf{A}^\top\bar{\boldsymbol{\Lambda}} + \mathbf{Q}\bar{\mathbf{X}} + \mathbf{S}\mathbf{U} + \left(\frac{1}{\omega_N}\right)(\mathbf{e}_N \otimes \mathbf{I}_n)\mathbf{T}\mathbf{X}_N = \mathbf{y}_4, \quad (6.4)$$

where  $\bar{\mathbf{D}}^\dagger = \mathbf{D}^\dagger \otimes \mathbf{I}_n$ ,  $\bar{\boldsymbol{\Lambda}}$  is obtained by vertically stacking  $\boldsymbol{\Lambda}_1$  through  $\boldsymbol{\Lambda}_N$ ,  $\mathbf{Q}$  and  $\mathbf{S}$  are block diagonal matrices with  $i$ -th diagonal blocks  $\mathbf{Q}_i$  and  $\mathbf{S}_i$  respectively, and

$\mathbf{e}_N \in \mathbb{R}^N$  is the vector whose components are all zero except for the  $N$ -th component which is 1. Similar to our manipulations of the state dynamics in (6.1), we use (A3) and (P3) to solve for  $\bar{\mathbf{\Lambda}}$ :

$$\bar{\mathbf{\Lambda}} = -(\mathbf{I} + \bar{\mathbf{D}}^\dagger^{-1} \mathbf{A}^\top)^{-1} \bar{\mathbf{D}}^\dagger^{-1} \left[ \mathbf{S}\mathbf{U} + \mathbf{Q}\bar{\mathbf{X}} + \left( \frac{1}{\omega_N} \right) (\mathbf{e}_N \otimes \mathbf{I}_n) \mathbf{T}\mathbf{X}_N - \mathbf{y}_4 \right] \quad (6.5)$$

Let  $p$  be the polynomial that is identically one, and let  $\mathbf{p}$  be the associated vector whose components are all one. Making this substitution in (1.11) gives  $\mathbf{D}^\dagger \mathbf{1} = -\mathbf{e}_N/\omega_N$ , which implies that

$$\mathbf{D}^\dagger^{-1} \mathbf{e}_N = -\omega_N \mathbf{1}. \quad (6.6)$$

Consequently, we have

$$\bar{\mathbf{D}}^\dagger^{-1} (\mathbf{e}_N \otimes \mathbf{I}_n) / \omega_N = [\mathbf{D}^\dagger^{-1} \otimes \mathbf{I}_n] (\mathbf{e}_N \otimes \mathbf{I}_n) / \omega_N = -\mathbf{1} \otimes \mathbf{I}_n.$$

With this substitution, (6.5) becomes

$$\bar{\mathbf{\Lambda}} = -(\mathbf{I} + \bar{\mathbf{D}}^\dagger^{-1} \mathbf{A}^\top)^{-1} [\bar{\mathbf{D}}^\dagger^{-1} (\mathbf{S}\mathbf{U} + \mathbf{Q}\bar{\mathbf{X}} - \mathbf{y}_4) - (\mathbf{1} \otimes \mathbf{I}_n) \mathbf{T}\mathbf{X}_N]. \quad (6.7)$$

Exactly as in Lemma 4.1, we have  $\|\bar{\mathbf{D}}^\dagger^{-1}\|_\infty \leq 2$  by (P2), and  $(\mathbf{I} + \bar{\mathbf{D}}^\dagger^{-1} \mathbf{A}^\top)^{-1} \|\infty \leq 2$  by (A3). As in the analysis of the state dynamics, all the terms on the right side of (6.7) are bounded by  $c\|\mathbf{y}\|_\infty$  for a suitable choice of  $c$ . Hence, we have  $\|\bar{\mathbf{\Lambda}}\|_\infty \leq c\|\mathbf{y}\|_\infty$ .

Now consider the 6th component of  $\nabla T^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$ , which can be written

$$\mathbf{S}_i^\top \mathbf{X}_i + \mathbf{R}_i \mathbf{U}_i + \mathbf{B}_i^\top \mathbf{\Lambda}_i = \mathbf{y}_{6i}, \quad 1 \leq i \leq N.$$

Previously, we have shown the existence of a constant  $c$ , independent of  $N$ , such that  $\|\mathbf{X}_i\|_\infty \leq c\|\mathbf{y}\|_\infty$  and  $\|\mathbf{\Lambda}_i\|_\infty \leq c\|\mathbf{y}\|_\infty$ ,  $1 \leq i \leq N$ . By (A2) the smallest eigenvalue of  $\mathbf{R}_i$  is bounded from below by  $\alpha$ . Hence, there also exists a constant  $c$  such that  $\|\mathbf{U}_i\|_\infty \leq c\|\mathbf{y}\|_\infty$ . Finally, the 3rd component of  $\nabla T^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$  can be written

$$\mathbf{\Lambda}_0 - \mathbf{T}\mathbf{X}_N - \sum_{i=1}^N \omega_i (\mathbf{A}_i^\top \mathbf{\Lambda}_i + \mathbf{Q}_i \mathbf{X}_i + \mathbf{S}_i \mathbf{U}_i) = \mathbf{y}_3.$$

By the uniform bounds on  $\|\mathbf{X}_i\|_\infty$ ,  $\|\mathbf{U}_i\|_\infty$ , and  $\|\mathbf{\Lambda}_i\|_\infty$ ,  $1 \leq i \leq N$ , there also exists a constant  $c$ , independent of  $N$ , such that  $\|\mathbf{\Lambda}_0\|_\infty \leq c\|\mathbf{y}\|_\infty$ . We summarize these results as follows:

**LEMMA 6.1.** *If (A2)–(A3) and (P1)–(P4) hold, then there exists a constant  $c$ , independent of  $N$ , such that the solution of  $\nabla T^*[\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}] = \mathbf{y}$  satisfies*

$$\|\mathbf{X}\|_\infty + \|\mathbf{U}\|_\infty + \|\mathbf{\Lambda}\|_\infty \leq c\|\mathbf{y}\|_\infty.$$

The proof of Theorem 1.1 follows exactly as in [16]. By Lemma 6.1,  $\mu = \|\nabla T(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)^{-1}\|_\infty$  is bounded uniformly in  $N$ . Choose  $\varepsilon$  small enough that  $\varepsilon\mu < 1$ . When we compute the difference  $\nabla T(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) - \nabla T(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)$  for  $(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda})$  near  $(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)$  in the  $\infty$ -norm, the  $\mathbf{D}$  and  $\mathbf{D}^\dagger$  constant terms cancel, and we are left with terms involving the difference of derivatives of  $\mathbf{f}$  or  $C$  up to second order at nearby points. By assumption, these second derivatives are uniformly continuous on the closure of  $\Omega$  and on a ball around  $\mathbf{x}^*(1)$ . Hence, for  $r$  sufficiently small, we have

$$\|\nabla T(\mathbf{X}, \mathbf{U}, \mathbf{\Lambda}) - \nabla T(\mathbf{X}^*, \mathbf{U}^*, \mathbf{\Lambda}^*)\|_\infty \leq \varepsilon$$

whenever

$$\max\{\|\mathbf{X} - \mathbf{X}^*\|_\infty, \|\mathbf{U} - \mathbf{U}^*\|_\infty, \|\boldsymbol{\Lambda} - \boldsymbol{\Lambda}^*\|_\infty\} \leq r. \quad (6.8)$$

By Lemma 3.1, it follows that  $\|\mathcal{T}(\mathbf{X}^*, \mathbf{U}^*, \boldsymbol{\Lambda}^*)\| \leq (1 - \mu\varepsilon)r/\mu$  for all  $N$  sufficiently large. Hence, by Proposition 2.1, there exists a solution to  $\mathcal{T}(\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}) = \mathbf{0}$  satisfying (6.8). Moreover, by (2.7) and (3.1), the estimate (1.12) holds. To complete the proof, we need to show that  $(\mathbf{X}, \mathbf{U})$  is a local minimizer for (2.1). After replacing the KKT multipliers by the transformed quantities given by (2.6), the Hessian of the Lagrangian is a block diagonal matrix with the following matrices forming the diagonal blocks:

$$\begin{aligned} \omega_i \nabla_{(x,u)}^2 H(\mathbf{X}_i, \mathbf{U}_i, \boldsymbol{\Lambda}_i), & \quad 1 \leq i < N, \\ \omega_i \nabla_{(x,u)}^2 H(\mathbf{X}_i, \mathbf{U}_i, \boldsymbol{\Lambda}_i) + \nabla_{(x,u)}^2 C(\mathbf{X}_i), & \quad i = N, \end{aligned}$$

where  $H$  is the Hamiltonian. In computing the Hessian, we assume that the  $\mathbf{X}$  and  $\mathbf{U}$  variables are arranged in the following order:  $\mathbf{X}_1, \mathbf{U}_1, \mathbf{X}_2, \mathbf{U}_2, \dots, \mathbf{X}_N, \mathbf{U}_N$ . By (A2) the Hessian is positive definite when evaluated at  $(\mathbf{X}^*, \mathbf{U}^*, \boldsymbol{\Lambda}^*)$ . By continuity of the second derivative of  $C$  and  $\mathbf{f}$ , and by the convergence result (1.12), we conclude that the Hessian of the Lagrangian, evaluated at the solution of  $\mathcal{T}(\mathbf{X}, \mathbf{U}, \boldsymbol{\Lambda}) = \mathbf{0}$  satisfying (6.8), is positive definite for  $N$  sufficiently large. Hence, by the second-order sufficient optimality condition [21, Thm. 12.6],  $(\mathbf{X}, \mathbf{U})$  is a strict local minimizer of (2.1). This completes the proof of Theorem 1.1.

**7. Numerical illustration.** Let us consider the unconstrained control problem previously introduced in [16]:

$$\min \left\{ -x(2) : \dot{x}(t) = \frac{5}{2}(-x(t) + x(t)u(t) - u(t)^2), x(0) = 1 \right\} \quad (7.1)$$

The optimal solution and associated costate are

$$\begin{aligned} x^*(t) &= 4/a(t), \quad a(t) = 1 + 3 \exp(2.5t), \\ u^*(t) &= x^*(t)/2, \\ \lambda^*(t) &= -a^2(t) \exp(-2.5t) / [\exp(-5) + 9 \exp(5) + 6]. \end{aligned}$$

Figure 7.1 plots the logarithm of the sup-norm error in the state, control, and costate as a function of the number of collocation points. Since these plots are nearly linear, the error behaves like  $c10^{-\alpha N}$  where  $\alpha \approx 0.6$  for either the state or the control and  $\alpha \approx 0.8$  for the costate. In Theorem 1.1, the dependence of the error on  $N$  is somewhat complex due to the connection between  $\eta$  and  $N$ . As we increase  $N$ , we can also increase  $\eta$  when the solution is infinitely differentiable, however, the norm of the derivatives also enters into the error bound as in (3.5). Nonetheless, in cases where the solution derivatives can be bounded by  $c^\eta$  for some constant  $c$ , it is possible to deduce an exponential decay rate for the error as observed in [12, Sect. 2].

**8. Conclusions.** A Radau collocation scheme is analyzed for an unconstrained control problem. For a problem with a smooth solution and a Hamiltonian which satisfies a strong convexity assumption at a local minimizer of the continuous problem, we show that the discrete approximation has a local minimizer in a neighborhood of the continuous solution, and as the number of collocation points increases, the distance in the sup-norm between the discrete solution and the continuous solution is  $O(N^{2-\eta})$  when the continuous solution has  $\eta + 1$  continuous derivatives,  $\eta \geq 3$ , and the number of collocation points  $N$  is sufficiently large. A numerical example is given which exhibits an exponential convergence rate.



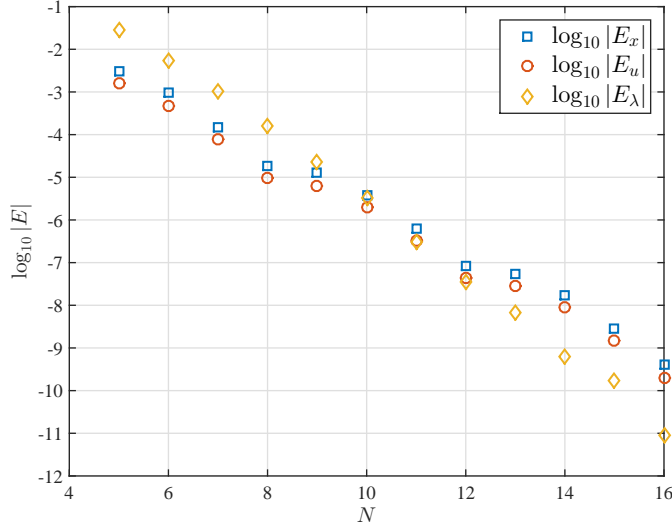


FIG. 7.1. The base 10 logarithm of the error in the sup-norm as a function of the number of collocation points.

**9. Appendix.** Before stating property (P3) in the Introduction, we showed that  $\mathbf{D}^\ddagger$  is an invertible matrix. In this section, we give an analytic formula for the inverse.

PROPOSITION 9.1. *The inverse of  $\mathbf{D}^\ddagger$  is given by*

$$\begin{aligned} D_{ij}^{\ddagger -1} &= \omega_N M_j(1) + \int_1^{\tau_i} M_j(\tau) d\tau, & 1 \leq i < N, \quad 1 \leq j < N, \\ D_{iN}^{\ddagger -1} &= -\omega_N, & 1 \leq i \leq N, \\ D_{Nj}^{\ddagger -1} &= \omega_N M_j(1), & 1 \leq j < N, \end{aligned}$$

where  $M_j$ ,  $1 \leq j < N$ , is the Lagrange interpolating basis relative to the point set  $\tau_1, \dots, \tau_{N-1}$ . That is,

$$M_j(\tau) = \prod_{\substack{i=1 \\ i \neq j}}^{N-1} \frac{\tau - \tau_i}{\tau_j - \tau_i}, \quad j = 1, \dots, N-1.$$

*Proof.* The relation (1.11) holds for any polynomial  $p$  of degree at most  $N-1$ . Let  $\mathbf{p} \in \mathbb{R}^N$  denote the vector with  $i$ -th component  $\dot{p}(\tau_i)$ . In vector form, the system of equations (1.11) can be expressed  $\mathbf{D}^\ddagger \mathbf{p} = \mathbf{p} - \mathbf{e}_N p(1)/\omega_N$ . Multiply by  $\mathbf{D}^{\ddagger -1}$  and exploit the identity  $\mathbf{D}^{\ddagger -1} \mathbf{e}_N = -\omega_N \mathbf{1}$  of (6.6) to obtain

$$\mathbf{D}^{\ddagger -1} \mathbf{p} = \mathbf{p} - \mathbf{1} p(1). \tag{9.1}$$

Since  $\dot{p}$  is a polynomial of degree at most  $N-2$ , we can only specify the derivative of  $p$  at  $N-1$  distinct points. Given any  $j$  satisfying  $1 \leq j < N$ , let us insert in (9.1) a polynomial  $p \in \mathcal{P}_{N-1}$  satisfying

$$\dot{p}(\tau_j) = 1 \quad \text{and} \quad \dot{p}(\tau_i) = 0 \quad \text{for all } i < N, \quad i \neq j.$$

A specific polynomial with this property is

$$p(\tau) = \int_1^\tau M_j(\tau) d\tau. \tag{9.2}$$

Since  $p_N = p(1) = 0$ , the last component of the right side of (9.1) vanishes to give the relation  $D_{Nj}^{\dagger -1} + D_{NN}^{\dagger -1} \dot{p}(1) = 0$ . In (6.6) we showed that all the elements in the last column of  $\mathbf{D}^{\dagger -1}$  are equal to  $-\omega_N$ , and by (9.2),  $\dot{p}(1) = M_j(1)$ . Hence, we obtain the relation

$$D_{Nj}^{\dagger -1} = -D_{NN}^{\dagger -1} \dot{p}(1) = \omega_N \dot{p}(1) = \omega_N M_j(1), \quad 1 \leq j < N. \quad (9.3)$$

Finally, let us consider  $D_{ij}^{\dagger -1}$  for  $i < N$  and  $j < N$ . We combine the  $i$ -th component of (9.1) for  $i < N$  with (9.2) to obtain

$$(\mathbf{D}^{\dagger -1} \dot{\mathbf{p}})_i = \int_1^{\tau_i} M_j(\tau) d\tau. \quad (9.4)$$

Recall that all components of  $\dot{\mathbf{p}}$  vanish except for the  $j$ -th, which is 1, and the  $N$ -th, which is  $M_j(1)$  by (9.2). Hence, (9.4) and the fact that the elements in the last column of  $\mathbf{D}^{\dagger -1}$  are all  $-\omega_N$  yield

$$D_{ij}^{\dagger -1} = \int_1^{\tau_i} M_j(\tau) d\tau - D_{iN}^{\dagger -1} M_j(1) = \omega_N M_j(1) + \int_1^{\tau_i} M_j(\tau) d\tau$$

This completes the proof.  $\square$

As noted in the Introduction, the elements in the last row of  $\mathbf{D}_{1:N}^{-1}$  are positive and sum to 2, and the last row of  $\mathbf{D}_{1:N}^{-1} \mathbf{W}^{-1/2}$  is positive and has Euclidean norm  $\sqrt{2}$ . Numerically, we observe that the absolute row sums for the first  $N - 1$  rows of  $\mathbf{D}_{1:N}^{-1}$  are always less than 2, while the Euclidean norm of the first  $N - 1$  rows of  $\mathbf{D}_{1:N}^{-1} \mathbf{W}^{-1/2}$  are always less than  $\sqrt{2}$ . Properties (P3) and (P4) are less tight in the sense that the bounds only hold as equalities asymptotically. Tables 9.1 and 9.2 show  $\|\mathbf{D}^{\dagger -1}\|_{\infty}$  and the maximum Euclidean norm of the rows of  $\mathbf{D}^{\dagger -1} \mathbf{W}^{-1/2}$  for an increasing sequence of dimensions.

$N$	25	50	75	100	125	150
norm	1.995376	1.998844	1.999486	1.999711	1.999815	1.999871
$N$	175	200	225	250	275	300
norm	1.999906	1.999928	1.999943	1.999954	1.999962	1.999968

TABLE 9.1  
 $\|\mathbf{D}^{\dagger -1}\|_{\infty}$

$N$	25	50	75	100	125	150
norm	1.412209	1.413691	1.413982	1.414083	1.414130	1.414156
$N$	175	200	225	250	275	300
norm	1.414171	1.414181	1.414188	1.414193	1.414196	1.414199

TABLE 9.2  
Maximum Euclidean norm for the rows of  $[\mathbf{W}^{1/2} \mathbf{D}^{\dagger}]^{-1}$

## REFERENCES

- [1] C. L. DARBY, W. W. HAGER, AND A. V. RAO, *Direct trajectory optimization using a variable low-order adaptive pseudospectral method*, AIAA Journal of Spacecraft and Rockets, 48 (2011), pp. 433–445.
- [2] ———, *An hp-adaptive pseudospectral method for solving optimal control problems*, Optim. Control Appl. Meth., 32 (2011), pp. 476–502.
- [3] A. L. DONTCHEV AND W. W. HAGER, *Lipschitzian stability in nonlinear control and optimization*, SIAM J. Control Optim., 31 (1993), pp. 569–603.
- [4] ———, *The Euler approximation in state constrained optimal control*, Math. Comp., 70 (2001), pp. 173–203.
- [5] A. L. DONTCHEV, W. W. HAGER, AND K. MALANOWSKI, *Error bounds for Euler approximation of a state and control constrained optimal control problem*, Numer. Funct. Anal. Optim., 21 (2000), pp. 653–682.
- [6] A. L. DONTCHEV, W. W. HAGER, AND V. M. VELIOV, *Second-order Runge-Kutta approximations in constrained optimal control*, SIAM J. Numer. Anal., 38 (2000), pp. 202–226.
- [7] G. ELNAGAR, M. KAZEMI, AND M. RAZZAGHI, *The pseudospectral Legendre method for discretizing optimal control problems*, IEEE Trans. Automat. Control, 40 (1995), pp. 1793–1796.
- [8] G. N. ELNAGAR AND M. A. KAZEMI, *Pseudospectral Chebyshev optimal control of constrained nonlinear dynamical systems*, Comput. Optim. Appl., 11 (1998), pp. 195–217.
- [9] F. FAHROO AND I. M. ROSS, *Costate estimation by a Legendre pseudospectral method*, J. Guid. Control Dyn., 24 (2001), pp. 270–277.
- [10] ———, *Direct trajectory optimization by a Chebyshev pseudospectral method*, J. Guid. Control Dyn., 25 (2002), pp. 160–166.
- [11] D. GARG, M. A. PATTERSON, C. L. DARBY, C. FRANÇOLIN, G. T. HUNTINGTON, W. W. HAGER, AND A. V. RAO, *Direct trajectory optimization and costate estimation of finite-horizon and infinite-horizon optimal control problems using a Radau pseudospectral method*, Comput. Optim. Appl., 49 (2011), pp. 335–358.
- [12] D. GARG, M. A. PATTERSON, W. W. HAGER, A. V. RAO, D. A. BENSON, AND G. T. HUNTINGTON, *A unified framework for the numerical solution of optimal control problems using pseudospectral methods*, Automatica, 46 (2010), pp. 1843–1851.
- [13] Q. GONG, I. M. ROSS, W. KANG, AND F. FAHROO, *Connections between the covector mapping theorem and convergence of pseudospectral methods for optimal control*, Comput. Optim. Appl., 41 (2008), pp. 307–335.
- [14] W. W. HAGER, *Runge-Kutta methods in optimal control and the transformed adjoint system*, Numer. Math., 87 (2000), pp. 247–282.
- [15] ———, *Numerical analysis in optimal control*, in International Series of Numerical Mathematics, K.-H. Hoffmann, I. Lasiecka, G. Leugering, J. Sprekels, and F. Tröltzsch, eds., vol. 139, Basel/Switzerland, 2001, Birkhauser Verlag, pp. 83–93.
- [16] W. W. HAGER, H. HOU, AND A. V. RAO, *Convergence rate for a Gauss collocation method applied to unconstrained optimal control*, J. Optim. Theory Appl., submitted (2015, [arxiv.org/abs/1507.08263](https://arxiv.org/abs/1507.08263)).
- [17] ———, *Lebesgue constants arising in a class of collocation methods*, IMA J. Numer. Anal., submitted (2015, [arxiv.org/abs/1507.08316](https://arxiv.org/abs/1507.08316)).
- [18] R. A. HORN AND C. R. JOHNSON, *Matrix Analysis*, Cambridge University Press, Cambridge, 2013.
- [19] S. KAMESWARAN AND L. T. BIEGLER, *Convergence rates for direct transcription of optimal control problems using collocation at radau points*, Comput. Optim. Appl., 41 (2008), pp. 81–126.
- [20] F. LIU, W. W. HAGER, AND A. V. RAO, *Adaptive mesh refinement method for optimal control using nonsmoothness detection and mesh size reduction*, J. Franklin Inst., 352 (2015), pp. 4081–4106.
- [21] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer, New York, 2nd ed., 2006.
- [22] M. A. PATTERSON, W. W. HAGER, AND A. V. RAO, *A ph mesh refinement method for optimal control*, Optim. Control Appl. Meth., 36 (2015), pp. 398–421.
- [23] G. W. REDDIEN, *Collocation at Gauss points as a discretization in optimal control*, SIAM J. Control Optim., 17 (1979), pp. 298–306.
- [24] L. N. TREFETHEN, *Approximation Theory and Approximation Practice*, SIAM Publications, Philadelphia, 2013.
- [25] P. VÉRTESI, *On Lagrange interpolation*, Period. Math. Hungar., 12 (1981), pp. 103–112.
- [26] P. WILLIAMS, *Jacobi pseudospectral method for solving optimal control problems*, J. Guid. Control Dyn., 27 (2004), pp. 293–297.