

**RECENT DEVELOPMENTS AND FUTURE PERSPECTIVES  
IN NONLINEAR SYSTEM THEORY**

John L. Casti

*International Institute for Applied Systems Analysis, Laxenburg, Austria*

RR-82-43

December 1982

Reprinted from *SIAM Review*, volume 24 number 3 (1982)

**INTERNATIONAL INSTITUTE FOR APPLIED SYSTEMS ANALYSIS**  
Laxenburg, Austria

*Research Reports*, which record research conducted at IIASA, are independently reviewed before publication. However, the views and opinions they express are not necessarily those of the Institute or the National Member Organizations that support it.

---

Reprinted with permission from *SIAM Review* 24(3):301–331, 1982.  
Copyright © 1982 by the Society for Industrial and Applied Mathematics, Philadelphia, Pennsylvania.

All rights reserved. No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopy, recording, or any information storage or retrieval system, without permission in writing from the copyright holder.

## FOREWORD

One of the opportunities for the International Institute for Applied Systems Analysis is to summarize and communicate the state of the art of theoretical tools of potential use in systems analysis, particularly when they suggest challenging new vistas of application.

In this article, the work of which was partially supported by IIASA, John L. Casti surveys recent developments in nonlinear system theory, focusing on problems of reachability, observability, and realization. He also discusses the major obstacles standing in the way of a comprehensive theory of nonlinear systems. He suggests that a more feasible and rewarding approach is to consider special classes of nonlinear problems motivated by applications and to use the structures in these classes as guides to useful and applicable results.

The article closes with a sketch of some important problem areas for future research.

HUGH J. MISER  
*Leader*  
The Craft of Systems Analysis



## RECENT DEVELOPMENTS AND FUTURE PERSPECTIVES IN NONLINEAR SYSTEM THEORY\*

JOHN L. CASTI†

**Abstract.** Results on controllability, observability and realization of input/output data for linear systems are well-known and extensively covered in a variety of books and papers. What is not so well-known is that substantial progress has been made in recent years on providing similarly detailed results for *nonlinear* processes. This paper represents a survey of the most interesting work on nonlinear systems, together with a discussion of the major obstacles standing in the way of a comprehensive theory of nonlinear systems.

**1. Basic problems and results in linear system theory.** The theory of linear dynamical processes and control has by now been developed to such an extent that it is only a slight exaggeration to term it a branch of applied mathematics, sharing equal rank with more familiar areas such as hydrodynamics, classical and quantum mechanics and electromagnetism, to name but a few. For those who doubt this assessment of linear system theory, a perusal of some of the more advanced recent literature [13], [20], [25], [44], [45] should prove to be an enlightening activity, showing how deeply imbedded system-theoretic concepts are in areas such as algebraic geometry, differential topology and Lie algebras. Conversely, the "purer" parts of mathematics have proven to be fruitful sources of inspiration for system theorists seeking more powerful tools with which to analyze and classify broad classes of problems.

Encouraged by the tremendous success in the study of linear processes, system theorists have been increasingly turning their attention and methods to the analysis of the same circle of questions for *nonlinear* systems. As one would suspect, the jungleground of nonlinearity is not easily tamed and so far no comprehensive theory has emerged capable of treating general nonlinear processes with the detail available in the linear case. Nonetheless, substantial progress has been made on several fronts and part of our story will be to survey some of the more interesting developments.

An equally important part of the picture we wish to present is to outline some of the reasons why a complete theory of nonlinear systems seems remote, at least at our current level of mathematical sophistication. All current indications point toward the conclusion that seeking a completely general theory of nonlinear systems is somewhat akin to the search for the Holy Grail: a relatively harmless activity full of many pleasant surprises and mild disappointments, but ultimately unrewarding. A far more profitable path to follow is to concentrate upon special classes of nonlinear problems, usually motivated by applications, and to use the structure inherent in these classes as a guide to useful (i.e., applicable) results. As we go along in this survey, we shall try to emphasize this approach by example, as well as by precept.

Before entering into the mainstream of nonlinear system theory and the problems inherent therein, let us briefly review some of the principal questions and results of the linear theory. We are concerned with a process described by the system of differential equations

$$\begin{aligned} (\Sigma) \quad \frac{dx}{dt} &= Fx(t) + Gu(t), & x(0) &= x_0, \\ y(t) &= Hx(t), \end{aligned}$$

\*Received by the editors June 24, 1981. This invited paper was prepared in part and published under contract DAAG-29-80-C-0091 with the U.S. Army Research Office.

†International Institute for Applied Systems Analysis, A-2361 Laxenburg, Austria and Department of Systems Engineering, University of Arizona, Tucson Arizona 85721.

where  $x$ ,  $u$  and  $y$  are  $n$ -,  $m$ - and  $p$ -dimensional vector functions, taking values in  $R^n$ ,  $R^m$  and  $R^p$ , respectively. For ease of exposition, we assume that the matrices  $F$ ,  $G$  and  $H$  are constant, although the theory extends easily to the time-varying case at the expense of more delicate notation and definitions.

The principal questions of mathematical system theory may be conveniently separated into three categories:

A. *Reachability/controllability.* Given an admissible set of input functions  $\Omega$ , determine the region  $\mathcal{R}$  of the system state space  $X$  which can be reached from the initial state  $x_0$  in some prescribed finite time  $T$  by application of inputs  $u \in \Omega$ . If  $x_0 \neq 0$  and  $\mathcal{R} = 0$ , then we have a problem of (null-) *controllability*; otherwise it is a question of *reachability*. In the case of constant  $F$  and  $G$  (the output matrix  $H$  plays no role in category A problems), with  $\Omega =$  piecewise continuous functions on  $[0, T]$ , the two notions coincide and the basic result is

THEOREM 1 [6], [14], [35]. *A state  $x$  is reachable (and controllable) if and only if  $x$  is contained in the subspace of  $X$  generated by the vectors*

$$\{G, FG, F^2G, \dots, F^{n-1}G\}.$$

The system  $\Sigma$  is said to be *completely reachable* if and only if  $\mathcal{R} = X$ , i.e.,  $x$  is reachable for every  $x \in X$ . An immediate consequence of Theorem 1 is

COROLLARY 1.  *$\Sigma$  is completely reachable if and only if the  $n \times nm$  matrix*

$$\mathcal{C} = [G|FG|F^2G|\dots|F^{n-1}G]$$

*has rank  $n$ .*

Many variations on the above theme are possible by changing  $\Omega$ ,  $\mathcal{R}$ ,  $T$  and/or admitting time-varying  $F$  and  $G$  (see [14] for details). However, the *algebraic* result given by Theorem 1 and its corollary forms the cornerstone for the study of almost all questions relating to reachability and controllability of linear systems. As we shall see below, analogous algebraic results can be obtained for large classes of *nonlinear* systems at the expense of a more elaborate mathematical machinery, further emphasizing the underlying algebraic nature of dynamical systems.

B. *Observability/constructibility.* Switching attention from inputs to outputs, we consider the class of questions centering upon what information can be deduced about the system state from the measured output. As in category A, the basic question comes in two forms, depending upon whether we wish to determine the *initial* state  $x_0$  from knowledge of *future* inputs and outputs (observability) or if we wish to determine the *current* state  $x(T)$  from knowledge of *past* inputs and outputs (constructibility). The linearity of the situation enables us to consider the case of no input ( $u = 0$ ) and, as in the controllability/reachability situation, the two basic concepts of observability and constructibility coincide if  $F$  and  $H$  are constant matrices. The main result for category B questions is

THEOREM 2 [6], [14], [35]. *A state  $x \in X$  is unobservable (unconstructible) if and only if  $x$  is of the form  $x = x_1 + \ker \theta$ , for some  $x_1 \in X$  with*

$$\theta = \begin{bmatrix} H \\ \hline HF \\ \hline HF^2 \\ \hline \cdot \\ \hline \cdot \\ \hline HF^{n-1} \end{bmatrix}.$$

Note that the basic test implicit in Theorem 2 is given in terms of *unobservable* states. Thus, any initial state  $x_0 \neq 0$  may be uniquely determined from the measured output  $y(t)$ ,  $0 \leq t \leq T$ ,  $T > 0$ , if and only if  $x_0 \neq x_1 + \ker \theta$ , for some  $x_1 \in X$ . An important corollary to Theorem 2, characterizing *complete observability/constructibility*, is

**COROLLARY 2.** *The system  $\Sigma$  is completely observable (constructible) if and only if the matrix  $\theta$  has rank  $n$ .*

The striking similarity in form between Theorems 1 and 2 suggests a duality between the concepts of reachability and observability. This idea can be made mathematically precise through the identifications

$$F \rightarrow F', \quad G \rightarrow H',$$

showing that any result concerning reachability may be transcribed into a dual result about observability, and conversely.

**C. Realizations/identification.** The basic questions subsumed under categories A and B assume for their statement that the system is given in the so-called *state-variable* form  $\Sigma$ . This leads to the basic system-theoretic problem of determining "good" state-variable models given only input/output (experimental) data.

Let  $\mathcal{U}(s)$  and  $\mathcal{Y}(s)$  denote the Laplace transforms of the input and output functions, respectively. It is then easy to see that  $\mathcal{U}$  and  $\mathcal{Y}$  are linearly related as

$$\mathcal{Y}(s) = W(s)\mathcal{U}(s),$$

where

$$W(s) = H(sI - F)^{-1}G,$$

is called the system *transfer matrix*. If  $\Sigma$  is reachable and observable,  $W(s)$  is a *strictly proper rational matrix* (i.e., the elements of  $W$  are ratios of relatively prime polynomials with the degree of the numerator less than that of the denominator), so we may expand  $W(\cdot)$  in a Laurent series about  $\infty$  obtaining

$$W(s) = \sum_{i=1}^{\infty} A_i s^{-i}.$$

The matrix  $W(s)$  or, equivalently, the infinite sequence  $\{A_1, A_2, A_3, \dots\}$  will be called the *input/output data* (or external description) of the system  $\Sigma$ . We can now state one of the central problems of mathematical system theory:

*The realization problem.* Given the input/output data of a linear system  $\Sigma$ , determine a state-variable model  $\Sigma$  such that

- (i) the input/output behavior of the model agrees *exactly* with the given data and
- (ii) the model is completely reachable and completely observable, i.e., the model is *canonical*.

*Remark.* Condition (ii), that the model be canonical, is mathematically equivalent to requiring that the dimension of the state space  $X$  of the model be minimal. However, for purposes of extension to the nonlinear case, where  $X$  may not even be a vector space, it is preferable to state the requirement as given above. Reachability and observability are natural requirements to impose on a model since unreachable and/or unobservable components of  $\Sigma$  are not implied by the data; they are pieces of the system which have been *arbitrarily* imposed by the modeler. Consequently, they have no claim to be part of a canonical, i.e., minimal model.

Perhaps surprisingly, the realization problem for linear systems has the following definitive solution.

**THEOREM 3 [35].** *For each input/output description of a system having a finite-dimensional realization there exists a canonical model  $\Sigma$ , which is unique up to a choice of coordinate system in the state space  $X$ .*

A weak form of the realization problem occurs when the dimension of  $\Sigma$  is fixed in advance, perhaps by a priori engineering or physical considerations, and only some of the components of  $F$ ,  $G$  and  $H$  need to be determined from the input/output data. This is the so-called *parameter identification* (or *structural realization*) problem and is tantamount to not only forcing the system upon the data (by fixing the dimension of  $X$ ), but also partially fixing the coordinate system in  $X$  (by demanding that certain elements of  $F$ ,  $G$  and  $H$  remain fixed). Nevertheless, much work has been done on parameter estimation, especially in the case where there are uncertainties in the data. See, for example [2], [46].

It will be noted that the realization problem demands *all* of the system input/output data before the internal model  $\Sigma$  can be chosen. In principle, this involves an *infinite* data string. Of somewhat more practical concern is the case in which only a *finite* behavior sequence

$$B_N = \{A_1, A_2, \dots, A_N\}$$

is available. The construction of a canonical model  $\Sigma_N$  from the sequence  $B_N$  constitutes the *partial realization* problem, which has only recently been definitively resolved. While a precise statement of the main result would take us too far afield, the basic conclusion is that each behavior sequence  $B_N$  has a canonical realization  $\Sigma_N$ , which may be unique (modulo a coordinate change in  $X$ ), or which may contain a certain number of undetermined parameters. Furthermore, it can be shown that as  $N$  increases (more data becomes available), the sequence of canonical realizations  $\{\Sigma_N\}$  is nested, i.e., the matrices  $F_N, G_N, H_N$  of the realization  $\Sigma_N$ , can be made to appear as submatrices in the realization  $\Sigma_{N+k}$ ,  $k \geq 1$ , if a suitable basis in  $X$  is chosen. A complete discussion of these matters is given in [32], [34].

In addition to the problems of categories A, B and C, two other broad areas are also usually considered to form part of the general field of mathematical system theory: stability theory and optimization. Generations of work on optimal control theory and stability is by now so well covered in the literature that we shall refrain from a discussion of these areas here. For the interested reader, the sources [1], [12], [49] can be recommended.

## 2. Linearization. Given a nonlinear internal model

$$(N) \quad \begin{aligned} \dot{x} &= f(x, u), & x(0) &= x_0, \\ y(t) &= h(x), \end{aligned}$$

the first temptation in analyzing questions of Type A or B is to linearize the process (N) by choosing some nominal input  $\bar{u}(t)$  and generating the corresponding reference trajectory  $\bar{x}(t)$ . Such a procedure yields the linearized dynamics

$$(\Sigma_L) \quad \begin{aligned} \dot{z} &= F(t)z + G(t)v, & z(0) &= x_0, \\ w(t) &= H(t)z, \end{aligned}$$

where

$$z(t) = x(t) - \bar{x}(t), \quad v(t) = u(t) - \bar{u}(t), \quad w(t) = y(t) - \bar{y}(t),$$



with

$$F(t) = \left( \frac{\partial f}{\partial x} \right), \quad G(t) = \left( \frac{\partial f}{\partial u} \right), \quad H(t) = \left( \frac{\partial h}{\partial x} \right),$$

with  $F(\cdot)$ ,  $G(\cdot)$  and  $H(\cdot)$  being evaluated at the pair  $(\bar{x}(t), \bar{u}(t))$ . The approach to studying reachability/observability issues is to now employ the time-varying analogues [14] of Theorems 1 and 2 for the analysis of the system  $\Sigma_L$ . We would clearly like to be able to conclude something about the controllability properties of  $(N)$  in a neighborhood of  $(\bar{x}, \bar{u})$  by studying the corresponding properties of  $\Sigma_L$ . A typical result in this direction is

**THEOREM 4** [38]. *Let the dynamics  $f(x, u)$  be  $C^1$  in a neighborhood  $U$  of  $(\bar{x}, \bar{u})$ . Then the system  $(N)$  is locally controllable if the pair  $(F(t), G(t))$  is controllable in  $U$ .*

Here "local controllability" means that for each  $x^*$  in some neighborhood of  $\bar{x}$ , there exists a piecewise continuous control  $u^*(t)$ , in some neighborhood of  $\bar{u}(t)$ ,  $0 \leq t \leq T$ , such that  $x(T) = 0$ .

The problem with the above type of linearized results is that they usually provide only sufficient conditions and are inherently local in character. As illustration of this point, consider the example  $\dot{x} = xu$  or the second order nonlinear problem

$$\dot{x}_1 = -x_1 + u, \quad \dot{x}_2 = -x_1^3 - x_2,$$

with  $|u(t)| \leq 1$ . Let  $\bar{x}(t) = 0$ ,  $\bar{u}(t) = 0$ , so that the linearized system is

$$\dot{x}_1 = Fx + Gu,$$

with

$$F = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}, \quad G = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

The pair  $(F, G)$  is not controllable since

$$\text{rank} \begin{bmatrix} 1 & \vdots & -1 \\ 0 & \vdots & 1 \end{bmatrix} = 1 < 2 = n.$$

Nevertheless, it can be shown [38] that each initial state  $(x_1^0, x_2^0)$  near  $(0, 0)$  can be transferred to the origin in finite time by a control of the above type. Thus, the system is locally controllable although the linearized approximation is not controllable.

Another obvious defect of linearization is the smoothness requirement on the dynamics  $f(x, u)$  and/or the output function  $h(x)$ . In order for the linearization to make sense, these functions must be at least continuously differentiable in each argument. While many practical processes obey this restriction, systems with switching points in the dynamics or other types of discontinuities frequently occur and would be outside the realm of straightforward linearization techniques.

**3. Nonlinear processes.** The inadequacies of linearization as outlined in the preceding section are far from the only reasons why we would like to develop a system theory for truly nonlinear processes. Some of the reasons are associated with intrinsic features of nonlinear dynamical processes, while others are more closely connected with the methods employed in the study of such processes. Let us consider the first of these aspects as it is somewhat more relevant to the issues raised in this survey.

Among the inherent difficulties associated with nonlinear processes which are not present in linear phenomena, we may cite nonuniqueness, singularities and critical dependence on parameters as features worthy of special attention.

*Nonuniqueness.* The simple scalar process

$$(1) \quad \dot{x} = ax + bx^3 + u$$

illustrates the fact that a nonlinear process may have multiple equilibria, even in the presence of no control input ( $u = 0$ ). In the event a feedback law

$$u = \phi(x)$$

is employed, the closed-loop dynamics

$$\dot{x} = ax + bx^3 + \phi(x)$$

may have an infinite (or even uncountable) number of equilibria, depending upon the form of  $\phi$ . Clearly, this situation is in stark contrast to the linear case where only the equilibrium  $x = 0$  can generically occur. Furthermore, no linearized version of (1) can possibly capture the global structure of the system equilibria manifold as a function of  $a$  and  $b$ .

*Singularities.* The solutions of many nonlinear systems may develop singularities, even though the systems themselves have smooth coefficients. The simple two-point boundary value problem

$$\ddot{x} + x\dot{x} = 0, \quad x(0) = 0, \quad x(T) = 0$$

possesses no solutions without singularities for any  $T > 0$ .

In a more system-theoretic direction, it can be shown [8] that the system

$$(2) \quad \begin{aligned} \dot{x}_1 &= x_2, & x_1(0) &= 1, \\ \dot{x}_2 &= -x_1 + u(t)x_1, & x_2(0) &= 0, \end{aligned}$$

with  $|u(t)| \leq \epsilon \ll 1$ , has a reachable set from  $x_0$  which is homeomorphic to a disk for  $T$  small, but encircles the origin for  $T$  large (see Fig. 1).

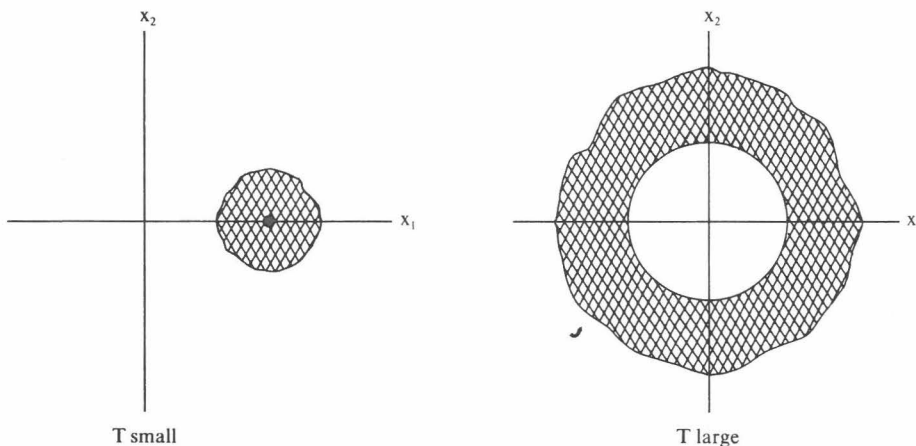


FIG. 1. The reachable set for the system (2).

The situation can be even worse than this as some nonlinear systems have a reachable set which is not even path-connected [8]. In the linear case, of course, Theorem 1 shows that the reachable set is a subspace of  $R^n$ , hence, not only simply connected but even convex. Again, no linearized version of the system (2) can hope to capture the global structure of the reachable set.

The simple bilinear system

$$\dot{x} = xu$$

also shows that a state may not be reachable from the origin with bounded control. Thus a more appropriate state space for this problem is the "punctured" region  $R^n - \{0\}$ , rather than  $R^n$  itself. In general, the "natural" state space for a nonlinear process is no longer the familiar vector space (or  $k[z]$ -module) of the linear theory, but a much more complicated mathematical object, usually some type of manifold in a Euclidean space of high dimension. Such facts account for the need to employ much more sophisticated machinery than simple linear algebra to study the structure of nonlinear processes.

*Critical dependence on parameters.* For the linear dynamical system

$$\dot{x} = Fx,$$

there are no parametric changes in the elements of  $F$  which can cause the system to have more than a single solution curve  $x(t)$ . However, this is far from the case for nonlinear processes. For example, consider the system

$$\ddot{x} + \frac{1}{2}\lambda e^x = 0, \quad x(0) = x(1) = 0.$$

For  $\lambda > \beta$  (a certain positive number), the system has no smooth solution. For  $\lambda = \beta$  there is exactly one smooth solution, while for  $0 < \lambda < \beta$  there are two solutions. Thus,  $\beta$  is a bifurcation point in the parameter space at which the character of the solution set changes radically.

To illustrate another point, consider the system

$$\dot{x}_1 = \mu x_1 + x_2 - x_1^3, \quad \dot{x}_2 = -x_1, \quad -1 \leq \mu \leq 1.$$

For each  $\mu$ ,  $-1 \leq \mu \leq 0$ , all solutions tend asymptotically to zero as  $t \rightarrow \infty$ . As  $\mu$  crosses 0, the system has a unique periodic solution  $p(\mu)$  and the origin becomes a source. For all  $\mu$ ,  $0 < \mu \leq 1$ , every nontrivial solution tends to  $p(\mu)$  as  $t \rightarrow \infty$ . Thus,  $\mu = 0$  is a bifurcation point at which the equilibrium at the origin changes suddenly from a sink to a source and a limit cycle  $p(\mu)$  is created. This so-called "Hopf bifurcation" is a consequence of the system nonlinearity and has no counterpart in linear problems.

Finally, consider the equilibria of the nonlinear system

$$\dot{x} = f(x, \alpha),$$

where  $\alpha$  is an  $m$ -dimensional vector of parameters. The equilibria  $x^*$  for which  $f(x^*, \alpha) = 0$  depend upon  $\alpha$  and we can define a multivalued map

$$\begin{aligned} \chi: A &\rightarrow X \\ \alpha &\mapsto x^*(\alpha), \end{aligned}$$

where  $A \subset R^m$ ,  $X \subset R^n$ . Under appropriate hypotheses on the function  $f$ , properties of the map  $\chi$  can be characterized using Thom's theory of catastrophes. In particular, it is of interest to categorize those submanifolds of  $A$  for which the map  $\chi$  is discontinuous, the so-called "catastrophe" manifold. Again, if  $f$  is linear the map  $\chi$  is continuous and there

is no interesting structure to analyze. Thus, no linearized version of the problem will suffice to study the geometry of the equilibrium manifold.

The above examples provide convincing evidence of the need to develop a nonlinear system theory capable of handling the same broad array of questions so successfully dealt with by the linear theory. In succeeding sections, we present some steps in this direction. As will become evident, almost everything remains to be done to complete such a program despite the impressive advances of recent years.

#### 4. Reachability and controllability.

**Smooth systems.** Certainly the area in which most progress has been made in understanding the system-theoretic behavior of nonlinear processes is in the effective characterization of reachable sets and in the determination of algebraic criteria for complete reachability. Since the mathematical apparatus involved goes somewhat beyond the elementary linear algebra which suffices for the study of linear systems, we make the following fairly standard definitions as given, for example, in [26].

Consider the nonlinear system

$$(N) \quad \begin{aligned} \dot{x} &= f(x, u), & x(0) &= x_0, \\ y(t) &= h(x), \end{aligned}$$

where  $u \in \mathcal{U} \subset \mathbb{R}^m$ ,  $x \in M$ , a  $C^\infty$ -connected manifold of dimension  $n$  and  $f$  and  $h$  are  $C^\infty$  functions of their arguments. To simplify notation, it is assumed that  $M$  admits globally defined coordinates  $x = (x_1, \dots, x_n)$ , allowing us to identify the points of  $M$  with their coordinate representations and to describe the control system (N) in the usual engineering form above. We also assume that (N) is *complete*, i.e., for every bounded measurable control  $u(t)$  and every  $x_0 \in M$ , there exists a solution of  $\dot{x} = f(x, u)$  satisfying  $x(0) = x_0$ ,  $x(t) \in M$  for all real  $t$ .

**DEFINITION 1.** Given a point  $x^* \in M$ , we say that  $x^*$  is *reachable* from  $x_0$  at  $T$  if there exists a bounded measurable control  $u(t)$ , satisfying  $u(t) \in \mathcal{U}$ , such that the system trajectory satisfies  $x(0) = x_0$ ,  $x(T) = x^*$ ,  $x(t) \in M$ ,  $0 \leq t \leq T$ . The set of states reachable from  $x_0$  is denoted as

$$\mathcal{R}(x_0) = \bigcup_{0 \leq T < \infty} \{x : x \text{ reachable from } x_0 \text{ at time } T\}.$$

We say (N) is *reachable at  $x_0$*  if  $\mathcal{R}(x_0) = M$  and *reachable* if  $\mathcal{R}(x) = M$  for all  $x \in M$ .

Since it may be necessary to travel either a long distance or a great time to reach points near  $x_0$ , the property of reachability from  $x_0$  is not always of practical use. This fact leads to a local version of reachability.

**DEFINITION 2.** (N) is *locally reachable at  $x_0$*  if for every neighborhood  $U$  of  $x_0$ ,  $\mathcal{R}(x_0) \cap U$  is also a neighborhood of  $x_0$  with the trajectory from  $x_0$  to  $\mathcal{R}(x_0) \cap U$  lying entirely within  $U$ . The system (N) is *locally reachable* if it is locally reachable for every  $x \in M$ .

The reachability concept detailed in Definition 1 is not symmetric:  $x^*$  may be reachable from  $x_0$  but not conversely (in contrast to the situation for autonomous linear systems). To remedy this situation, we need a weaker notion of reachability. This is provided by

**DEFINITION 3.** Two states  $x^*$  and  $\bar{x}$  are *weakly reachable* from each other if and only if there exist states  $x^0, x^1, \dots, x^k$  such that  $x^0 = x^*$ ,  $x^k = \bar{x}$  and either  $x^i$  is reachable from  $x^{i-1}$  or  $x^{i-1}$  is reachable from  $x^i$ ,  $i = 1, 2, \dots, k$ . The system (N) is said to be *weakly reachable* if it is weakly reachable from every  $x \in M$ . Since weak reachability is a

global concept like reachability, we can define a local version of it in correspondence to Definition 2.

Among the various reachability concepts, we have the following chain of implications

$$\begin{array}{ccc} \text{locally reachable} & \implies & \text{reachable} \\ \Downarrow & & \Downarrow \\ \text{locally weakly reachable} & \implies & \text{weakly reachable} \end{array}$$

For autonomous linear systems it can be shown that all four of the above notions coincide.

The advantage of local weak reachability over the other concepts defined above is that it lends itself to a simple algebraic test. For this, however, we need a few additional notions.

DEFINITION 4. Let  $p(x), q(x)$  be two  $C^\infty$  vector fields on  $M$ . Then the *Jacobi bracket* of  $p$  and  $q$ , denoted  $[p, q]$ , is given by

$$[p, q](x) = \left( \frac{\partial q}{\partial x} \right) p - \left( \frac{\partial p}{\partial x} \right) q \Big|_x.$$

The set of all  $C^\infty$  vector fields on  $M$  is an infinite-dimensional vector space denoted by  $X(M)$  and becomes a Lie algebra under the multiplication defined by the Jacobi bracket.

Each constant control  $u \in \Omega$  defines a vector field  $f(x, u) \in X(M)$ . We let  $\mathcal{F}_0$  denote the subset of all such vector fields, i.e.,  $\mathcal{F}_0$  is the set of all vector fields generated from  $f(x, \cdot)$  through use of constant controls.  $\mathcal{F}$  denotes the smallest subalgebra of  $X(M)$  containing  $\mathcal{F}_0$ . The elements of  $\mathcal{F}$  are linear combinations of elements of the form

$$[f^1[f^2 \cdots [f^i, f^{i+1}] \cdots]],$$

where  $f^i(x) = f(x, u^i)$  for some constant  $u^i \in \Omega$ . We let  $\mathcal{F}(x)$  be the space of tangent vectors spanned by the vector fields of  $\mathcal{F}$  at  $x$ .

DEFINITION 5. (N) is said to satisfy the *reachability rank condition* at  $x_0$  if the dimension of  $\mathcal{F}(x_0)$  is  $n$ . If this is true for every  $x \in M$ , then (N) satisfies the *reachability rank condition*.

The following theorems illustrate the importance of the reachability rank condition. The proofs may be found, for instance, in [26].

THEOREM 5. If (N) satisfies the reachability rank condition at  $x_0$ , then (N) is weakly locally reachable at  $x_0$ .

For  $C^\infty$ -systems, the converse is not quite true, but we do have

THEOREM 6. If (N) is locally weakly reachable then the reachability rank condition is satisfied on an open dense subset of  $M$  (i.e., the rank condition is satisfied generically).

In the event we strengthen the smoothness requirement on (N) from  $C^\infty$  to analytic, we can strengthen Theorems 5 and 6 to

THEOREM 7 [26]. If (N) is analytic then (N) is weakly reachable if and only if it is locally weakly reachable if and only if the reachability rank condition is satisfied.

The simplest illustration of the use of these results is to recapture the linear result of Theorem 1. In this case

$$\mathcal{F}_0 = \{Fx + Gu : u \in \Omega\}$$

so the Lie algebra is generated by the vector fields  $\{Fx, g_1, g_2, \dots, g_m\}$ , where  $g_i$  denotes the  $i$ th column of  $G$  regarded as a constant vector field. Computing brackets yields

$$\begin{aligned} [Fx, g_j] &= -Fg_j, & [g_i, g_j] &= 0, \\ [Fx, [Fx, g_j]] &= F^2g_j, & [g_i, [Fx, g_j]] &= 0, \text{ etc.} \end{aligned}$$

The Cayley–Hamilton theorem implies that  $\mathcal{F}$  is spanned by the vector fields  $Fx$  and the constant vector fields  $F^i g_j$ ,  $i = 0, 1, \dots, n-1$ ,  $j = 1, 2, \dots, m$ . Thus in this context the reachability rank condition reduces to the condition of Theorem 1, namely, (N) is locally reachable if and only if

$$\text{rank } [G | FG | F^2G | \dots | F^{n-1}G] = n.$$

However, for linear systems local reachability and reachability are equivalent, so the usual results are obtained.

The practical problem with applying the preceding results is that we have no nonlinear version of the Cayley–Hamilton theorem insuring that the test for complete reachability can be concluded in a finite number of steps. In principle, we could compute bracket after bracket in the Lie algebra generated by the  $\{f^i\}$  with no assurance that the next bracket might not yield a vector field linearly independent of those already computed.

In order to rule out the above type of behavior, we introduce the following definition.

**DEFINITION 6.** A set of vector fields  $\{f^i\}_{i=1}^r$  is called *involutive* if there exist smooth functions  $\gamma_{ijk}(x)$  such that

$$[f^i, f^j](x) = \sum_{k=1}^r \gamma_{ijk}(x) f^k(x).$$

The property of being involutive is a necessary condition in order to be able to “integrate” the vector fields  $f^1, \dots, f^r$  to obtain a solution manifold. The following theorem of Frobenius shows that this property is (with mild regularity assumptions) also sufficient to assert the existence of maximal solutions.

**THEOREM 8 [9].** Let  $\{f^i\}_{i=1}^r$  be an involutive collection of vector fields which are

a) analytic on an analytic manifold  $M$ . Then given any point  $x_0 \in M$ , there exists a maximal submanifold  $N$  containing  $x_0$  such that  $\{f^i\}$  spans the tangent space of  $N$  at each point of  $N$ .

b)  $C^\infty$  on a  $C^\infty$  manifold  $M$  with the dimension of the span of  $\{f^i\}$  constant on  $M$ . Then given any point  $x_0 \in M$ , there exists a maximal submanifold  $N$  containing  $x_0$  such that  $\{f^i\}$  spans the tangent space of  $N$  at each point of  $N$ .

As an illustration of the Frobenius theorem, consider the analytic vector fields in  $R^3$

$$f^1(x) = \begin{bmatrix} 0 \\ x_3 \\ -x_2 \end{bmatrix}, \quad f^2(x) = \begin{bmatrix} -x_3 \\ 0 \\ x_1 \end{bmatrix}, \quad f^3(x) = \begin{bmatrix} x_2 \\ -x_1 \\ 0 \end{bmatrix}.$$

It is easily verified that this collection is involutive, and if we look at any point  $x \in R^3$  then we can integrate the distribution through that point. For instance, if  $x = \frac{1}{3}(\sqrt{3}, \sqrt{3}, \sqrt{3})$ , then we obtain the set

$$N = \{x : \|x\| = 1\}$$

as the corresponding integral manifold. In fact, in this example the vectors  $f^1, f^2, f^3$  are tangent to the spherical shell  $N$  at each point. Additional details on this example are provided in [9].

In terms of the Frobenius theorem, if we allow positive *and* negative time, the problem of complete reachability for an involutive system of vector fields may be restated: does the maximal submanifold  $N = M$ ? In order to answer this question, it is necessary to have a more explicit characterization of the submanifold  $N$ . This is provided by a theorem of Chow, which also provides the underpinning for our earlier results, Theorems 5-7. But first a bit of additional notation.

Given a vector field  $f$  on  $M$ , for each  $t \exp tf$  defines a map of  $M \rightarrow M$ , which is the mapping produced by the flow on  $M$  defined by the differential equation  $\dot{x} = f(x)$ . We denote by  $\text{diff}(M)$  the group of diffeomorphisms of  $M$  and let  $\{\exp\{f^i\}\}_G$  be the smallest subgroup of  $\text{diff}(M)$  which contains  $\exp tf$  for all  $f \in \{f^i\}$ . Finally,  $\{f^i\}_{LA}$  denotes the Lie algebra of vector fields generated by  $\{f^i\}$  under the Jacobi bracket multiplication defined above. We are now in a position to state the following control-theoretic version of Chow's theorem.

**THEOREM 9** [9]. *Let  $\{f^i(x)\}_{i=1}^r$  be a collection of vector fields such that  $\{f^i(x)\}_{LA}$  is*  
 a) *analytic on an analytic manifold  $M$ . Then given any  $x_0 \in M$ , there exists a maximal submanifold  $N \subset M$  containing  $x_0$  such that*

$$\{\exp\{f^i\}\}_G x_0 = \{\exp\{f^i\}_{LA}\}_G x_0 = N.$$

b)  *$C^\infty$  on a  $C^\infty$  manifold  $M$  with  $\dim(\text{span}\{f^i(x)\}_{LA})$  constant on  $M$ . Then given any  $x_0 \in M$ , there exists a maximal submanifold  $N \subset M$  containing  $x_0$  such that*

$$\{\exp\{f^i\}\}_G x_0 = \{\exp\{f^i\}_{LA}\}_G x_0 = N.$$

**Linear-analytic systems.** The conclusions of Chow's theorem enable us to effectively resolve the reachability problem for systems of the form

$$\dot{x} = \sum_{i=1}^r u_i f_i(x), \quad x(0) = x_0.$$

However, in applications we are often confronted with systems of the form

$$(3) \quad \dot{x} = p(x) + \sum_{i=1}^r u_i(t) g^i(x), \quad x(0) = x_0.$$

In this situation, Chow's theorem has the serious drawback that it does not distinguish between positive and negative time. Thus, the submanifold  $N$  may include points which can only be reached by passing backward along the vector field  $p(x)$ . This means that the reachable set will, in general, only be a proper subset of  $N$ .

If we let  $(\exp tp)(x_0)$  denote the solution to (3) at time  $t$  corresponding to all  $u_i \equiv 0$ , while  $\mathcal{R}(t, x_0)$  denotes the reachable set at time  $t$ , then the problem of *local reachability* is to find necessary and sufficient conditions that  $(\exp tp)(x_0) \in \text{interior } \mathcal{R}(t, x_0)$  for all  $t > 0$ . Denoting  $(\text{ad } X, Y) = [X, Y]$ ,  $(\text{ad}^{k+1} X, Y) = [X, (\text{ad}^k X, Y)]$ , the basic known results on this problem are contained in

**THEOREM 10** [27], [58].

a) *A necessary and sufficient condition that for any  $T > 0$ ,  $\text{int } \bigcup_{0 \leq t \leq T} \mathcal{R}(t, x_0) \neq \emptyset$  is that  $\dim(\{p, g^i\}_{LA})(x_0) = n$ .*

b) *A necessary and sufficient condition that interior  $\mathcal{R}(t, x_0) \neq \emptyset$  for all  $t > 0$  is that*

$$\dim(\{\text{ad}^k p, g^i : k = 0, 1, \dots; i = 1, \dots, r\}_{LA})(x_0) = n.$$

c) *A sufficient condition that  $(\exp tp)(x_0) \in \text{interior } \mathcal{R}(t, x_0)$  for all  $t > 0$  is that*

$$\{(\text{ad}^j p, g^i) : j = 0, 1, 2, \dots; i = 1, 2, \dots, r\}$$

contain  $n$  linearly independent elements.

*Remark.* Condition c) of Theorem 10 is also necessary in the case  $n = 2$ . In general, though, more stringent hypotheses are required for the “rank condition” to be necessary.

To illustrate the application of the foregoing results, consider the dynamical system

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = - \begin{bmatrix} x_2 \\ \frac{\sin x_1}{x_3} \\ 0 \end{bmatrix} + u \begin{bmatrix} 0 \\ \frac{-2x_2}{x_3} \\ 1 \end{bmatrix} = p(x) + ug(x).$$

Computing the Lie brackets, we have

$$[p, g] = \begin{bmatrix} \frac{2x_2}{x_3} \\ \frac{\sin x_1}{x_3^2} - \frac{2x_2}{x_3} \\ 0 \end{bmatrix},$$

so that  $p, g$  and  $[p, g]$  span  $R^3$  unless  $x_1 = 0$  or  $\pi$  or  $x_2 = 0$ . That is, the system satisfies the reachability rank condition for all nonzero  $x_0$ .

Let us return now to the problem of local reachability. If we assume that the origin is an equilibrium point for the vector field  $p(x)$ , i.e.,  $p(0) = 0$ , and if we measure the system to be in some state  $q$  at a future time  $t_1$ , then we can consider the local controllability problem to consist in determining the existence of a stabilizing control which would drive the trajectory of the system  $x(t)$  in the “direction”  $-q$ .

To be more explicit, consider the system

$$(4) \quad \dot{x} = p(x) + u(t)g(x),$$

where  $|u(t)| \leq 1$ . Further, assume that

$$\dim \text{span} \{(\text{ad}^k p, g) : k = 0, 1, \dots\}(0) = n,$$

so that a stabilizing control law exists, at least locally (Theorem 10b)). The problem in the construction of such a law is that the directions that are “instantaneously” possible are  $p(q) + \mu g(q)$ ,  $-1 \leq \mu \leq 1$ , and  $-q$  need not be among these directions. Let us write  $q$  as

$$q = \sum_{i=1}^n \alpha_i (\text{ad}^{i-1} p, g)(0).$$

Then if we can generate the directions  $\pm(\text{ad}^j p, g)(0)$  via compositions of solutions of (4) with controls  $|u| \leq 1$ , it follows that we can generate the direction  $-q$ .

A specific illustration of how to construct the locally stabilizing law is the following taken from [27]. Let  $n = 3$  and define

$$(5) \quad q(\epsilon)x = q_3^{\epsilon_3}(|\alpha_3|\epsilon) \circ q_1^{\epsilon_1}(|\alpha_2|\epsilon) \circ q_0^{\epsilon_0}(|\alpha_1|\epsilon)x,$$

where



$$e_j = \begin{cases} + & \text{if } \alpha_j < 0, \\ - & \text{if } \alpha_j > 0, \end{cases}$$

and

$$\begin{aligned} q_0^\pm(s)(x) &= (\exp s(p \pm g))(x), \\ q_1^\pm(s)(x) &= (\exp \sqrt{s}(p \mp g)) \circ (\exp \sqrt{s}(p \pm g))(x) \\ &= (\exp (2\sqrt{s}p \pm s[p, g] + O(s)))(x), \\ q_2^\pm(s)(x) &= (\exp (2s^{1/4}p \pm s(\text{ad}^2 p, g) + O(s)))(x). \end{aligned}$$

These flows are chosen so that if  $p(0) = 0$  and  $|p(x)| \leq c|x|$ , then

$$\frac{dq^\pm}{ds}(s)(x)|_{s=0} = \pm(\text{ad}^j p, g)(x).$$

Thus, if  $x = \sum_{i=1}^3 \alpha_i(\text{ad}^{i-1} p, g)(0)$ , then

$$q(s)x - x = (s + s^{1/2} + (s/2)^{1/4})p(x) - s \sum_{i=1}^3 \alpha_i(\text{ad}^{i-1} p, g)(x) + O(s).$$

Hence, if  $x$  is near 0 and  $s$  is sufficiently small,  $q(s)x - x = -sx + O(s)$  and the above formula shows how to choose a control over the time interval  $[0, \sum_{i=1}^3 |\alpha_i|s]$  so as to move the state essentially in the same direction  $-x$ , i.e., toward the origin. Summarizing, the steps in the process are:

- (i) measure the state  $x$ ;
- (ii) express  $x = \sum_{i=1}^3 \alpha_i(\text{ad}^{i-1} p, g)(x)$ ;
- (iii) use (4) to determine an "open-loop" control  $u(t, x)$  on the interval  $0 \leq t \leq \sum_{i=1}^3 |\alpha_i|s$ ;
- (iv) remeasure the state and repeat the process.

(Note that even though the measured state  $x$  is used to compute the control, the law  $u$  is still open-loop since no state over the interval  $0 \leq t \leq \sum_{i=1}^3 |\alpha_i|s$  is measured.) The formulae for the general case of the above result are given in [27] along with a report on the convergence of the algorithm sketched in steps (i)–(iii) above.

The formulae given above for generating  $\pm(\text{ad}^k p, g)(x)$  are but one of many possible schemes. The question (as yet unanswered) arises as to whether a different scheme can be derived in which the terms  $O(s)$  are actually *insignificant* when compared to  $\pm s(\text{ad}^k p, g)$  for large  $k$ . (In the formulae given above the term  $O(s)$  in  $q_k^\pm(s)(x)$  is of the form  $(s^{1+k/2})w$ , for some vector field  $w$  in  $\{(\text{ad}^i p, g) : i = 0, 1, \dots\}_{LA}$ . Numerically, this is *not* insignificant when compared to  $\pm s(\text{ad}^k p, g)$  for  $k$  large.)

Before moving on to results for important special classes of nonlinear systems, it is of value to cite the works [24], [28], [58] for additional reachability results. Of special note is [28] in which global results are obtained for systems in which the Lie algebra  $\{p, g\}_{LA}$  is not necessarily finite dimensional. See also [40] for an excellent survey of positive-time reachability and its connection with the topological structure of the state manifold  $M$ .

**Bilinear systems.** By far the most detailed and explicit results for the reachability of nonlinear systems are those developed for bilinear processes. Bilinear systems are characterized by the equations

$$(6) \quad \dot{x} = Fx + Gu + \sum_{i=1}^m N_i x u_i(t),$$

where  $F$  and  $N_i$  are  $n \times n$  real matrices and  $G$  is an  $n \times m$  real matrix.

There are a number of theoretical and practical motivations for the study of bilinear processes which are well detailed in [48]. For now we only note that the type of nonlinearity (multiplicative) makes the system structure in some sense "closest" to the linear case. This fact enables us to employ many of the techniques and procedures already set up for linear systems.

For studying the reachability properties of (6), we consider the case  $G = 0$  (homogeneous-in-the-state systems) since the inhomogeneous case ( $G \neq 0$ ) is in a somewhat less settled state. However, it should be noted that by adding extra components to the state and/or to the control, and constraining them to be equal to 1, an inhomogeneous bilinear system may be formally studied as a homogeneous-in-the-state system.

Given a homogeneous-in-the-state system

$$(7) \quad \dot{x} = \left( F + \sum_{i=1}^m N_i u_i(t) \right) x, \quad x(0) = x_0,$$

we may write the solution as  $x(t) = X(t)x_0$ , where  $X(t) \in GL(n)$ , the nonsingular  $n \times n$  real matrices. Thus, the reachability properties of (7) are directly related to those of the system

$$(8) \quad \dot{X} = FX + \sum_{i=1}^m N_i X u_i(t), \quad X(0) = I.$$

Here the system state space is taken to be  $M = GL(n)$ . To study reachability properties of (8), we need the notion of a matrix Lie algebra.

DEFINITION 7. Given two  $n \times n$  matrices  $A$  and  $B$ , their *Lie product* is defined as

$$[A, B] \doteq AB - BA.$$

A *Lie algebra* of  $n \times n$  matrices is a subspace of  $n \times n$  matrices closed under the Lie product operation.

Let  $\mathcal{L}$  denote the Lie algebra generated by the matrices  $\{F, N_1, N_2, \dots, N_m\}$  and let  $\mathcal{R}(t, I)$  denote the reachable set for (8) at time  $t$ . Then the main reachability result for homogeneous-in-the-state bilinear systems is

THEOREM 11 [57]. For the system (8), if

$$GL(n)(\mathcal{L}) \doteq \{ \Gamma \in GL(n) : \Gamma = e^{A_1} e^{A_2} \dots e^{A_m}, A_i \in \mathcal{L}, \\ i = 1, \dots, m, \quad m = 1, 2, \dots \}$$

is compact, then

- a)  $\bigcup_{t \geq 0} \mathcal{R}(t, I) = GL(n)(\mathcal{L})$ ;
- b) there exists a  $0 < T < \infty$  such that

$$\mathcal{R}(T, I) = \bigcup_{t \geq 0} \mathcal{R}(t, I) = GL(n)(\mathcal{L}).$$

In short, Theorem 11 says that the reachable set for (8) from the identity is  $GL(n)(\mathcal{L})$  and that all points that can be reached will be attained after some finite time  $T$ .

*Remarks.* (1) In the strictly bilinear case ( $F = 0$ ), the compactness can be dropped.

(2) If  $F = 0$  the system (8) is completely reachable on  $R^n - \{0\}$ , if and only if  $\mathcal{L}$  has rank  $n$  [54].

For the inhomogeneous system (6), a convenient sufficient condition for controllability is given by the following result.

THEOREM 12 [29]. The inhomogeneous system (6) is controllable from the state  $x_0$  if the sequence of vectors  $\{S_0^1, \dots, S_0^m, S_1^1, \dots, S_{n-1}^1, \dots, S_{n-1}^m\}$  contains  $n$  linearly

independent elements, where

$$S_k^i = F^k g_i + (\text{ad}_F^k N_i) x_0, \quad k = 0, 1, \dots, n-1, \quad i = 1, 2, \dots, m,$$

$$\text{ad}_F^k N_i = [F, \text{ad}_F^{k-1} N_i], \quad \text{ad}_F^0 N_i = N_i,$$

$g_i = i$ th column of  $G$ .

An alternate approach to the study of controllability of bilinear processes is to study the equilibrium points of (6). Let  $\bar{u}$  be a constant control in the unit hypercube  $H$ . Then the equilibrium point  $x^*(\bar{u})$  is the solution of the equation

$$Fx + Nx\bar{u} + G\bar{u} = 0.$$

(Note that here we adopt the more compact notation  $\sum_{i=1}^m N_i x u_i \equiv Nx u$ .) Let us assume that whenever  $F + N'\bar{u}$  is singular,  $G\bar{u}$  is not in its range. Then the expression

$$(9) \quad x^*(\bar{u}) = -(F + N'\bar{u})^{-1} G\bar{u}$$

is the form of all possible equilibrium points, and as  $\bar{u}$  ranges over  $H$ , (9) describes the equilibrium set.

A sufficient condition for the controllability of (6) is now given by

**THEOREM 13** [38]. *The bilinear system (6) is completely controllable using piecewise-continuous inputs if*

a) *there exist constant controls  $u^+$  and  $u^-$  in  $H$  such that  $\text{Re}[\lambda_i(F + N'u^+)] > 0$  and  $\text{Re}[\lambda_i(F + N'u^-)] < 0$ , with  $x^*(u^+)$  and  $x^*(u^-)$  contained in a connected subset of the equilibrium set and*

b) *for each  $x^*(\bar{u})$ , there exists a  $v \in R^m$  such that the pair  $\{F + N'\bar{u}, [Nx^*(\bar{u}) + G]v\}$  is controllable.*

A more thorough investigation of the above criterion, together with many auxiliary results and examples, is given in the book [48].

Important properties of the reachable set for a compact control set are that it be convex and closed, regardless of the initial state. These properties are important for understanding the time-optimal control problem and for generating computational algorithms for determining optimal controls. For bilinear systems the reachable set is usually not convex (or even closed unless the control set is both compact and convex).

Since the general case is not yet settled, we consider the special case of (7) when the matrices  $N_i$  have rank 1, i.e., we can write  $N_i = b_i c_i'$ , where  $b_i$  and  $c_i$  are  $n$ -dimensional vectors. The first convexity result involves the case of small  $t$ .

**THEOREM 14** [8]. *Let  $x_0$  be given and assume that  $c_i' x_0 \neq 0, i = 1, 2, \dots, m$ . Then there exists a  $T > 0$  such that for each  $t, 0 \leq t \leq T$ , the reachable set for (7) is convex for bounded controls  $u_i(t)$ .*

In order to "globalize" this result to the case  $T = \infty$ , additional conditions on  $F, b_i$  and  $c_i$  are needed.

**THEOREM 15** [8]. *Suppose each component of  $c_i$  is nonnegative and that for all  $t > 0$  the matrix  $F + \sum_{i=1}^m u_i(t) b_i c_i'$  has nonnegative off-diagonal entries. Then the reachable set at time  $t$  is convex for  $t > 0$  for bounded controls  $u_i(t)$ .*

Another very important class of nonlinear systems of which fairly explicit reachability results have been obtained is systems governed by the polynomial dynamics

$$\dot{x}(t) = f(x) + u(t)g(x),$$

where  $f$  and  $g$  are vector fields having components which are polynomials in the entries of  $x$ . It will be useful for us to assume that  $x(t) \in k^n$ , where  $k = R$  or  $\mathbb{C}$ , with  $u(\cdot)$  being a

$k$ -valued piecewise smooth function. The extension to the case of vector controls is straightforward, at the expense of a more elaborate notation.

Since  $f$  and  $g$  are polynomial maps, it should come as no surprise that concepts from elementary algebraic geometry play a fundamental role in studying reachability. Let us recall a few basic definitions. We let  $k[s_1, \dots, s_n]$  be the ring of polynomials in the indeterminates  $s_1, \dots, s_n$  with coefficients in  $k$ , abbreviated  $k[s]$ . An *algebraic set* in  $k^n$  is the zero set for some collection of polynomials in  $k[s]$ . Thus, if  $Q \subseteq k(s)$ , then we have the natural algebraic set

$$V(Q) = \{x \in k^n : f(x) = 0 \text{ for all } f \in Q\}.$$

We let  $\mathcal{V}_Q =$  smallest ideal in  $k[s]$  containing  $Q$ . Dually, if  $S \subseteq k^n$ , then we define the ideal

$$\mathcal{V}(S) = \{f \in k[s] : f(x) = 0 \text{ for all } x \in S\}.$$

Obviously,  $S \subseteq V(\mathcal{V}(S))$ . Also if  $\mathcal{V}$  if any ideal in  $k[s]$ , we have  $\mathcal{V} \subseteq \mathcal{V}(V(\mathcal{V}))$ . The ideal  $\mathcal{V}(V(\mathcal{V}))$  is called the *radical* of  $\mathcal{V}$  and is the largest ideal defining  $V(\mathcal{V})$ .

If  $f \in k[s]$ ,  $x \in k^n$ , the *differential of  $f$  at  $x$*  is the linear function  $d_x f : k^n \rightarrow k$  given by

$$d_x f(v) = \sum_{i=1}^n \frac{\partial f}{\partial s_i}(x) v_i.$$

If  $F(s)$  is a column vector with entries in  $k[s]$ , the *Lie derivative of  $f$  with respect to  $F$* ,  $L_F(f(s))$  is given by

$$L_F(f(s)) = d_s f(F(s)) = \sum_{i=1}^m \frac{\partial f}{\partial s_i}(s) F_i(s).$$

Finally, given a set  $Q \subseteq k[s]$  and a set  $P$  whose elements are column vector of polynomials, we define

$$I(Q; P) = \text{smallest polynomial ideal in } k[s] \text{ containing } Q \text{ and closed under Lie differentiation with respect to elements of } P.$$

The ideal  $I(Q; P)$  provides the key ingredient for the following important result.

**THEOREM 16** [4], [5]. *Let  $V$  be an algebraic set in  $k^n$ . If  $\mathcal{R}(x_0) \subseteq V$  for each  $x_0 \in V$ , then  $I(\mathcal{V}(V); \{f, g\}) = \mathcal{V}(V)$ . If for any ideal  $\mathcal{V}$  defining  $V$  we have  $I(\mathcal{V}; \{f, g\}) = \mathcal{V}$ , then  $\mathcal{R}(x_0) \subseteq V$  for each  $x_0 \in V$ .*

The above theorem gives a basis for testing whether or not a given algebraic set  $V$  contains points reachable from  $x_0$ . More importantly, it also provides a procedure for constructing reachable points from  $x_0$ , namely find any ideal  $\mathcal{V}$  such that  $I(\mathcal{V}; \{f, g\}) = \mathcal{V}$ . Then the associated algebraic set  $V = \{x \in k^n : \phi(x) = 0 \text{ for all } \phi \in \mathcal{V}\}$  is certainly contained in  $\mathcal{R}(x_0)$ . Note, however, that the statement " $\mathcal{R}(x_0) \subseteq V$  for each  $x_0 \in V$  implies  $I(\mathcal{V}; \{f, g\}) = \mathcal{V}$ " is *not* true for arbitrary  $\mathcal{V}$  defining  $V$ , e.g., let  $\mathcal{V} =$  ideal in  $k[s_1, s_2]$  generated by  $\phi_1(s_1, s_2) = s_1^2, \phi_2(s_1, s_2) = -s_2$  with  $f(s_1, s_2) = \begin{pmatrix} 0 \\ s_1 \end{pmatrix}, g(s_1, s_2) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ . Then  $V(\mathcal{V}) = \{\begin{pmatrix} 0 \\ 0 \end{pmatrix}\}$  and  $\mathcal{R}(x_0) \subseteq V$  for each  $x_0 \in V$ . However,  $I(\mathcal{V}; \{f, g\}) = \mathcal{V}(V) \supset \mathcal{V}$ .

As a very useful consequence of the foregoing theorem, we can provide a computable algebraic criterion for the interior of  $\mathcal{R}(x_0)$  to be nonempty for the special polynomial system

$$\dot{x} = Ax^{[p]} + bu, \quad x(0) = x_0,$$

where  $x^{[p]}$  denotes the  $(n + p + 1)$ -tuple of weighted  $p$ -forms in the components of  $x$ , i.e.,

$$x^{[p]} = (x_1^p, \alpha_1 x_1^{p-1} x_2, \alpha_2 x_2^{p-1} x_3, \dots, x_n^p),$$

with the entries ordered lexicographically and the weights chosen so that  $\|x^{[p]}\| = \|x\|^p$ ,  $\|\cdot\|$  = Euclidean norm. If we write  $f(x) = Ax^{[p]}$ , then the  $p$ th differential  $d^p f$  defines a symmetric  $p$ -linear mapping  $k^n \times k^n \times \dots \times k^n \rightarrow k$ . Consider now a set of vectors  $\mathcal{B}$  generated as follows:

- (i) let  $b \in \mathcal{B}$ ;
- (ii) if  $v_1, \dots, v_p \in \mathcal{B}$ , then let the vector  $d^p f(v_1, \dots, v_p)$  be added to  $\mathcal{B}$ .

Define the *order* of an element  $v_{p+1} \in \mathcal{B}$  as

$$\text{order } v_{p+1} = 1 + \sum_{i=1}^p \text{order } v_i.$$

By definition, order  $b = 1$ . The connection between the set  $\mathcal{B}$  and the set  $\mathcal{R}(x_0)$  is the following result.

**THEOREM 17** [5]. *The system  $\dot{x} = Ax^{[p]} + bu$ ,  $x(0) = x_0$ , has  $\text{int } \mathcal{R}(x_0) \neq \emptyset$  if and only if the elements of  $\mathcal{B}$  of order  $\leq 1 + p + \dots + p^{n-1}$  generate  $k^n$ .*

Theorem 17 improves upon the result of Theorem 10a in that the number of elements needed to check the dimensionality condition is finite and computable *in advance*. Thus, Theorem 17 is a generalization of the standard result for constant-coefficient linear systems (Theorem 1). Proofs of Theorems 16 and 17, together with many additional results for observability and optimal control may be found in the papers [3], [5].

Other reachability/controllability results for nonlinear systems have been reported, but space precludes their inclusion. Specifically, we refer to [41] for global controllability results for perturbed linear systems and in a highly algebraic treatment, the case of systems governed by *discrete-time* polynomial dynamics is covered in detail in [54].

**5. Observability and constructability.** The general notion of observability can be stated in the following terms: given the model (N) of an input/output map  $f$ , and an input function  $u \in \Omega$  applied after  $t = t_0$ , determine the state  $x_0$  of (N) at  $t = t_0$  from knowledge of the output function  $y(t)$ ,  $t_0 \leq t \leq T$ . Another way of looking at the question is to ask if every possible pair of initial states  $x_0, x'_0$  can be distinguished by every admissible input  $u \in \Omega$ .

There are several delicate issues which arise in the theory of nonlinear observability which are masked in the linear case discussed earlier. Let us consider two of the technical considerations.

i) *Choice of inputs.* In the linear case, it is easy to show that if *any* input distinguishes points then *every* input does. So, it suffices to consider the case  $u \equiv 0$ . However, for nonlinear systems this is not the case. There may be certain inputs which do not separate points. Thus, we must be critically aware of the observability definition employed. A thorough treatment of these issues is given in [53] and [56].

ii) *Length of observation.* For continuous-time linear systems, observing the output  $y(t)$  over *any* interval  $t_0 \leq t \leq t_0 + \epsilon$ ,  $\epsilon$  arbitrary, suffices to separate points for a completely observable system. However, it may be necessary to observe  $y(t)$  over a long, even infinite, interval in order to determine  $x_0$  for a nonlinear process. Thus, it is desirable to modify the *global* concept of observability by introducing a *local* version involving only the separation of points "near"  $x_0$  in either a spatial or temporal sense.

In what follows, we shall adopt definitions to deal with the foregoing difficulties, motivated by a desire to obtain a simple algebraic test for observability analogous to that given earlier for controllability.

We consider the system

$$(N) \quad \begin{aligned} \dot{x} &= f(x, u), & x(t_0) &= x_0, \\ y(t) &= h(x), \end{aligned}$$

as given in §4.

DEFINITION 8. Two initial states  $x^0, x^1 \in M$  are termed *indistinguishable* if the systems  $(N, x^0)$  and  $(N, x^1)$  realize the same input/output map, i.e., under the same input  $u \in \Omega$ , the system (N) produces the same output  $y(t)$  for the initial states  $x^0$  and  $x^1$ . The system (N) is termed *observable* if for all  $x \in M$ , the only state indistinguishable from  $x$  is  $x$  itself.

*Remark.* Observability of (N) does not imply that every input in  $\Omega$  distinguishes all points of  $M$ . This is true, however, if the output  $y$  is a sum of a function of the initial state and a function of the input, as in the linear case.

Since observability is a global concept, we localize the concept with the following definitions.

DEFINITION 9. (N) is *locally observable* at  $x^0 \in M$  if for every open neighborhood  $U$  of  $x^0$ , the set of points indistinguishable from  $x^0$  by trajectories in  $U$  consists of  $x^0$  itself. (N) is *locally observable* if it is locally observable for every  $x \in M$ .

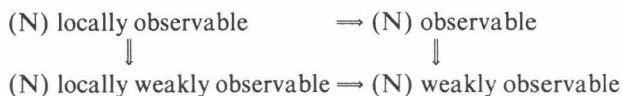
Practical considerations suggest that it may be sufficient only to distinguish points which are near to  $x^0$ , leaving open the possibility of  $x^0$  being equivalent to states  $x^1$  which are far removed. This heuristic idea motivates

DEFINITION 10. (N) is *weakly observable* at  $x^0$  if there exists an open neighborhood  $U$  of  $x^0$  such that the only point in  $U$  which is indistinguishable from  $x^0$  is  $x^0$  itself. The system (N) is *weakly observable* if it is weakly observable at every  $x \in M$ .

Again, weak observability may require that we travel far from  $U$  in order to distinguish the points of  $U$ . The following definition deals with this problem.

DEFINITION 11. (N) is *locally weakly observable* at  $x^0$  if there exists an open neighborhood  $U$  of  $x^0$  such that for every open neighborhood  $V$  of  $x^0$  contained in  $U$ , we have that the set of points indistinguishable from  $x^0$  in  $V$  is  $x^0$  itself. The system (N) is *locally weakly observable* if it is locally weakly observable for all  $x \in M$ .

As for controllability, the following diagram of implications exists:



For linear systems, all four concepts coincide.

As noted in §1, reachability and observability are dual concepts for linear systems in the precise meaning of vector space duality. In order to partially generalize this result to the manifold setting, additional machinery is required. In essence, we shall employ the duality between the space  $X(M)$  of vector fields on a manifold  $M$  and the space  $X^*(M)$  of the one-forms on  $M$ . This duality, coupled with the role  $X(M)$  played in the controllability situation, strongly suggests that the space of one-forms  $X^*(M)$  will be the appropriate vehicle for the study of nonlinear observability.

DEFINITION 12. Let  $\phi(x)$  be a  $C^\infty$  function on  $M$  with  $q$  an element of  $X(M)$ . Then the *Lie derivative* of  $\phi$  (in the direction  $q$ ),  $L_q(\phi)$ , is defined as

$$L_q(\phi)(x) = \frac{\partial \phi}{\partial x}(x) q(x).$$

(Note that the gradient  $d\phi = \partial\phi/\partial x$  is an  $n$ -dimensional row vector.)

Now let  $\mathcal{G}_0$  denote the subset of  $C^\infty(M)$  consisting of the functions  $h_1(x), h_2(x), \dots, h_p(x)$ , i.e., the components of the observation vector function  $h(x)$ . Further, we let  $\mathcal{G}$  denote the smallest vector space generated by  $\mathcal{G}_0$  and elements obtained from  $\mathcal{G}_0$  by Lie differentiation in the direction of elements of  $\mathcal{F}_0$  (recall that  $\mathcal{F}_0$  is the set of all vector fields generated from  $f(x, \cdot)$  using constant controls). A typical element of  $\mathcal{G}$  is a finite linear combination of elements of the form

$$L_{f^i}(\dots(L_{f^i}(h_i))\dots),$$

where  $f^i(x) = f(x, u^i)$  for some constant  $u^i \in \Omega$ . It is easily verified that  $\mathcal{G}$  is closed under Lie differentiation by elements of  $\mathcal{F}$  also.

Define  $X^*(M)$  as the real vector space of one-forms on  $M$ , i.e., all finite  $C^\infty(M)$  linear combinations of gradients of elements of  $C^\infty(M)$ . Further, let  $d\mathcal{G}_0 = \{d\phi : \phi \in \mathcal{G}_0\}$ ,  $d\mathcal{G} = \{d\phi : \phi \in \mathcal{G}\}$ . From the well-known identity

$$L_q(d\phi) = d(L_q(\phi)),$$

it follows that  $d\mathcal{G}$  is the smallest linear space of one-forms containing  $d\mathcal{G}_0$  and which is closed with respect to Lie differentiation by elements of  $\mathcal{F}$ . The elements of  $d\mathcal{F}$  are finite linear combinations of elements of the form

$$d(L_{f^i}(\dots(L_{f^i}(h_i))\dots)) = L_{f^i}(\dots(L_{f^i}(dh_i))\dots),$$

where  $f^i(x) = f(x, u^i)$  for some constant  $u^i \in \Omega$ . Let  $d\mathcal{G}(x)$  denote the space of vectors obtained by evaluating the elements of  $d\mathcal{G}$  at  $x$ .

DEFINITION 13. (N) is said to satisfy the *observability rank condition at  $x^0$*  if the dimension of  $d\mathcal{G}(x^0)$  equals  $n$ . If  $\dim d\mathcal{G}(x) = n$  for all  $x \in M$ , then (N) is said to satisfy the *observability rank condition*.

The observability rank condition provides an *algebraic* test for local weak observability as the next result demonstrates.

THEOREM 18 [26]. *If (N) satisfies the observability rank condition at  $x^0$  then (N) is locally weakly observable at  $x^0$ .*

The observability rank condition is "almost" a necessary condition for local weak controllability, as well, as is seen from

THEOREM 19 [26]. *If (N) is locally weakly observable then the observability rank condition is satisfied generically.*

We refer to [26] for the precise meaning of "generic" in Theorem 19. Intuitively, the set of states for which the observability rank condition fails is a "thin" set in the state space  $M$ .

For analytic systems (N), we have the stronger result

THEOREM 20 [26]. *If (N) is an analytic system then the following conditions are equivalent:*

- i) (N) satisfies the observability rank condition;
- ii) (N) is weakly observable;
- iii) (N) is locally weakly observable.

*Example.* To show that the observability rank condition generalizes Theorem 2, consider the linear system

$$\dot{x} = Fx + Gu, \quad y = Hx.$$

In this case, the space of vector fields  $\mathcal{F}$  is generated by the elements

$$\{Fx, F^i g_j, i = 1, 2, \dots, n - 1; j = 1, 2, \dots, m\}.$$

If we let  $h_j$  denote the  $j$ th row of  $H$ , then the relevant Lie derivatives are

$$\begin{aligned} L_{F^i}(h_j F^i)(x) &= h_j F^{i+1} x, \\ L_{F^i g_j}(h_k F^i)(x) &= h_k F^{i+1} g_j, \\ L_{F^i}(h_j F^i g_k) &= 0. \end{aligned}$$

Thus, by the Cayley–Hamilton theorem  $\mathcal{G}$  is generated by the set

$$\{h_i F^k x, h_i F^k g_j : i = 1, \dots, p, j = 1, \dots, m, k = 0, 1, \dots, n-1\}$$

and  $d\mathcal{G}(x)$  is generated by

$$\theta = \{h_i F^k : i = 1, \dots, p, k = 0, 1, \dots, n-1\}.$$

Since  $d\mathcal{G}(x)$  is independent of  $x$ , it is of constant dimension and the observability rank condition reduces to the requirement that the set  $\theta$  consists of  $n$  linearly independent elements.

Other important observability results for general systems are given in [21], [36], [37]. Now we consider some specific classes of nonlinear processes.

**Bilinear systems.** As in the case of controllability, considerably more detailed results are available on the observability question when we impose a bilinear structure upon the system dynamics  $f$ . For instance, consider the homogeneous system

$$(10) \quad \begin{aligned} \dot{x} &= \left( F + \sum_{i=1}^m u_i(t) N_i \right) x, & x(0) &= x_0, \\ y(t) &= Hx(t). \end{aligned}$$

We have the following result for testing whether or not indistinguishable initial states exist.

**THEOREM 21** [7]. *The homogeneous bilinear system (10) has indistinguishable initial states if and only if there exists a state coordinate transformation  $T$  such that*

$$TFT^{-1} = \begin{bmatrix} F_{11} & 0 \\ F_{21} & F_{22} \end{bmatrix}, \quad TN_i T^{-1} = \begin{bmatrix} N_{11}^i & 0 \\ N_{21}^i & N_{22}^i \end{bmatrix}, \quad HT^{-1} = [H_1 \quad 0].$$

An alternate characterization of the same result is given by

**THEOREM 22** [31]. *The set of all unobservable (i.e., indistinguishable) states of the system (10) is the largest subspace  $\theta$  of  $R^n$  invariant under  $F, N_1, \dots, N_m$ , which contains the kernel of  $H$ .*

Theorem 22 suggests the following computational algorithm for calculating the subspace  $\theta$ :

- i) Let  $U_1 = \text{range}(H')$ .
- ii) Calculate the subspace  $U_{i+1} = U_i + N_1' U_i + \dots + N_m' U_i$ .
- iii) There exists an integer  $k^*$  such that  $U_{k^*} = U_{k^*-1}$ . Continue step ii) until  $k^*$  is determined and set  $Z = \text{range } U_{k^*}$ .
- iv)  $\theta = Z^\perp$ , the orthogonal complement of  $Z$ .

Additional results on observability for bilinear systems may be found in the papers already cited in the previous section.

**Factorable systems.** An interesting class of nonlinear systems is that composed of linear systems connected in parallel with outputs multiplied. Such “factorable” systems are surprisingly general since a broad class of systems with separable Volterra kernels



may be expressed as finite sums of factorable systems. Thus, the factorable systems might be thought of as comprising the basic building blocks for the representation of constant parameter nonlinear systems. In fact, over a finite time interval, any continuous-time systems can be arbitrarily closely approximated by a factorable system.

The mathematical form of a factorable system is

$$(11) \quad \dot{x} = Fx + gu(t), \quad y(t) = \prod_{i=1}^K h_i x_i(t),$$

where we adopt the notation

$$x(t) = (x_1(t), \dots, x_K(t))', \quad g = (g_1, \dots, g_K)',$$

$$F = \begin{bmatrix} F_1 & & & 0 \\ & F_2 & & \\ & & \dots & \\ 0 & & & F_K \end{bmatrix},$$

with  $x_i$  being an  $n_i$ -dimensional vector, and the elements  $h_i, g_i, F_i$  being of corresponding sizes. Thus, the overall state vector  $x(t)$  is of dimension  $n = n_1 + \dots + n_K$ .

Since the nonlinearity occurs only in the system output, the usual reachability test from the linear theory shows that the factorable system (11) is completely reachable if and only if  $W_i(\lambda)$  and  $W_j(\lambda)$  have no poles in common for  $i \neq j$ , where  $W_k(\lambda)$  is the transfer matrix associated with the  $k$ th component subsystem. Thus, we turn attention to study of the observability properties of the system (11).

It turns out to be convenient to investigate observability for the system (11) by using the Kronecker product of the component subsystems comprising (11). Letting

$$x^\otimes(t) = x_1(t) \otimes x_2(t) \otimes \dots \otimes x_K(t),$$

where  $\otimes$  denotes the usual Kronecker product, it can be seen that  $x^\otimes(t)$  serves as a state vector for a linear system (with  $u \equiv 0$ ). We have

$$(12) \quad \frac{d}{dt} x^\otimes(t) = F^\otimes x^\otimes(t), \quad y(t) = h^\otimes x^\otimes(t),$$

with

$$F^\otimes = F_1 \otimes I_{n_2} \otimes \dots \otimes I_{n_K} + I_{n_1} \otimes F_2 \otimes I_{n_3} \otimes \dots \otimes I_{n_K} \\ + \dots + I_{n_1} \otimes I_{n_2} \otimes \dots \otimes I_{n_{K-1}} \otimes F_K,$$

$$h^\otimes = h_1 \otimes h_2 \otimes \dots \otimes h_K.$$

Knowledge of the initial state  $x^\otimes(0)$  enables us to compute (up to certain ambiguities in sign) the state  $x(0)$ . So, we say that the system (11) is completely observable if its associated linear system (12) is observable in the usual sense.

A convenient characterization of the observability of (12) is possible if we define the vector  $\Lambda_i$  of distinct characteristic roots of the matrix  $F_i$ , i.e.,

$$\Lambda_i = (\lambda_{i1}, \dots, \lambda_{i1,p_i}),$$

where  $i = 1, 2, \dots, K, p_i \leq n_i$ . The Kronecker sum of two such vectors is given by

$$\Lambda_1 \otimes \Lambda_2 = \begin{bmatrix} \lambda_{11} + \lambda_{21} \\ \lambda_{11} + \lambda_{22} \\ \vdots \\ \lambda_{11} + \lambda_{2,p_2} \\ \vdots \\ \lambda_{1,p_1} + \lambda_{21} \\ \lambda_{1,p_1} + \lambda_{22} \\ \vdots \\ \lambda_{1,p_1} + \lambda_{2,p_2} \end{bmatrix}$$

In terms of the Kronecker sum of the  $\{\Lambda_i\}$ , we characterize observability of (12) by the following result.

**THEOREM 23** [23]. *The factorable system (11) is completely observable if and only if the vector  $\Lambda_1 \otimes \Lambda_2 \otimes \dots \otimes \Lambda_K$  has distinct entries and at most one of the subsystems has multiple characteristic values.*

**Polynomial systems.** Very few results exist on the observability question for general continuous-time polynomial systems, i.e., systems of the form

$$\dot{x} = P(x, u), \quad y(t) = h(x),$$

where  $P(\cdot, \cdot)$  and  $h(\cdot)$  are polynomial functions of their arguments. However, in the *discrete time* case a considerable body of knowledge has been reported in [34]. For brevity, let us consider a representative case, the so-called (polynomial) *state-affine* system

$$(13) \quad x(t+1) = F(u(t))x(t) + G(u(t)), \quad y(t) = Hx(t),$$

where  $F(\cdot)$  and  $G(\cdot)$  are polynomial functions of  $u$  and  $H$  is a constant matrix. A particular case is that of internally-bilinear systems, when  $F$  and  $G$  are themselves linear functions of  $u$ . The observability of the state-affine system (13) is settled by the following test, which is a restatement of a result taken from [53].

**THEOREM 24** [53]. *The input sequence  $w = u_1, u_2, \dots, u_{n-1}$  distinguishes all pairs of initial states for the state-affine system (13) if and only if the matrix*

$$\theta(w) = \begin{bmatrix} H \\ HF(u_1) \\ HF(u_2)F(u_1) \\ \vdots \\ HF(u_{n-1}) \dots F(u_1) \end{bmatrix}$$

has rank  $n$ .

Thus Theorem 24 shows that any input sequence  $w$  such that the observability matrix  $\theta(w)$  is of full rank suffices to distinguish initial states for the system (13).

For a more complete discussion of various observability concepts for discrete-time polynomial systems and their interrelations, the work [53] should be consulted. Also, for continuous-time polynomial and analytic systems, the paper [56] shows that *universal* inputs  $w^*$  exist, i.e. the single universal input  $w^*$  distinguishes all initial states which are distinguishable by any input.

**6. Realization theory.** The specification of the realization problem for linear systems is simplified by the fact that it is easy to parametrize the input, output and state spaces via a globally defined coordinate system. This fact enables us to reduce the problem of construction of a canonical model from input/output data to a problem of linear algebra involving matrices. In the nonlinear case no such global coordinate system exists, in general, and it is necessary to take considerable care in defining what we mean by the problem "data." We can no longer regard the input/output data as being represented by an object as simple as an infinite sequence of matrices or, equivalently, a matrix transfer function. So, the first step in the construction of an effective nonlinear realization procedure is to develop a generalization of the transfer matrix suitable for describing the input/output behavior of a reasonably broad class of nonlinear processes.

If we consider the nonlinear system (N)

$$\begin{aligned} \dot{x} &= f(x, u), & x(0) &= x_0, \\ y &= h(x), \end{aligned}$$

then it is natural to attempt to represent the output of (N) in terms of the input as a series expansion

$$y(t) = w_0(t) + \int_0^t w_1(t, s)u(s) ds + \int_0^t \int_0^{s_1} w_2(t, s_1, s_2)u(s_2)u(s_1) ds_2 ds_1 + \dots$$

Formally, the above *Volterra series* expansion is a generalization of the linear variation of constant formula

$$y(t) = He^{Ft}x_0 + \int_0^t He^{F(t-s)}Gu(s) ds.$$

Arguing by analogy with the linear case, the realization problem for nonlinear systems may be expressed as: *given the sequence of Volterra kernels*  $\mathcal{W} = \{w_0, w_1, w_2, \dots\}$ , *find a canonical model*  $N = (f, h)$  *whose input/output behavior generates*  $\mathcal{W}$ .

Without further hypotheses on the analytic behavior of  $f, h$ , together with a suitable definition of "canonical model," the realization problem as stated is much too ambitious and, in general, unsolvable. So, let us initially consider conditions under which the Volterra series exists and is unique. Further, we restrict attention to the class of linear-analytic systems, i.e.,  $f(x, u) = f(x) + u(t)g(x)$ , where  $f(\cdot), g(\cdot)$  and  $h(\cdot)$  are analytic vector fields. The basic result for Volterra series expansions is

**THEOREM 25** [39]. *If  $f, g$  and  $h$  are analytic vector fields and if  $\dot{x} = f(x)$  has a solution on  $[0, T]$  with  $x(0) = x_0$ , then the input/output behavior of (N) has a unique Volterra series representation on  $[0, T]$ .*

In the case of a bilinear system where  $f(x) = Fx, g(x) = Gx, h(x) = x, u(\cdot) =$  scalar control, the Volterra kernels can be explicitly computed as

$$w_n(t, s_1, \dots, s_n) = e^{Ft} e^{-Fs_1} G e^{Fs_1} e^{-Fs_2} G e^{Fs_2} \dots e^{-Fs_n} G e^{Fs_n} x_0.$$

It can be shown [16] that for bilinear systems the Volterra series converges *globally* for all locally bounded  $u$ .

The global convergence of the Volterra series for bilinear processes suggests an approach to the construction of a Volterra expansion in the general case. First, expand all functions in their Taylor series, forming a sequence of bilinear approximations of increasing accuracy. We then compute the Volterra series for each bilinear approximation. However, the simple system

$$\dot{x} = u^2x, \quad x(0) = 1, \quad h(x) = x$$

shows that, in general, no Volterra expansion exists which is valid for all  $u$  such that  $\|u\|$  is sufficiently small. Further details on the above bilinear approximation technique can be found in [9].

By taking the Laplace transform of the Volterra kernels  $\{w_i\}$ , it is possible to develop a nonlinear analogue of the standard matrix transfer function of the linear theory. Such an approach as carried out, for example, in [47] provides an alternate "frequency-domain" approach to the realization problem. See also the work of Fliess [18] in this regard. We shall forego the details of such a procedure here due to space considerations, and focus our attention solely upon nonlinear systems whose input/output data is given in terms of the infinite sequence of Volterra kernels  $\{w_i\}$ .

Now let us turn to the definition of a canonical model for a nonlinear process. As noted earlier, in the linear case we say a model is canonical if it is both reachable (controllable) and observable (constructible). Such a model is also minimal in the sense that the state space has smallest possible dimension (as a vector space) over all such realizations. In order to preserve this minimality property, we make the following

DEFINITION 14. A system (N) is called *locally weakly minimal* if it is locally weakly controllable and locally weakly observable.

The relevance of Definition 14 to the realization problem is seen from the following result.

THEOREM 26 [26]. Let (N), ( $\hat{N}$ ) be two nonlinear systems with input sets  $\Omega = \hat{\Omega}$ , and state manifolds  $M$  and  $\hat{M}$  of dimensions  $m, \hat{m}$ , respectively. Suppose (N,  $x_0$ ) and ( $\hat{N}, \hat{x}_0$ ) realize the same input/output map. Then if ( $\hat{N}$ ) is locally weakly minimal,  $\hat{m} \leq m$ .

Thus, we see that two locally weakly minimal realizations of the same input/output map must be of the same state dimension which is minimal over all possible realizations.

Remark. Two locally weakly minimal realizations need not be diffeomorphic, in contrast to the linear case. This is seen from the two systems

$$(N) \quad \dot{x} = u, \quad y_1 = \cos x, \quad y_2 = \sin x,$$

$$(\hat{N}) \quad \dot{\theta} = u, \quad y_1 = \cos \theta, \quad y_2 = \sin \theta,$$

with  $\Omega = \hat{\Omega} = R, M = R, \hat{M} = S^1$ , the unit circle,  $y \in R^2, x_0 = 0, \theta_0 = 0$ . Here (N) and ( $\hat{N}$ ) realize the same input/output map. Furthermore, both systems are locally weakly controllable and observable.

The above result leaves open the question if two canonical realizations are isomorphic, i.e., given two nonlinear systems (N) and ( $\hat{N}$ ), with state manifolds  $M$  and  $\hat{M}$ ,

$$(N) \quad \dot{x} = f(x, u), \quad y = h(x),$$

$$(\hat{N}) \quad \dot{z} = \hat{f}(z, u), \quad y = \hat{h}(z),$$

when does there exist a diffeomorphism  $\phi : M \rightarrow \hat{M}$  such that  $x = \phi(z), z = \phi^{-1}(x)$  or

$$\frac{\partial \phi}{\partial z} f(\phi(z), u) = \hat{f}(\cdot, u), \quad h(\phi(z)) = \hat{h}(\cdot).$$

The answer to this question is provided by the following restatement of a result of Sussman.

**THEOREM 27** [55]. *Let there be given a mapping  $G_{x_0, u}$  which to each input  $u(t)$ ,  $0 \leq t \leq T$ , assigns a curve  $y(t)$  and assume that there exists a finite dimensional analytic complete system*

$$\begin{aligned} \dot{x} &= f(x, u), & x(0) &= x_0, \\ y &= h(x), & x &\in M, \end{aligned}$$

which realizes the map  $G_{x_0, u}$ . Then  $G_{x_0, u}$  can also be realized by a system which is weakly controllable and observable. Furthermore, any two such realizations are isomorphic.

*Remark.* In all the results above, as well as those to follow, the conditions of analyticity and completeness of the defining vector fields is crucial. The reason is clear: analyticity forces a certain type of "rigidity" upon the system; i.e., the global behavior of the system is determined by its behavior in an arbitrarily small open set. Completeness is also a natural condition since without this property the system is not totally specified, as it is then necessary to speak about the type of behavior exhibited in the neighborhood of the vector field singularity. Fortunately, analyticity and completeness are properties possessed by any class of systems defined by sets of algebraic equations, having a reasonable amount of homogeneity. For instance, linear systems and bilinear systems are included in this class, together with any other type of system which is both finite-dimensional, "algebraic," and bounded.

Now let us turn to some realization results for specific classes of nonlinear systems. For ease of notation, we consider only single-input, single-output systems citing the references for the more general case.

**Bilinear systems.** Given a sequence of Volterra kernels  $\{w_i\}_{i=0}^{\infty}$ , the first question is to determine conditions under which the sequence may be realized by a bilinear system. For this we need the concept of a factorizable sequence of kernels.

**DEFINITION 15.** A sequence of kernels  $\{w_i\}_{i=2}^{\infty}$  is said to be *factorizable* if there exist three matrix functions  $F(\cdot)$ ,  $G(\cdot)$ ,  $H(t, \cdot)$  of sizes  $n \times n$ ,  $n \times 1$ ,  $1 \times m$ , respectively such that

$$w_i(t, s_1, \dots, s_i) = H(t, s_1)F(s_2 - s_1) \cdots F(s_{i-1} - s_{i-2})G(s_i - s_{i-1}),$$

$$s_1 \leq s_2 \leq \dots \leq s_i.$$

The set  $\{F, G, H\}$  is called the *factorization* of  $\{w_i\}$  and the number  $n$  is its *dimension*. A factorization  $\{F_0, G_0, H_0\}$  of minimal dimension is called a *minimal* factorization.

We can now characterize those Volterra kernels which can be realized by a bilinear system.

**THEOREM 28** [16]. *The sequence of Volterra kernels  $\{w_i\}_{i=0}^{\infty}$  is realizable by a bilinear system if and only if  $w_1$  has a proper rational Laplace transform and  $\{w_i\}_{i=2}^{\infty}$  is factorizable by functions  $F, G, H$  with proper rational Laplace transforms.*

Let us assume that a given sequence of kernels  $\{w_i\}$  is bilinearly realizable. We then face the question of the construction of a minimal realization and its properties. The main result in this regard is

**THEOREM 29** [16]. *For a sequence of bilinearly realizable kernels  $\{w_i\}$ , the minimal realizations are such that*

- i) the state space dimension  $n_0$  is given by the dimension of the linear system whose

impulse response matrix is

$$W(s) = \begin{bmatrix} w_1(t, s) & H_0(t, s) \\ G_0(s) & F_0(s) \end{bmatrix};$$

ii) any two minimal realizations

$$\begin{aligned} \dot{x} &= Ax + Bu + Nxu, & y &= Cx, \\ \dot{z} &= \hat{A}z + \hat{B}u + \hat{N}zu, & y &= \hat{C}z, \end{aligned}$$

are related by a linear transformation of their state spaces, i.e., there exists an  $n_0 \times n_0$  matrix  $T$  such that

$$\hat{A} = TAT^{-1}, \quad \hat{B} = TB, \quad \hat{N} = TNT^{-1}, \quad \hat{C} = CT^{-1}.$$

Theorem 29 provides the basic information needed in order to actually construct the matrices  $A, B, C, N$  of a minimal realization. Since  $W(s)$  is the impulse response of a linear system of dimension  $n_0$ , there must exist three matrices  $P, Q, R$  of sizes  $n_0 \times n_0, n_0 \times (n+1), (n+1) \times n_0$  such that

$$W(s) = Re^{Ps}Q.$$

By partitioning  $Q$  and  $R$  as

$$R = \begin{pmatrix} R_1 \\ R_2 \end{pmatrix}, \quad Q = (Q_1 \quad Q_2),$$

where  $R_1$  is  $1 \times n_0$  and  $Q_1$  is  $n_0 \times 1$ , we obtain

$$\begin{aligned} w_1(t, s) &= R_1 e^{Ps} Q_1, & H_0(t, s) &= R_1 e^{Ps} Q_2, \\ G_0(s) &= R_2 e^{Ps} Q_1, & F_0(s) &= R_2 e^{Ps} Q_2. \end{aligned}$$

We now define the matrices of our minimal realization as

$$A = P, \quad B = Q_1, \quad C = R_1, \quad N = Q_2 R_2.$$

Thus the surprising conclusion is that the realization procedure for bilinear systems can be carried out using essentially the same techniques as those employed in the linear case once the minimal factorization  $\{F_0, G_0, H_0\}$  has been found.

Other approaches to the construction of bilinear realizations are discussed in [30], while results for the discrete-time case are given in [14]. The case of multilinear systems is similar to the bilinear situation and is discussed in detail in [43].

**Linear-analytic systems.** The general question of when a given Volterra series  $\{w_i\}_{i=0}^{\infty}$  admits realization by a finite-dimensional linear-analytic system  $\{f, g, h\}$  of the form

$$\dot{x} = f(x) + ug(x), \quad y = h(x),$$

has no easily computable answer, although some difficult to test conditions are given in [10]. On the other hand, if the Volterra series is finite then the results are quite easy to check and reasonably complete. For their statement, we make

DEFINITION 16. A Volterra kernel  $w(t, s_1, \dots, s_r)$  is called *separable* if it can be expressed as a finite sum

$$w(t, s_1, \dots, s_r) = \sum_{i=1}^m \gamma_i^1(t) \gamma_i^2(s_1) \dots \gamma_i^r(s_r).$$

It is called *differentiably separable* if each  $\gamma_i$  is differentiable and *stationary* if

$$w(t, s_1, \dots, s_r) = w(0, s_1 - t, s_2 - t, \dots, s_r - t).$$

The main theorem characterizing the realization of finite Volterra series by a linear-analytic system is

**THEOREM 30** [10]. *A finite Volterra series is realizable by a (stationary) linear-analytic system if and only if each term in the series is individually realizable by a (stationary) linear-analytic system. Furthermore, this will be the case if and only if the kernels are (stationary and differentiably) separable.*

The above result leaves open the question of actual computation of the vector fields  $\{f, g, h\}$  defining the linear-analytic realization of a finite Volterra series. However, this problem is formally bypassed by the following result.

**THEOREM 31** [10]. *A finite Volterra series has a (stationary) linear-analytic realization if and only if it has a (stationary) bilinear realization.*

From Theorem 31 it is tempting to conclude that there is no necessity to study linear-analytic systems when given a finite Volterra series, since we can always realize the data with a bilinear model. Unfortunately, the situation is not quite so simple since the dimension of the canonical bilinear realization will usually be somewhat greater than that of the corresponding linear-analytic model. To illustrate this point, consider the finite Volterra series

$$\begin{aligned} w_0(t) &= 0, & w_1(t, s_1) &= \exp(t - s_1), & w_2(t, s_1, s_2) &= 0, \\ w_3(t, s_1, s_2, s_3) &= \frac{1}{3} \exp[3(t - s_1)] \exp[2(s_1 - s_2)] \exp[(s_2 - s_3)], \\ w_i &= 0, & i &\geq 4. \end{aligned}$$

This series is realized by the three-dimensional bilinear model

$$\begin{aligned} \dot{x} &= Fx + Gu + Nxu, \text{ where } y(t) = x(t), \\ F &= \begin{bmatrix} 1 & 0 & \frac{1}{6} \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}, & N &= \begin{bmatrix} 0 & 0 & 0 \\ 2 & 0 & 0 \\ 0 & 3 & 0 \end{bmatrix}, & G &= \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}. \end{aligned}$$

However, the same set of kernels is also realized by the one-dimensional linear-analytic system

$$\begin{aligned} \dot{x} &= \sin x + u(t), & x(0) &= 0, \\ y(t) &= x(t). \end{aligned}$$

Another interesting example is  $\dot{x} = u, y = x^n$ , which requires an  $n$ th order bilinear realization.

**Polynomial systems.** If the system input/output map is of polynomial type, i.e., each term in the Volterra series is a polynomial function of its arguments, then an elegant realization theory for such maps has been developed by Sontag [54] in the *discrete-time* case. Since presentation of the details would entail too large an excursion into algebraic geometry, we loosely summarize the main results referring to the references for a more complete account.

For simplicity, we restrict our account to *bounded* polynomial input/output maps  $f$ , which means that there exists an integer  $\alpha$  such that the degree of each term in the

Volterra series for  $f$  is uniformly bounded by  $\alpha$ . The main realization result for bounded polynomial input/output maps is

**THEOREM 32** [54]. *If a bounded input/output map is at all realizable by a polynomial system, then it is realizable by an observable state-affine system of the form*

$$\begin{aligned}x(t+1) &= F(u(t))x(t) + G(u(t)), & x(0) &= 0, \\ y(t) &= Hx(t),\end{aligned}$$

where  $F(\cdot)$  and  $G(\cdot)$  are polynomial matrices,  $H$  is a linear map and the system state space is  $R^n$ .

An observable state-affine realization is termed *span-canonical* if the subspace of reachable states is all of  $R^n$ . Then it can be shown that a span-canonical realization of a given bounded finitely realizable  $f$  always exists and any two such realizations are related by a state coordinate change. Furthermore, a realization is span-canonical if and only if its dimension  $n$  is minimal among all state-affine realizations of the same input/output map.

Somewhat less complete results are also reported in [54] for unbounded polynomial input/output maps. The relationship between the foregoing discrete-time results and the continuous-time case is still far from clear, due mainly to the nonreversibility of difference (as opposed to differential) equations and to the different algebraic properties of difference and differential operations. To bridge this gap may turn out to be a nontrivial task, as is seen by the recent result [15] that a "finite" continuous-time map has its canonical state space unconstrained, which is far from true in the discrete-time setting.

Some additional work on polynomial systems taking a functional-analytic, rather than algebraic, approach is reported in [51].

**"Almost"-linear systems.** By imposing special types of nonlinearities upon a standard linear system, it is possible to employ techniques similar to the usual linear methods for realization of input/output maps. In this regard we note the "factorable" Volterra systems considered earlier, having the internal form

$$\dot{x} = Fx + gu(t), \quad y(t) = \prod_{i=1}^K c_i x_i(t).$$

Here the nonlinearities enter only through the system output. Utilizing tensor products, it can be shown [23] that the input/output behavior of such a process can be described by a so-called Volterra transfer function  $H(s_1, \dots, s_k)$ . Since a factorable Volterra system consists of  $K$  linear subsystems connected in parallel, with the outputs multiplied, the realization problem reduces to determining the transfer functions  $H_1(s), \dots, H_K(s)$  of each subsystem from  $H(s_1, \dots, s_k)$ . If the  $H_i(s)$  are known, then standard linear theory provides the overall system realization. Techniques for solving this problem are reported in [23].

In another direction, we could consider cascade combinations of linear subsystems and static power nonlinearities as in [52]. For inputs of the form

$$u(t) = \sum_{i=1}^P a_i e^{\lambda_i t},$$

the output of such a system is

$$y(t) = \sum_{i_1=1}^P \sum_{i_2=1}^P \cdots \sum_{i_m=1}^P \left[ \prod_{j=1}^m a_{i_j} \right] \exp \left( \sum_{k=1}^m \lambda_{i_k} t \right),$$



where  $m > 0$  is an integer defining the degree of the static nonlinearity, i.e., the block diagram of the system is

$$u(t) \rightarrow \boxed{H_0(s)} \rightarrow \boxed{(\cdot)^{P_1}} \rightarrow \boxed{H_1(s)} \rightarrow \dots \rightarrow \boxed{H_{q-1}(s)} \rightarrow \boxed{(\cdot)^{P_q}} \rightarrow \boxed{H_q(s)} \rightarrow y(t),$$

where  $\prod_{j=1}^q P_j = m$  and  $H_j(s)$  is a strictly proper rational function of degree  $\leq n$ ,  $j = 0, 1, \dots, q$ . In the work [52] an algorithm is given for solution of the minimal realization problem for such a system.

**7. Conclusions and future research.** The foregoing results leave little doubt that substantial progress has been made in nonlinear system theory over the past decade. As noted in the introduction, we have focused only upon problems of reachability, observability and realization, omitting the more well-known areas of stability and optimal control. Advances in these areas have also been impressive as can be seen from the works [11], [22]. Thus, the inescapable conclusion is that nonlinear system theory is alive and well and it is to be expected that progress on outstanding issues will be rapid in the years to come.

By way of closing remarks, let us now engage in a bit of crystal ball gazing and sketch some problem areas which seem to be most important for future research in nonlinear systems.

1) *Computational methods.* The effective employment of any of the results given here relies upon efficient computational algorithms. For those procedures which mimic the linear case (e.g., bilinear realization), good methods already exist for computing the necessary quantities. However, much remains to be done to develop comparable methods for, say, computing the reachable set for a nonlinear process or determining the Volterra series of a given input/output map from measured data.

2) *Stochastic effects.* A cornerstone of linear system theory is the Kalman filter and its associated apparatus for determining the "best" estimate of system parameters in the presence of noise. This is a special case of the more general stochastic realization problem, in which the input/output data itself is corrupted by noise and "best" estimates of the system model must be made. Again in the linear case results are available [50]. However, almost nothing has been accomplished along these lines for nonlinear processes. It seems likely, though, that with the increased understanding now available good progress can be made. We should note the works [42], [50], [59] as promising initial forays in this area.

3) *Nonanalytic systems.* Almost all interesting results for nonlinear systems are for processes whose defining vector fields are analytic. As pointed out earlier, there is good reason for this since the local behavior of analytic systems entirely determines the global behavior. However, there are interesting and important processes which do not fall into this category (e.g., systems with threshold effects, processes with phase transitions, and so on). A concerted attempt at relaxation of the analyticity assumptions can be expected to yield substantial dividends in furthering our ability to tackle a variety of problems in the social and biological sciences.

4) *Infinite dimensional processes.* In general, systems whose underlying dynamics are governed by partial differential equations or processes involving time-lag effects cannot be modeled by a finite set of ordinary differential or difference equations. Even in the linear case such processes lead to thorny analytical questions which are, as yet, far from being well under control. So, it is perhaps wildly optimistic to think that substantial advances can be made in this direction for nonlinear processes. Nonetheless, we have seen that many of the results and techniques of the linear theory can be extended to classes of nonlinear systems with modest additional effort. So, it seems reasonable to attempt an

investigation of those nonlinear problems which are the counterparts of the corresponding infinite dimensional linear processes.

**Acknowledgments.** The content and form of this paper have benefited substantially from the conversations and critiques of H. Hermes, A. Krener, R. Kalman and especially E. Sontag.

#### REFERENCES

- [1] J. AGGARWAL AND M. VIDYASAGAR, eds., *Nonlinear Systems: Stability Analysis*, Dowden, Hutchinson and Ross, Stroudsburg, PA, 1977.
- [2] K. ASTROM AND P. EYKHOFF, *System identification: A survey*, *Automatica*, 7 (1971), pp. 123–162.
- [3] J. BAILLIEUL, *Geometric methods for nonlinear optimal control problems*, *J. Optim. Theory Appl.* 25 (1978), pp. 519–548.
- [4] ———, *The geometry of homogeneous polynomial dynamical systems*, *Nonlinear Analysis-TMA*, 4 (1980), pp. 879–900.
- [5] ———, *Controllability and observability of polynomial dynamical systems*, *Nonlinear Analysis-TMA*, 5 (1981), pp. 543–552.
- [6] R. BROCKETT, *Finite-Dimensional Linear Systems*, John Wiley, New York, 1970.
- [7] ———, *On the algebraic structure of bilinear systems*, in R. Mohler and A. Ruberti, eds., *Theory and Application of Variable Structure Systems*, Academic Press, New York, 1972.
- [8] ———, *On the reachable set for bilinear systems*, in R. Mohler and A. Ruberti, eds., *Proc. 1974 Conference on Bilinear Systems*, Springer, New York, 1975.
- [9] ———, *Nonlinear systems and differential geometry*, *Proc. IEEE*, 64 (1976), pp. 61–72.
- [10] ———, *Volterra series and geometric control theory*, *Automatica*, 12 (1976), pp. 167–176.
- [11] P. BRUNOVSKY, *On the structure of optimal feedback systems*, *Proc. International Congress of Mathematicians*, Helsinki, 1978.
- [12] A. BRYSON AND Y. C. HO, *Applied Optimal Control*, Blaisdell, Waltham, MA, 1969.
- [13] A. BYRNES AND N. HURT, *On the moduli of linear dynamical systems*, *Advances in Mathematics*, 4 (1979), pp. 83–122.
- [14] J. CASTI, *Dynamical Systems and Their Applications: Linear Theory*, Academic Press, New York, 1977.
- [15] P. CROUCH, *Finite Volterra series*, doctoral dissertation, Harvard University, Cambridge, MA, 1977.
- [16] P. D'ALESSANDRO, A. ISIDORI AND A. RUBERTI, *Realization and structure theory of bilinear systems*, *SIAM J. Control*, 12 (1974), pp. 517–535.
- [17] D. ELLIOTT, *A consequence of controllability*, *J. Differential Equations*, 10 (1971), p. 364.
- [18] M. FLIESS, *Sur la réalisation des systèmes dynamiques bilinéaires*, *C. R. Acad. Sci., Paris, Ser. A*, 277 (1973), pp. 243–247.
- [19] E. FORNASINI AND G. MARCHESINI, *Algebraic realization theory of bilinear discrete-time input/output maps*, *J. Franklin Institute*, 301 (1976), pp. 143–159.
- [20] P. FUHRMANN, *Algebraic system theory: An analyst's point of view*, *J. Franklin Institute*, 301 (1976), pp. 521–540.
- [21] E. GRIFFITH AND K. S. P. KUMAR, *On the observability of nonlinear systems—I*, *J. Math. Anal. Appl.*, 35 (1971), pp. 135–147.
- [22] O. GUREL AND O. RÖSSLER, eds., *Bifurcation Theory and Applications in Scientific Disciplines*, N.Y. Academy of Sciences, Vol. 316, 1979.
- [23] T. HARPER AND W. RUGH, *Structural features of factorable Volterra systems*, *IEEE Trans. Automat. Control*, AC-21 (1976), pp. 822–832.
- [24] G. W. HAYNES AND H. HERMES, *Nonlinear controllability via Lie theory*, *SIAM J. Control*, 8 (1970), pp. 450–460.
- [25] R. HERMANN, *Linear System Theory and Introductory Algebraic Geometry*, *Interdisciplinary Mathematics*, Vol. 8, Math. Science Press, Brookline, MA, 1974.
- [26] R. HERMANN AND A. KRENER, *Nonlinear controllability and observability*, *IEEE Trans. Automat. Control*, AC-22 (1977), pp. 728–740.
- [27] H. HERMES, *On the synthesis of a stabilizing feedback control via Lie algebraic methods*, *SIAM J. Control Optim.*, 18 (1980), pp. 352–361.
- [28] R. HIRSCHORN, *Global controllability of nonlinear systems*, *SIAM J. Control Optim.*, 14 (1976), pp. 700–711.

- [29] N. IMBERT, M. CLIQUE AND A.-J. FOSSARD, *Un critère de gouvernabilité des systèmes bilinéaires*, RAIRO, J-3 (1979), pp. 55–64.
- [30] A. ISIDORI, *Direct construction of minimal bilinear realizations from nonlinear input/output maps*, IEEE Trans. Automat. Control, AC-18 (1973), pp. 626–631.
- [31] A. ISIDORI, AND A. RUBERTI, *Realization theory of bilinear systems*, in [45].
- [32] R. E. KALMAN, *On partial realizations of a linear input/output map*, in Systems and Networks, N. de Claris and R. Kalman, eds., Holt, New York, 1968.
- [33] ———, *Pattern recognition problems of multilinear machines*, in Proc. IFAC Symposium on Technical and Biological Problems of Control, Yerevan, Armenia, USSR, 1968.
- [34] ———, *On partial realizations, transfer functions and canonical forms*, Acta Polytechnica Scandinavia, 31 (1979), pp. 9–32.
- [35] R. KALMAN, P. FALB AND M. ARBIB, *Topics in Mathematical System Theory*, McGraw-Hill, New York, 1969.
- [36] Y. KOSTYNKOVSKII, *Observability of nonlinear controlled systems*, Automat. and Remote Control, 9 (1968), pp. 1384–1396.
- [37] S. R. KOU, D. ELLIOTT AND T. J. TARN, *Observability of nonlinear systems*, Inform. and Control, 22 (1973), pp. 89–99.
- [38] E. LEE AND L. MARKUS, *Foundations of Optimal Control*, John Wiley, New York, 1967.
- [39] C. LESIAK AND A. KRENER, *The existence and uniqueness of Volterra series for nonlinear systems*, IEEE Trans. Automat. Control, AC-23 (1978), pp. 1090–1095.
- [40] C. LOBRY, *Controllability of nonlinear control dynamical systems in Control Theory and Topics in Functional Analysis*, Vol. 1, International Atomic Energy Agency, Vienna, Austria, 1976.
- [41] D. L. LUKES, *On the global controllability of nonlinear systems*, in L. Weiss, ed., *Ordinary Differential Equations*, Academic Press, New York, 1971.
- [42] S. MARCUS, *Modelling and analysis of linear systems with multiplicative Poisson white noise*, in [44].
- [43] G. MARCHESINI AND G. PICCI, *Some results on the abstract realization theory of multilinear systems*, in R. Mohler and A. Ruberti, eds., *Theory and Application of Variable Structure Systems*, Academic Press, New York, 1972.
- [44] C. MARTIN AND R. HERMANN, eds., *The 1976 Ames Research Center (NASA) Conference on Geometric Control Theory*, Math. Science Press, Brookline, MA, 1977.
- [45] D. MAYNE AND R. BROCKETT, eds., *Geometric Methods in System Theory*, Reidel, Dordrecht, 1973.
- [46] R. MEHRA AND D. LAINIOTIS, eds., *System Identification Advances and Case Studies*, Academic Press, New York, 1976.
- [47] G. MITZEL, S. CLANCY AND W. RUGH, *On transfer function representation for homogeneous nonlinear systems*, IEEE Trans. Automat. Control, AC-24 (1979), pp. 252–249.
- [48] R. MOHLER, *Bilinear Control Processes*, Academic Press, New York, 1973.
- [49] M. PEIXOTO, ed., *Dynamical Systems*, Academic Press, New York, 1973.
- [50] G. PICCI, *Stochastic realization of Gaussian processes*, Proc. IEEE, 64 (1976), pp. 112–122.
- [51] W. A. PORTER, *An overview of polynomial system theory*, Proc. IEEE, 64 (1976), pp. 18–23.
- [52] W. SMITH AND W. RUGH, *On the structure of a class of nonlinear systems*, IEEE Trans. Automat. Control, AC-19 (1974), pp. 701–706.
- [53] E. SONTAG, *On the observability of polynomial systems, I: finite-time problems*, SIAM J. Control Optim., 17 (1979), pp. 139–151.
- [54] ———, *Polynomial Response Maps*, Lecture Notes in Control, 13, Springer-Verlag, Berlin, 1979.
- [55] H. SUSSMAN, *Minimal realizations of nonlinear systems*, in [45].
- [56] ———, *Single-input observability of continuous-time systems*, Math. Syst. Theory, 12 (1979), pp. 31–52.
- [57] H. SUSSMAN AND V. JURDJEVIC, *Control systems on Lie groups*, J. Differential Equations, 12 (1972), pp. 313–329.
- [58] ———, *Controllability of nonlinear systems*, J. Differential Equations, 12 (1972), pp. 95–116.
- [59] A. WILLISKY AND S. MARCUS, *Analysis of bilinear noise models in circuits and devices*, J. Franklin Institute, 301 (1976), pp. 103–122.

