

Further remarks on KKL observers

L. Brivadis^a, V. Andrieu^b, P. Bernard^c, U. Serres^b

^a*Université Paris-Saclay, CNRS, CentraleSupélec, Laboratoire des Signaux et Systèmes, 91190, Gif-sur-Yvette, France*

^b*Univ. Lyon, Université Claude Bernard Lyon 1, CNRS, LAGEPP UMR 5007, 43 bd du 11 novembre 1918, F-69100 Villeurbanne, France*

^c*Centre Automatique et Systèmes, Mines Paris, Université PSL, 60 boulevard Saint-Michel, Paris, France*

Abstract

We extend the theory of Kazantzis-Kravaris/Luenberger (KKL) observers. These observers consist in immersing the system into a linear stable filter of the output with sufficiently large dimension and appropriate structure. After discussing the uniqueness of such an immersion, we provide two main results about its existence. The first one extends a known existence result by generalizing the structure of the target linear filter and reducing its dimension. While this approach relies on a generic choice of a sufficiently large number of distinct eigenvalues in the filter, we then propose a second existence result in the novel symmetric case where instead, the target filter is a cascade of a sufficiently large number of one-dimensional filters sharing the same eigenvalue. Finally, we propose a new cascaded procedure for the design of KKL observers. This method can be used in two ways: either to pre-filter a noisy output before using it in the observer, or to simplify the construction of the observer when the system can be written as the cascade of a nonlinear system and a linear one.

Key words: observers, nonlinear systems, KKL observers, Luenberger observers

* Corresponding author : Vincent.Andrieu@gmail.com

This research was partially supported by the French Grant ANR ODISSE (ANR-19-CE48-0004-01).

Email addresses: lucas.brivadis@gmail.com (L. Brivadis),
vincent.andrieu@gmail.com (V. Andrieu),
pauline.bernard@minesparis.psl.eu (P. Bernard),
ulysse.serres@univ-lyon1.fr (U. Serres).

1 Introduction

The synthesis of observers is a standard problem in control and automation. Over the last four decades, many methods have been developed allowing the design of these estimation algorithms. The interested reader can refer to [6] which is a survey on the various methods allowing to design such algorithms for nonlinear dynamical systems. Among these listed methods, the *Kazantzis-Kravaris/Luenberger* (KKL) approach or *Nonlinear Luenberger approach* initially developed in [36,22,23,4] is one of the most powerful one from a theoretical point of view. Indeed, the so-called *backward-distinguishability* assumption guaranteeing its existence is very weak and does not require any particular normal form.

When D. Luenberger published his first results concerning the design of observers for linear systems in [25], his idea was to look for a linear change of coordinates T transforming the linear plant dynamics

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

with state x in \mathbb{R}^n , output y in \mathbb{R} , and F and H matrices in $\mathbb{R}^{n \times n}$ and $\mathbb{R}^{1 \times n}$ respectively, into a form

$$\dot{z} = Az + By \tag{1}$$

with z in \mathbb{R}^n and A Hurwitz, for which a trivial observer is simply made of a copy of its dynamics. Indeed, since Tx verifies (1), the estimation error $e = z - Tx$ for any solution z of (1) evolves along the contracting dynamics $\dot{e} = Ae$, so that any solution z converges to Tx . It follows that an estimate \hat{x} of x can be obtained from z by inverting the transformation T . Luenberger proved that when the pair (F, H) is observable, this is always possible for any Hurwitz matrix A in $\mathbb{R}^{n \times n}$ with no common eigenvalues with F , and any vector B in \mathbb{R}^n such that the pair (A, B) is controllable. This is based on the fact that the Sylvester equation

$$TF = AT + BH \tag{2}$$

ensuring that Tx follows (1) admits in this case a solution that is unique and invertible.

Some researchers have then tried to reproduce Luenberger's methodology on nonlinear systems in the form

$$\dot{x} = f(x) \quad , \quad y = h(x) \tag{3}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $h : \mathbb{R}^n \rightarrow \mathbb{R}$ are two sufficiently smooth functions. Following [36,22,23,4], a *nonlinear Luenberger* observer or *Kazantzis-*

Kravaris/Luenberger (KKL) observer is a dynamical system of the form

$$\dot{z} = Az + By \quad , \quad \hat{x} = T^{\text{inv}}(z) \quad , \quad (4)$$

with state z in \mathbb{R}^m (or \mathbb{C}^m), a Hurwitz matrix A in $\mathbb{R}^{m \times m}$ (or $\mathbb{C}^{m \times m}$), a vector B in \mathbb{R}^m such that the pair (A, B) is controllable and $T^{\text{inv}} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ a continuous map.

Given \mathcal{X} an open and bounded subset of \mathbb{R}^n containing the system trajectories of interest, following [36,22,23,4], the goal is to design the mapping T^{inv} as a uniformly continuous left inverse¹ of a C^1 mapping² $T : \text{cl}(\mathcal{X}) \rightarrow \mathbb{R}^m$ satisfying

$$\frac{\partial T}{\partial x}(x)f(x) = AT(x) + Bh(x) \quad , \quad \forall x \in \mathcal{X} \quad . \quad (5)$$

In other words, we look for a solution T to (5) for a pair (A, B) to be chosen with A Hurwitz, and, if possible, design T^{inv} uniformly continuous verifying

$$T^{\text{inv}}(T(x)) = x \quad , \quad \forall x \in \mathcal{X} \quad . \quad (6)$$

Indeed, (5) is a direct extension of the Sylvester equation (2) and says that along trajectories $t \mapsto x(t)$ of system (3) remaining in \mathcal{X} , $T(x)$ is solution to (1) with $y = h(x)$. Then, any other solution z to (1) with $y = h(x)$ converges to $T(x)$, so that $T^{\text{inv}}(z)$ asymptotically provides an estimate of x thanks to (6) by (uniform) continuity of T^{inv} . Hence the observer given by (4) converges asymptotically. This is summed up in the following theorem which is a direct rephrasing of [4, Theorem 2.2] in the case where \mathcal{X} is bounded. A proof of this theorem can be found in [11, Appendix A] .

Theorem 1.1 ([4, Theorem 2.2]) *Assume there exist m in \mathbb{N} , A in $\mathbb{R}^{m \times m}$ Hurwitz, B in \mathbb{R}^m and an injective function $T : \text{cl}(\mathcal{X}) \rightarrow \mathbb{R}^m$ satisfying (5). Then, there exists a continuous function $T^{\text{inv}} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ verifying (6), and for any such map, we have that for any solution $x : [0, +\infty) \rightarrow \mathcal{X}$ to (3) and any solution to (4) with input $y = h(x)$, the output \hat{x} is defined on $[0, +\infty)$ and $\lim_{t \rightarrow +\infty} |\hat{x}(t) - x(t)| = 0$.*

From there, different questions can be raised:

- (1) For which choice of m , A and B does an injective solution T to (5) exist?
- (2) Is this solution unique?
- (3) How to construct such a solution?

¹ In practice, $T^{\text{inv}}(z) = \text{argmin}_{x \in \text{cl}(\mathcal{X})} |z - T(x)|$ can be employed even though it is not continuous.

² As shown in [4], we do not need T to be C^1 as long as the Lie derivative of T along f exists.

- (4) How does the choice of the pair (A, B) impact the observer performance and how to optimize this choice ?

The existence question (1) was first considered in [36], [22] and [24] in the analytic context and around an equilibrium point. Then, the localness was dropped following another perspective in [23] where a global existence result was proposed based on a strong observability assumption which unfortunately did not provide an indication on the necessary dimension of the pair (A, B) . This problem was solved in [4] by proving the existence of the injective map T under a weak *backward-distinguishability* condition, for A complex diagonal of dimension $n + 1$, with a generic choice of $n + 1$ *distinct* complex eigenvalues. Those results have then been extended to non autonomous systems [5], discrete-time autonomous systems [12], and to the problem of *functional* observer design when the full state is not observable and only a function of the state needs to be estimated [38].

In terms of design, an explicit expression of the map T can sometimes be found in particular contexts such as parameter identification [2], state/parameter estimation for electrical machines [20,8]. Otherwise, when an expression for T , or its left-inverse T^{inv} is not available, approximation approaches have been proposed as in [26]. More recently, numerical methods based on neural networks are being developed to learn a model of the maps T and T^{inv} based on the generation of a data pairs (x, z) approximating $(x, T(x))$ through backward and forward integration of the dynamics [35,14,13,32]. Insight about how those methods work is given in Example 1.2. Note that while the difficulty in knowing the transformation T is a peculiarity of the KKL observer, its left-inversion on the other hand is a problem appearing in most observer designs, as soon as the observer is designed in other coordinates [6].

Once a method to compute the map T and, more importantly T^* , is available, the question of the link between the choice of (A, B) and the performance of the observer comes naturally. In particular, the impact of the choice of (A, B) on the robustness to noise has been observed to be significant in some applications (see e.g., [7] in the time-varying context). It is then tempting, and sometimes crucial, to optimize this choice as done in [7] in a particular application, and proposed in [14] in a general setting (see also [21]). But the quantification of performance and the formulation of the optimization problem is not trivial in the nonlinear context and is a novel active field of research. In order to allow a maximum of flexibility, efficiency and validity of the optimization process, it is helpful to clarify theoretically

- (i) the dimension and the class of admissible pairs (A, B) allowing the existence of an injective map T ,
- (ii) whether any map T_a obtained numerically has any chance to be this injective map T , which is related to uniqueness of solutions.

Hence, it is worth investigating theoretically questions (1) and (2) above. In particular, as noticed in [7,14], the result of [4] considering diagonal complex matrices A with $n+1$ generic complex (independent) eigenvalues is restrictive: we should be able to pick generic real pairs (A, B) with both complex and real eigenvalues, maybe of multiple multiplicity and sufficiently large dimension. It is thus the goal of this paper to provide existence results in such wider contexts.

Example 1.2 *In order to get a better grip on the theoretical contributions given all along the paper, we give an illustrative observation problem example. Consider a Wilson-Cowan model [40] of the form*

$$\dot{x} = -x + S(Wx), \quad y = Cx \quad (7)$$

where the state $x \in \mathbb{R}^n$ quantifies the neural activity of each neuron group, $W \in \mathbb{R}^{n \times n}$ describes the interconnections, $S : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a nonlinear activation function to be chosen of the form

$$S(x_1, \dots, x_n) = (s_1(x_1), \dots, s_n(x_n))$$

with $s_i : \mathbb{R} \rightarrow \mathbb{R}$, and $C \in \mathbb{R}^{p \times n}$ indicates which neural activity is measured. In this case, it is not clear how to obtain an explicit analytical expression of a map T solving (5), let alone its left-inverse T^* . Instead, we follow [35] and proceed as follows:

- Pick $m \in \mathbb{N}$, larger than n , a pair $(A, B) \in \mathbb{R}^{m \times m} \times \mathbb{R}^m$ and a grid of initial conditions $x_0 \in \mathbb{R}^n$.
- Simulate forward the interconnection of (3) and (1) from each x_0 , with z initialized arbitrarily, typically at zero, and for a time t_c ten times larger than $1/|\operatorname{Re}(\lambda_m(A))|$, where $\lambda_m(A)$ is the eigenvalue with largest real part.
- Remove the first part of the simulation where $t \in [0, t_c]$, and keep a certain amount of data pairs $(x(t), z(t))$ at times larger than t_c . Because t_c is large compared to the eigenvalues of A , we know that $z(t) \approx T(x(t))$, with T solution to (5).
- Fit numerical models T_a and T_a^* such that $z(t) \approx T_a(x(t))$ and $x(t) \approx T_a^*(z(t))$.

We are then ready to implement the KKL observer (4) with T^* replaced by the learned model T_a^* . Results of a simulation on (7) with $n = 4$, $W = \begin{pmatrix} 2 & 2 & 0 & 0 \\ -2 & 2 & 2 & 0 \\ 0 & -2 & 2 & 2 \\ 0 & 0 & -2 & 2 \end{pmatrix}$, $y = x_1$, $s_i(x) = \tanh(x)$, are shown on Figures A.1 and A.2 in appendix. In terms of observer parameters, we picked $\dim z = n + 1 = 5$, A a block diagonal matrix with eigenvalues equal to those of a Bessel filter of dimension 5 and cut-off frequency $\omega_c = 0.4$, and $B = (1, 1, 1, 1)^\top$. In terms of data, we took 100 random initial conditions in the forward-invariant compact set $\mathcal{X} = [-2, 2]^4$ and 200 random points in each generated trajectory. For the

learning of T_a and T_a^* , we used a neural network with 5 hidden layers of 50 neurons, SiLU activation function, learning rate 5×10^{-4} , weight decay 10^{-8} , and scheduler with factor 0.1, patience 3 and threshold 10^{-4} .

Note that observability of the system is here guaranteed since the knowledge of y and its 3 first successive derivatives determines the state uniquely. However, the choice of A does not follow the recommendations of [4]. Indeed, the constraint of a real matrix forces to pick conjugate eigenvalues, so that A contains only 3 independent degrees of freedom in \mathbb{C} , instead of $n+1 = 5$ recommended in [4]. In other words, the real matrices of dimension 5 could very well be in the zero-measure set for which the injectivity of T is not guaranteed. To follow [4], one should pick $n+1$ eigenvalues in \mathbb{C} , with thus a real implementation of dimension $2n+2$. In this paper, we will show that actually, there is always a generic choice of a real pairs (A, B) of dimension $2n+1$ guaranteeing the existence of an injective map T . Of course, as illustrated in this example, smaller dimensions may be used but the theory provides at least an upper-bound. We will also show that other structures with eigenvalues with multiplicity are admissible.

Finally, one may wonder whether the learned map T_a has any chance to be close to the injective map T provided by the theorem, a question linked to uniqueness of the solution to (5), also treated in this paper.

Organization and contributions of the paper. As justified above, we give further and more precise answers to the questions (1) and (2) of existence and uniqueness, while we leave aside the implementation-related questions (3) and (4), since they essentially require improvement of the numerical methods developed in [35,14,13,32], as well as a better understanding of the link between the pair (A, B) and the observer performance, which are both interesting problems in their own right and out of the scope of this paper. More precisely, we start by discussing uniqueness in Section 2 with sufficient conditions described in Theorem 2.1. Then, two novel existence results are provided in Section 3:

- one refining [4] with almost any choice of controllable pair (A, B) having A real diagonalizable of dimension $2n+1$ (see Theorem 3.4),
- the second in the analytic context for a different structure of the pair (A, B) with only one eigenvalue of sufficiently large multiplicity (see Theorem 4.2).

Then, in Section 4, we introduce a cascaded procedure which facilitates the synthesis of such an observer for certain cascaded nonlinear systems and which allows the use of a filtered version of y in the observer (see Theorem 4.2).

Notations. A map $\rho : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is of class- \mathcal{K} if it is continuous, increasing and such that $\rho(0) = 0$. For a differential equation $\dot{x} = f(x)$ with $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ locally Lipschitz, we denote by $X(x, t)$ the value at time t of the solution initialized at $x \in \mathbb{R}^n$ at time 0, and by $(\sigma_{\mathcal{O}}^-(x), \sigma_{\mathcal{O}}^+(x))$ the maximal domain

of existence of $t \mapsto X(x, t)$ in an open set \mathcal{O} . When $\mathcal{O} = \mathbb{R}^n$, we just denote $(\sigma^-(x), \sigma^+(x))$. Given a subset $\mathcal{X} \subseteq \mathbb{R}^n$ and a positive real number δ , $\mathcal{X} + \delta$ is the open set defined as

$$\mathcal{X} + \delta = \{x \in \mathbb{R}^n \mid \exists x_{\mathcal{X}} \in \mathcal{X}, |x - x_{\mathcal{X}}| < \delta\}. \quad (8)$$

The real and imaginary parts of a complex number are denoted by Re and Im respectively, and

$$\mathbb{R}_{\rho} = \{\lambda \in \mathbb{R} : \lambda < -\rho\} \quad , \quad \mathbb{C}_{\rho} = \{\lambda \in \mathbb{C} : \text{Re}(\lambda) < -\rho\}. \quad (9)$$

We say that a map $g : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ is in $C^{\infty}(\mathbb{R}; C^1(\mathbb{R}^n; \mathbb{R}))$ if $\lambda \mapsto g(\lambda, x)$ is C^{∞} for all $x \in \mathbb{R}^n$ and $x \mapsto \frac{\partial^k g}{\partial \lambda^k}(\lambda, x)$ is C^1 for all $\lambda \in \mathbb{R}$ and all $k \in \mathbb{N}$.

2 Remarks on the uniqueness of the map T

Typical KKL theorems as in [4] or in the next section, provide the existence of an injective solution T to the partial differential equation (PDE) (5). However, we might find other solutions of this PDE (via exact computations or a numerical approach [35,13]) and it is legitimate to wonder how much any such maps differ.

Theorem 2.1 *Let \mathcal{O} be a subset of \mathbb{R}^n that is backward invariant³ by f and consider A a Hurwitz matrix in $\mathbb{R}^{m \times m}$ and B a vector in \mathbb{R}^m . Let $T_a : \mathcal{O} \mapsto \mathbb{R}^m$ and $T_b : \mathcal{O} \mapsto \mathbb{R}^m$ be two C^1 solutions of*

$$\frac{\partial T}{\partial x}(x)f(x) = AT(x) + Bh(x) \quad , \quad \forall x \in \mathcal{O}.$$

If either

- (i) T_a and T_b are bounded on \mathcal{O} ,
- (ii) or there exist positive constants $\kappa_f, \rho_f, \kappa_a, q_a, \rho_a, \kappa_b, q_b$ and ρ_b such that for all $x \in \mathcal{O}$,

$$|f(x)| \leq \kappa_f |x| + \rho_f, \quad |T_a(x)| \leq \kappa_a |x|^{q_a} + \rho_a, \quad |T_b(x)| \leq \kappa_b |x|^{q_b} + \rho_b$$

with $\kappa_f q_a < |\text{Re}(\lambda_m(A))|$ and $\kappa_f q_b < |\text{Re}(\lambda_m(A))|$, where $\lambda_m(A)$ is the eigenvalue of A having the largest real part,

- (iii) or there exist positive constants $\kappa_a, \rho_a, \kappa_b$ and ρ_b such that for all $x \in \mathcal{O}$,

$$\left| \frac{\partial T_a}{\partial x}(x)f(x) \right| \leq \kappa_a |T_a(x)| + \rho_a \quad \text{and} \quad \left| \frac{\partial T_b}{\partial x}(x)f(x) \right| \leq \kappa_b |T_b(x)| + \rho_b$$

³ That is, for all $x \in \mathcal{O}$, $\sigma^-(x) = -\infty$ and $X(x, t) \in \mathcal{O}$ for all $t \leq 0$.

with $\kappa_a < |\operatorname{Re}(\lambda_m(A))|$ and $\kappa_b < |\operatorname{Re}(\lambda_m(A))|$,

then $T_a(x) = T_b(x)$ for all x in \mathcal{O} . In particular, if \mathcal{O} is compact, then (i) is satisfied and $T_a(x) = T_b(x)$ for all x in \mathcal{O} .

Proof : See Appendix C. □

If multiple solutions to the PDE (5) exist, the injectivity of one solution may not imply the injectivity of all solutions. In the following we give sufficient conditions for a particular bounded injective solution denoted T to exist. If another solution T_a is found by other means on a backward invariant set of the system and if this map T_a is bounded on that set, then it is actually unique and coincides with the theoretical injective solution. Otherwise, injectivity of T_a is not ensured a priori, and must be checked on each individual example.

Example 2.2 *For illustration purpose only, we consider the trivial one-dimensional example*

$$\dot{x} = -ax \quad , \quad y = x$$

with $a > 0$. This example falls into the original linear Luenberger context [25] where an injective solution to PDE (5) is known to exist with dimension $m = 1$. Taking $A = -\lambda$ and $B = 1$ with $\lambda > 0$ and $\lambda \neq a$, the map T_λ^0 defined by $T_\lambda^0(x) = \frac{1}{\lambda-a}x$ verifies the PDE

$$\frac{\partial T_\lambda}{\partial x}(x)f(x) = -\lambda T_\lambda(x) + h(x)$$

everywhere. Clearly, T_λ^0 is injective. However, for any real number α ,

$$T_\lambda(x) = \alpha \operatorname{sign}(x)|x|^{\frac{\lambda}{a}} + \frac{1}{\lambda-a}x$$

is also a C^1 solution to the PDE everywhere and clearly T_λ can be non injective for some values of α . Note that in this example, there is no backward-invariant set apart from $\{0\}$, where the maps T_λ indeed agree.

Actually, given x in \mathcal{X} such that $X(x, t)$ belongs to the bounded set \mathcal{X} for all $t \geq 0$, the ω -limit set of x $\omega(x) = \bigcap_{t \geq 0} \operatorname{cl}(\bigcup_{s \geq t} \{X(x, s)\})$ is a non-empty compact backward and forward invariant set. Hence, with the former proposition, T and T_a coincide on this set. In other words, T_a coincides on the set $\omega(x)$ with an injective map. Besides, because T and T_a are both solutions to the PDE, the state z of the KKL observer converges both towards $T_a(X(x, t))$ and $T(X(x, t))$, which tend to each other. However, this does not mean that $X(x, t)$ can be uniquely determined asymptotically from the knowledge of T_a , since there could be other $x' \in \mathcal{X}$ such that $T(x^*) = T_a(x^*) = T_a(x')$ for $x^* \in \omega(x)$ and $x' \notin \omega(x)$ (this is possible since T_a is only injective on $\omega(x)$).

Example 2.3 Consider the case of the Wilson-Cowan model (see Example 1.2) with bounded activation function, i.e., assume that there exists a positive constant \bar{S} such that $|s_i(x)| \leq \bar{S}$ for all $x \in \mathbb{R}$ and all $i \in \{1, \dots, n\}$. Then, for each $x \in \mathbb{R}^n$, the corresponding trajectory $t \mapsto X(x, t)$ is such that

$$X(x, t) = e^{-t}x + \int_0^t e^{-(t-s)}S(WX(x, s))ds$$

Hence $|X(x, t)| \leq |x| + \bar{S}$ for all $t \geq 0$. Thus, according to Theorem 2.1, solutions of (5) are unique on $\omega(x)$. When $n = 2$ and S and W are chosen such that 0 is an unstable equilibrium of (7), the system has a unique asymptotically stable limit cycle with basin of attraction $\mathbb{R}^n \setminus \{0\}$. Then, the compact set containing 0 and having this limit cycle as boundary is backward invariant. Then, by Theorem 2.1, solutions of (5) are unique on this set.

Note however that because the observer is supposed to estimate trajectories remaining in \mathcal{X} for all positive times, the map f may be modified outside of \mathcal{X} as long as the observability properties given below are preserved. It is thus usually possible to replace f by modified dynamics $\dot{x} = \chi(x)f(x)$, which admit a backward invariant compact set (by making f vanish outside of a larger open set containing $\text{cl}(\mathcal{X})$ and ensuring observability, for instance as in (D.12) below). Once this regularization has been done, any solution T_a to the PDE found on that set is unique and thus injective on \mathcal{X} if the required observability properties are preserved. For instance, in a numerical KKL design [35,13], where T_a is learned on a compact set, a trick to ensure injectivity is to apply the learning procedure to the modified f on the whole backward invariant set. This has the additional advantage to make the solutions well-defined and bounded in backward-time, which is crucial in the learning procedure.

3 Remarks on the existence of an injective map T

3.1 Existence result based on A diagonalizable with $2n + 1$ eigenvalues

As shown in [4], one of the main interests of the KKL observer is that its existence is guaranteed under a very weak observability assumption. Indeed, assume that for any $x \in \mathcal{X}$, the past output path $t \mapsto h(X(x, t))$ of (3) restricted to the time in which the trajectory remains in a certain set determines x uniquely. Then, from [4], it is sufficient to choose $m = 2n + 2$ and A the real representation of a diagonal Hurwitz complex matrix in \mathbb{C}^{n+1} to get the existence of an injective map T solving (5). The specific observability condition made is the following.

Definition 3.1 ($((\mathcal{O}, \delta_d)$ -backward distinguishability) For a given open set

\mathcal{O} of \mathbb{R}^n containing $\text{cl}(\mathcal{X})$ and a given positive real number δ_d , system (3) is (\mathcal{O}, δ_d) -backward distinguishable if for each pair of distinct points x_a and x_b in \mathcal{O} , there exists a negative time t in $(\max\{\sigma_{\mathcal{O}+\delta_d}^-(x_a), \sigma_{\mathcal{O}+\delta_d}^-(x_b)\}, 0]$ such that $h(X(x_a, t)) \neq h(X(x_b, t))$.

This distinguishability assumption says that the present state x can be distinguished from other states in \mathcal{O} by looking at the past output path restricted to the time in which the solution remains in $\mathcal{O} + \delta_d$.

Example 3.2 *Let us investigate under which condition the Wilson-Cowan model (7) is backward distinguishable. Assume that S is either globally Lipschitz or locally Lipschitz and bounded, so that according to the Picard–Lindelöf theorem, all maximal solutions are global (i.e. $\sigma^-(x) = -\infty$ and $\sigma^+(x) = +\infty$ for all $x \in \mathbb{R}^n$). Make the following observability hypotheses:*

- (i) The pair (C, W) is observable and its observability matrix $\begin{pmatrix} C \\ CW \\ \vdots \\ CW^{k-1} \end{pmatrix}$ is lower triangular;
- (ii) For all $k \in \{1, \dots, n-1\}$, the map s_k is injective.

All these assumptions are satisfied in the case considered in Example 1.2 where $n = 4$, $s_i = \tanh$, $W = \begin{pmatrix} 2 & 2 & 0 & 0 \\ -2 & 2 & 2 & 0 \\ 0 & -2 & 2 & 2 \\ 0 & 0 & -2 & 2 \end{pmatrix}$ and $C = \begin{pmatrix} 1 & 0 & 0 & 0 \end{pmatrix}$. Under these assumptions, let us show that the system is (\mathbb{R}^n, δ_d) -backward distinguishable for any $\delta_d > 0$. Let $x_a, x_b \in \mathbb{R}^n$. Set $y_{a,b}(t) = CX(x_{a,b}, t)$ and assume that $y_a(t) = y_b(t)$ for all $t \leq 0$. Let us show by induction that $CW^k(X(x_a, t) - X(x_b, t)) = 0$ for all $k \in \{0, n-1\}$ and all $t \leq 0$. Indeed, since the pair (C, W) is observable, applying this at $t = 0$ yields $x_a = x_b$ and thus backward-distinguishability. First, $y_a(t) = y_b(t)$ yields $C(X(x_a, t) - X(x_b, t)) = 0$. Assume that for some $k \in \{0, n-1\}$, $CW^j(X(x_a, t) - X(x_b, t)) = 0$ for all $j \in \{0, \dots, k\}$. Taking the derivative with respect to t yields $CW^j(S(WX(x_a, t)) - S(WX(x_b, t))) = 0$. Due to the invertibility and lower triangular structure of the observability matrix of (C, W) , we get that $s_j(W_j(X(x_a, t))) - s_j(W_j(X(x_b, t))) = 0$ for all $j \in \{0, \dots, k\}$, where W_j denotes the j -th line of W . Since s_j is injective, this yields $W_j(X(x_a, t) - X(x_b, t)) = 0$. Hence, due to the triangular structure of the observability matrix of (C, W) , CW^{k+1} being a linear combination of W_j 's for $j \leq k$, we have that $CW^{k+1}(X(x_a, t) - X(x_b, t)) = 0$ for all $t \leq 0$, which ends the induction. Note that the previous reasoning exploits the successive derivatives of $y_{a,b}$ at time zero, and thus it still holds if $y_a(t) = y_b(t)$ for $t \in [-\epsilon, 0]$ only, with ϵ arbitrarily small. It follows that for any open set \mathcal{O} containing $\text{cl}(\mathcal{X})$ and for any $\delta_d > 0$, the Wilson-Cowan model is (\mathcal{O}, δ_d) -backward distinguishable.

One of the results obtained in [4] can be reformulated as follows.

Theorem 3.3 ([4]) *Assume that system (3) is (\mathcal{O}, δ_d) -backward distinguishable for some open bounded set \mathcal{O} containing $\text{cl}(\mathcal{X})$ and some $\delta_d > 0$. Then there exist a positive real number ρ and a zero Lebesgue measure subset \mathcal{I} of $(\mathbb{C}_\rho)^{n+1}$ with \mathbb{C}_ρ defined in (9) such that for each $(\lambda_1, \dots, \lambda_{n+1})$ in $(\mathbb{C}_\rho)^{n+1} \setminus \mathcal{I}$, there exists an injective map $T : \mathcal{O} \rightarrow \mathbb{C}^{n+1}$ verifying (5) with*

$$A = \text{diag}(\lambda_1 \dots, \lambda_{n+1}) , B = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^\top . \quad (10)$$

In [4], this result was not stated in this way. However, it is a direct consequence of the fact that we restrict our analysis to a bounded set \mathcal{X} and that the output is one-dimensional. If the output is multi-dimensional then the same result holds but with the filter (1) applied to each output and thus T concatenating the solutions to (5) for each output.

Note that this observer can be realized in \mathbb{R}^{2n+2} by picking

$$A_{\text{real}} = \text{diag} \left(\begin{bmatrix} \text{Re}(\lambda_i) & -\text{Im}(\lambda_i) \\ \text{Im}(\lambda_i) & \text{Re}(\lambda_i) \end{bmatrix} \right) , B_{\text{real}} = \begin{bmatrix} b_{\text{real}} \\ \vdots \\ b_{\text{real}} \end{bmatrix} , b_{\text{real}} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} . \quad (11)$$

However, we see that the existence result imposes strong constraints on the matrices A and B . This is different from the result of Luenberger for linear systems for which no assumptions besides controllability and a spectrum different from F is required. The result we obtain in this paper is the following one.

Theorem 3.4 *Assume that system (3) is (\mathcal{O}, δ_d) -backward distinguishable for some open bounded set \mathcal{O} containing $\text{cl}(\mathcal{X})$ and some $\delta_d > 0$. Then, there exist a positive real number ρ and a zero Lebesgue measure subset \mathcal{J} of $\mathbb{R}^{(2n+1) \times (2n+1)} \times \mathbb{R}^{2n+1}$ such that for any pair (A, B) in $(\mathbb{R}^{(2n+1) \times (2n+1)} \times \mathbb{R}^{2n+1}) \setminus \mathcal{J}$ with $A + \rho I$ Hurwitz, there exists an injective map $T : \mathcal{O} \rightarrow \mathbb{R}^{2n+1}$ verifying (5).*

Proof : See Appendix D. □

Remark 3.5 *Theorem 3.4 generalizes Theorem 3.3, in several directions :*

- (1) *The observer matrices do not need to be with complex eigenvalues.*
- (2) *The dimension of the observer is $2n+1$, whereas the observer in Theorem 3.3 is of real dimension $2n+2$. This allows to recover some well known fact in observability theory that it is generically sufficient to extract $2n+1$ pieces of information from the output path to observe a state of dimension n (see for instance [1,39,18,16,37]).*

(3) The matrices A and B do not need to have a particular structure unlike in (11). In particular, we show that they may be almost any controllable pair (A, B) with A diagonalizable. This “almost any” pair actually comes from an “almost any” choice of distinct $p_{\mathbb{C}}$ complex conjugate eigenvalues and $p_{\mathbb{R}}$ real eigenvalues in A such that $2p_{\mathbb{C}} + p_{\mathbb{R}} \geq 2n + 1$. Indeed, we show that for any such $p_{\mathbb{C}}$ and $p_{\mathbb{R}}$, the set of eigenvalues in $\mathbb{C}_{\rho}^{p_{\mathbb{C}}} \times \mathbb{R}_{\rho}^{p_{\mathbb{R}}}$ which do not provide injectivity of T for (A, B) defined in (10) is of zero-Lebesgue measure in $\mathbb{C}_{\rho}^{p_{\mathbb{C}}} \times \mathbb{R}_{\rho}^{p_{\mathbb{R}}}$. This generalizes Theorem 3.3 where $p_{\mathbb{C}}$ is fixed to $n + 1$ and $p_{\mathbb{R}} = 0$. Then, in the case where $2p_{\mathbb{C}} + p_{\mathbb{R}} = 2n + 1$, we show that the set of matrices in $\mathbb{R}^{(2n+1) \times (2n+1)}$ having eigenvalues in the union of those zero-measure sets is of zero-measure in $\mathbb{R}^{(2n+1) \times (2n+1)}$. We refer the reader to the proof for more details on this genericity result.

Remark 3.6 Following what has been done in [5], this result can be extended to time-varying systems in the form

$$\dot{x} = f(x, t) , \quad y = h(x, t) \tag{12}$$

with $f : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}^n$ and $h : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}$. In that case, combining the arguments of [5] with the proof of Theorem 3.4, one obtains that the dimension of the observer needs to be $2n + 2$. Note that this was already the real dimension of the observer in [5], but Theorem 3.4 shows that no structural constraints on A and B need to be imposed.

Example 3.7 It has been shown in Example 3.2 that the Wilson-Cowan model (7) satisfies the backward-distinguishability condition for any open set \mathcal{O} and any $\delta_d > 0$ under observability assumptions on (C, W) and S . Therefore, Theorem 3.4 can be applied to (7) under these same hypotheses. Hence, Theorem 3.4 shows that a generic choice of real matrices (A, B) of dimension $2n + 1$ guarantees the existence of an injective map T satisfying (5). Actually, it was shown in Example 1.2 that a smaller dimension may be used in this particular case.

In the following subsection, another existence result is given for some particular structures of matrices A and B which are not covered by Theorem 3.4.

3.2 Existence result for A triangular with single eigenvalue

In this section, inspired by [9], we consider the case in which the pair (A, B) in $\mathbb{R}^{m \times m} \times \mathbb{R}^m$ is in the form

$$A_{\lambda, m} = \begin{pmatrix} \lambda & 0 & \cdots & \cdots & 0 \\ 1 & \lambda & \ddots & & \vdots \\ 0 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 & \lambda \end{pmatrix}, \quad B_m = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{pmatrix} \quad (13)$$

for some negative real number λ . This corresponds to a chain of filters, successively filtering m times the output with the same eigenvalue λ . In other words, instead of parallelizing the filters with different eigenvalues, this choice rather exploits the *depth* of the filter.

This case is not covered by Theorem 3.4 because its proof relies on diagonalizable matrices A with a generic choice of *distinct* eigenvalues. Instead, the choice of (13) is parameterized by a single real parameter λ , which is typically in the zero-measure set of Theorem 3.4. It thus requires another type of analysis which leads to the following result in the analytic context.

Theorem 3.8 *Assume that system (3) is (\mathcal{O}, δ_d) -backward distinguishable for some open backward invariant set \mathcal{O} containing $\text{cl}(\mathcal{X})$ and some $\delta_d > 0$. Let Θ be a non-empty open subset of $\mathbb{R}_{<0}$. Assume there exists a C^∞ map $T_0 : \Theta \times \mathcal{O} \mapsto \mathbb{R}$, such that for each λ in Θ , $x \mapsto T_0(\lambda, x)$ is an analytic bounded function on \mathcal{O} which satisfies*

$$\frac{\partial T_0}{\partial x}(\lambda, x)f(x) = \lambda T_0(\lambda, x) + h(x). \quad (14)$$

Assume moreover that h is bounded on \mathcal{O} . Then, for each λ in Θ , for any compact subset $\mathcal{C} \subset \mathcal{O}$, there exists $m^ \in \mathbb{N}$ such that for all $m \geq m^*$, the (unique) solution $T_{\lambda, m}$ of (5) with $(A, B) = (A_{\lambda, m}, B_m) \in \mathbb{R}^{m \times m} \times \mathbb{R}^m$ given in (13) is injective on \mathcal{C} .*

Proof : See Appendix E. □

Remark 3.9 *In Theorem 3.8, the existence of an analytic solution $x \mapsto T_0(\lambda, x)$ of (14) is assumed for $\lambda \in \Theta$. In the proof, it is shown that T_0 actually coincides with the map S , defined by (E.1). By the Lebesgue dominated convergence theorem, if $x \mapsto h \circ X(x, s) \in C^\infty(\mathcal{O}, \mathbb{R})$ for all $s \leq 0$ and if for each multi-index α there exist a continuous map $M_\alpha : \mathcal{O} \rightarrow \mathbb{R}_+$ and an*

integrable map $\varphi_\alpha : \mathbb{R}_- \rightarrow \mathbb{R}_+$ such that

$$\exp(-\lambda s) \left| \frac{\partial^\alpha (h \circ X)}{\partial x^\alpha}(x, s) \right| \leq \varphi_\alpha(s) M_\alpha(x), \quad \forall s \in \mathbb{R}_-, \forall x \in \mathcal{O}, \quad (15)$$

then $S(\lambda, \cdot) \in C^\infty(\mathcal{O}, \mathbb{R})$ and its partial derivatives are given by

$$\frac{\partial^\alpha S}{\partial x^\alpha}(\lambda, x) = \int_{-\infty}^0 \exp(-\lambda s) \frac{\partial^\alpha (h \circ X)}{\partial x^\alpha}(x, s) ds. \quad (16)$$

Moreover,

$$\left| \frac{\partial^\alpha S}{\partial x^\alpha}(\lambda, x) \right| \leq M_\alpha(x) \int_{-\infty}^0 \varphi_\alpha(s) ds. \quad (17)$$

Hence, if for all compact sets $\mathcal{C} \subset \mathcal{O}$ there exists a positive constant γ such that $M_\alpha(x) \int_{-\infty}^0 \varphi_\alpha(s) ds \leq \gamma^{|\alpha|+1} \alpha!$ for all $x \in \mathcal{C}$ and all multi-indices α , then $S(\lambda, \cdot)$ is analytic.

Remark 3.10 Note that Theorem 3.8 does not readily extend to time varying systems in the form (12). This is mainly due to the fact that the t component is not bounded and the dimension m (if it exists) may depend on time.

Example 3.11 Consider an harmonic oscillator with unknown frequency investigated in [34] and modelled as

$$\dot{x}_1 = x_2, \dot{x}_2 = -x_3 x_1, \dot{x}_3 = 0, y = x_1. \quad (18)$$

In that case, for any $\rho > 0$ and $\varpi > 0$, the bounded set

$$\mathcal{O} = \left\{ x \in \mathbb{R}^3, \frac{1}{\rho} < x_3 x_1^2 + x_2^2 < \rho, \frac{1}{\varpi} < x_3 < \varpi \right\} \quad (19)$$

is backward invariant along the dynamics. Besides, the map defined on $\mathbb{R}_{<0} \times \mathcal{O}$ by

$$T_0(\lambda, x) = \frac{-\lambda x_1 - x_2}{\lambda^2 + x_3} \quad (20)$$

solves the PDE (14). According to Theorem 3.8, we know that for any $\lambda < 0$ and for any compact subset \mathcal{C} of \mathcal{O} , there exists an integer m such that the (unique) solution to (5) with $(A, B) = (A_{\lambda, m}, B_m) \in \mathbb{R}^{m \times m} \times \mathbb{R}^m$ given in (13) is injective on \mathcal{C} .

As shown in the proof of Theorem 3.8, $T_{\lambda, m}$ is built by successively differentiating T_0 with respect to λ as defined in (E.2) until obtaining an injective map. On this example, it is shown in the long version of this paper [11, Appendix E] that we can pick $m = 4$ and that the associated map $T_{\lambda, 4}$ defined by

$$T_{\lambda, 4}(x) = \left(T_0(\lambda, x), \frac{\partial T_0}{\partial \lambda}(\lambda, x), \frac{\partial^2 T_0}{\partial \lambda^2}(\lambda, x), \frac{\partial^3 T_0}{\partial \lambda^3}(\lambda, x) \right) \quad (21)$$

is actually injective on $(\mathbb{R}^2 \setminus \{(0, 0)\}) \times \mathbb{R}_{\geq 0}$ and a KKL observer can be designed with $(A, B) = (A_{\lambda, m}, B_m) \in \mathbb{R}^{m \times m} \times \mathbb{R}^m$ given in (13) for $m = 4$.

Remark 3.12 In practice, T_0 satisfying (25) is usually computed numerically by means of neural networks approaches (see Example 1.2), hence is not necessarily analytic with respect to x . Actually, even if the neural network is analytic, the learned map is only an approximation of the real map T_0 , which thus may not be analytic.

Example 3.13 It has been shown in Example 3.2 that the Wilson-Cowan model (7) is (\mathbb{R}^n, δ_d) -backward-distinguishable. Besides, when S is bounded, $\sigma^-(x) = -\infty$ for all $x \in \mathbb{R}^n$ so that \mathbb{R}^n is backward-invariant. However, h is not bounded on \mathbb{R}^n and T_0 cannot be checked to be analytic as explained in the previous remark. Nevertheless, numerical simulations show that the result of Theorem 3.8 holds with $\lambda = -2$ and $m = 5$, see [11]. Note that, since (7) is actually (\mathcal{O}, δ_d) -backward-distinguishable for any open set \mathcal{O} , we could pick \mathcal{O} bounded containing $\text{cl}(\mathcal{X})$ and modify f outside of $\mathcal{O} + \delta_d$ in order to make a compact set backward-invariant, as explained in the end of Section 2. h would then be bounded on \mathcal{O} , but it is still not straightforward to check that T_0 is analytic.

4 A cascaded design procedure for T

Consider now a dynamical system in the cascade form

$$\dot{x} = f(x), \quad \dot{\xi} = F\xi + Gh(x) \quad (22)$$

with $(x, \xi) \in \mathbb{R}^n \times \mathbb{R}^{n_\xi}$ and output $y_\xi = H\xi$ and with (F, G, H) in the normal controllability form

$$F = \begin{pmatrix} 0 & 1 & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & \cdots & \cdots & 0 & 1 \\ -a_0 & -a_1 & \cdots & \cdots & -a_{n_\xi-1} \end{pmatrix}, \quad G = \begin{pmatrix} 0 \\ \vdots \\ \vdots \\ 0 \\ \gamma \end{pmatrix} \quad (23)$$

$$H = \begin{pmatrix} 1 & 0 & \cdots & \cdots & 0 \end{pmatrix} \quad (24)$$

with $\gamma \neq 0$ and $(a_0, \dots, a_{n_\xi-1})$ in \mathbb{R}^{n_ξ} .

Assuming we know a KKL observer for $\dot{x} = f(x)$ from the output $y = h(x)$,

we would like to deduce an observer for (x, ξ) from the output y_ξ . This covers two cases of practical interest:

- $y = h(x)$ is not available but is used as an intermediary step in the design of an observer for (x, ξ) from the real output y_ξ ;
- $y = h(x)$ is available, but noisy, and we want to use a filtered version y_ξ of y in the KKL observer.

Due to the controllability form of (F, G, H) , if the system $\dot{x} = f(x)$ with output $y = h(x)$ is backward-distinguishable, then, the extended system (22) with state (x, ξ) is also backward-distinguishable. Indeed, intuitively speaking, the past values of y_ξ determine ξ and y uniquely and therefore also x . We could thus use Theorem 3.4 to show the existence of a KKL observer for this extended system. However, the goal of this section is rather to provide an explicit design method when a solution T_0 to the PDE (5) is available for the initial system (f, h) with $A = \lambda$ and $B = 1$ for each λ . More precisely, we exhibit a solution if the following assumption holds for some open set \mathcal{O} containing $\text{cl}(\mathcal{X})$.

Assumption 4.1 *There exist a mapping T_0 in $C^\infty(\mathbb{R}; C^1(\mathbb{R}^n; \mathbb{R}))$ and an open subset $\Theta_0 \subset \mathbb{R}_{<0}$ such that for all (λ, x) in $\Theta_0 \times \mathcal{O}$,*

$$\frac{\partial T_0}{\partial x}(\lambda, x)f(x) = \lambda T_0(\lambda, x) + h(x) , \quad (25)$$

and for all (λ, x_a, x_b) in $\Theta_0 \times \mathcal{O}^2$ verifying $x_a \neq x_b$, there exists $k \geq n_\xi$ in \mathbb{N} such that

$$\frac{\partial^k T_0}{\partial \lambda^k}(\lambda, x_a) - \frac{\partial^k T_0}{\partial \lambda^k}(\lambda, x_b) \neq 0 . \quad (26)$$

Note that if a solution T_0 to (25) is analytic with respect to λ , then Assumption 4.1 holds as long as for any $(x_a, x_b) \in \mathcal{O}^2$ verifying $x_a \neq x_b$, $T_0(\cdot, x_a) - T_0(\cdot, x_b)$ is not a polynomial of degree strictly less than n_ξ . Under the assumption of backward-distinguishability of Theorem 3.4, it is shown in its proof (see Section D.2.2) that such an analytic map T_0 always exists in the form (D.14). Indeed, $T_0(\cdot, x_a) - T_0(\cdot, x_b)$ takes the form of a Laplace transform of some non zero causal signal, hence cannot be a polynomial. For example, the Wilson-Cowan model considered in Example 3.2 is backward distinguishable under some observability assumptions on (C, W) and S , hence satisfies Assumption 4.1 under the same conditions.

The following theorem shows that this assumption is sufficient to give an explicit expression of an injective mapping T , allowing to obtain an observer for the entire system with state (x, ξ) .

Theorem 4.2 *Suppose that Assumption 4.1 holds. Let $\Theta_0^{\text{ext}} = \Theta_0 \setminus \sigma(F)$. Then there exists a zero Lebesgue measure subset $\mathcal{J} \subset (\Theta_0^{\text{ext}})^{2(n+n_\xi)+1}$ such that, for*

each $\lambda_1, \dots, \lambda_{2(n+n_\xi)+1}$ in $(\Theta_0^{\text{ext}})^{2(n+n_\xi)+1} \setminus \mathcal{J}$, the map $T : \mathcal{O} \times \mathbb{R}^{n_\xi} \rightarrow \mathbb{R}^{2(n+n_\xi)+1}$ defined by

$$\begin{aligned} T(x, \xi) &= (T_0^{\text{ext}}(\lambda_1, x, \xi), \dots, T_0^{\text{ext}}(\lambda_{2(n+n_\xi)+1}, x, \xi)), \\ T_0^{\text{ext}}(\lambda, x, \xi) &= H(\lambda I - F)^{-1}(GT_0(\lambda, x) - \xi) \end{aligned}$$

is injective and verifies

$$\frac{\partial T}{\partial(x, \xi)}(x, \xi) \begin{bmatrix} f(x) \\ F\xi + Gh(x) \end{bmatrix} = AT(x, \xi) + BH\xi, \quad (27)$$

with

$$A = \text{diag}(\lambda_1, \dots, \lambda_{2(n+n_\xi)+1}), \quad B = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^\top. \quad (28)$$

Proof : See Appendix F. □

Remark 4.3 Theorem 4.2 could be extended to time varying systems. However, the filter needs to remain autonomous. More precisely, instead of system (22), we could consider

$$\dot{x} = f(x, t), \quad \dot{\xi} = F\xi + Gh(x, t) \quad (29)$$

with $(x, \xi, t) \in \mathbb{R}^n \times \mathbb{R}^{n_\xi} \times \mathbb{R}$ and same output as before. In that case, the result holds by increasing by 1 the dimension of the observer.

Example 4.4 Consider again an harmonic oscillator with unknown frequency given in (18), but this time with a simple filter in the form

$$\dot{\xi} = -a\xi + y \quad (30)$$

with $a > 0$. Note that the function T_0 given in (20) is solution of (25), is analytic and satisfies Assumption 4.1. With the former theorem, we know that for almost all 9 negative real numbers λ_i (different from $-a$) the system

$$(\hat{x}, \hat{\xi}) = T^{\text{inv}}(z_1, \dots, z_9), \quad \dot{z}_i = \lambda_i z_i + \xi,$$

where T^{inv} is any continuous function which satisfies

$$T^{\text{inv}}(T_0^{\text{ext}}(\lambda_1, x, \xi), \dots, T_0^{\text{ext}}(\lambda_9, x, \xi)) = (x, \xi)$$

where $T_0^{\text{ext}}(\lambda, x, \xi) = \frac{1}{\lambda+a} \left[\frac{-\lambda x_1 - x_b}{\lambda^2 + x_3} - \xi \right]$, is an observer.

In fact, on this example, 9 different eigenvalues are not required to get injectivity of the mapping $(x, \xi) \mapsto (T_0^{\text{ext}}(\lambda_1, x, \xi), \dots, T_0^{\text{ext}}(\lambda_9, x, \xi))$. Indeed, we have

for all (x_a, x_b) in \mathcal{O} defined in (19) and all λ

$$T_0(\lambda, x_a) - T_0(\lambda, x_b) = \frac{1}{(\lambda^2 + x_{a,3})(\lambda^2 + x_{b,3})} \begin{bmatrix} 1 & \lambda & \lambda^2 & \lambda^3 \end{bmatrix} v(x_a, x_b)$$

with

$$v(x_a, x_b) = \begin{bmatrix} x_{a,3}x_{b,2} - x_{b,3}x_{a,2} \\ x_{a,3}x_{b,1} - x_{b,3}x_{a,1} \\ x_{b,2} - x_{a,2} \\ x_{b,1} - x_{a,1} \end{bmatrix}$$

It yields for all $\lambda_i \neq -a$, $i = 1, \dots, 5$, for all (x_a, ξ_a, x_b, ξ_b) in $(\mathcal{O} \times \mathbb{R})^2$,

$$\begin{aligned} T(x_a, \xi_a) - T(x_b, \xi_b) &= \begin{bmatrix} T_0^{\text{ext}}(\lambda_1, x_a, \xi_a) - T_0^{\text{ext}}(\lambda_1, x_b, \xi_b) \\ \vdots \\ T_0^{\text{ext}}(\lambda_5, x_a, \xi_a) - T_0^{\text{ext}}(\lambda_5, x_b, \xi_b) \end{bmatrix} \\ &= \mathfrak{D}(\lambda_1, \dots, \lambda_5) \mathfrak{V}(\lambda_1, \dots, \lambda_5) \begin{bmatrix} v(x_a, x_b) + w(x_a, x_b)(\xi_b - \xi_a) \\ \xi_b - \xi_a \end{bmatrix}. \end{aligned}$$

where \mathfrak{V} is the Vandermonde matrix

$$\mathfrak{V}(\lambda_1, \dots, \lambda_5) = \begin{bmatrix} 1 & \lambda_1 & \lambda_1^2 & \lambda_1^3 & \lambda_1^4 \\ \vdots & & & & \\ 1 & \lambda_5 & \lambda_5^2 & \lambda_5^3 & \lambda_5^4 \end{bmatrix},$$

which is invertible as soon as the λ_i 's are all different and

$$\mathfrak{D}(\lambda_1, \dots, \lambda_5) = \text{diag} \left\{ \frac{1}{(\lambda_i + a)(\lambda_i^2 + x_{a,3})(\lambda_i^2 + x_{b,3})} \right\}.$$

is also invertible and well defined for (x_a, ξ_a, x_b, ξ_b) in $(\mathcal{O} \times \mathbb{R})^2$. Note that injectivity of the mapping T is obtained since from the former expression $T(x_a, \xi_a) - T(x_b, \xi_b) = 0$ implies that $\xi_b = \xi_a$ and $v(x_a, x_b) = 0$. Moreover, for (x_a, x_b) in \mathcal{O}^2 , $v(x_a, x_b) = 0$ implies $x_a = x_b$.

Example 4.5 Consider the Wilson-Cowan model (7) and assume either 1) the available measurement is actually an averaged neural activity ξ modeled by filtering y ; or 2) the measurement of y is available but very noisy and we would like to filter y before using it in the observer. In both scenarios, the model takes the form (22) and we assume for instance $\dim \xi = 1$ and $a_0 = \gamma = 1$. We proceed as follows: first, we compute an approximate solution T_a to (5) for the

nominal dynamics (7) with diagonal pair (A, B) given in (28) with eigenvalues $(-2, -3, -4, -5, -6, -7)$ following the process of Example 1.2; second, we compute T_a^{ext} solution to the extended PDE (27) according to the explicit formula given by Theorem 4.2; then, we generate data points (x, ξ) by simulating (22) on $[0, 30]$ from 100 initial conditions in $[-1, 1]^5$; computing $z = T_a^{\text{ext}}(x, \xi)$ finally provides a data set of pairs (x, ξ, z) , from which a numerical left-inverse T_a^{ext} of T_a^{ext} can be learned. Theorem 4.2 thus avoids the learning step of T_a^{ext} by directly exploiting T_a learned in smaller dimension. Results of simulations are provided in Figure A.3, A.4 in the scenario 2) where the measurement y_m of y is noisy to illustrate the filtering capabilities of the KKL observer. Note that the colored white noise ν is obtained from a first order filter of a uniform white noise of amplitude 5 with eigenvalue 100.*

5 Conclusion

In this paper, we pursue the theoretical study on the nonlinear version of the observers initially introduced by Luenberger in his seminal paper of 1964 [25]. In a first part, we refine some existing results obtained previously by relaxing constraints on the structure of the observer and showing that genericity is achieved with an algorithm of (real) dimension $2n + 1$ and not of dimension $2n + 2$ as initially indicated in [4]. Also, three other results have been established allowing to improve methods to design these observers. The first one, is related to the uniqueness of the immersion on which is based the design. This allows to certify injectivity (and consequently observer convergence) if one succeeds in finding the immersion. A second result shows that it is possible to restrict ourselves to mono-parametric dynamics while preserving the hope of the convergence of the observer. Finally, we give a method to design filtered versions of these observers which may improve the behavior of the estimate in presence of measurement noise.

We believe that these theoretical results will allow an improvement of the learning techniques for the synthesis of KKL observers which are found in many current studies (see for instance [35,14,13,21,32,31]). In particular, existence results for larger classes of pairs (A, B) without constraints on their structure paves the way towards the optimization of this pair for better performance. For that, further research is needed to understand how to quantify the impact of the pair (A, B) on the considered performance.

Finally, although this paper is about state estimation, it is important to note that they will have a direct implication in the field of output regulation. Indeed, these results fit perfectly into the synthesis of the regulators introduced in [27].

The following tabular summarizes the results obtained so far for the existence of a KKL observer in the form (4).

<i>Observers</i>	<i>Assumptions</i>	<i>Dimension</i>	<i>Structure</i>
“Generic” KKL Theorem 3.4	Backward distinguishability	$2n + 1$	Almost all pair (A, B)
“Deep” KKL Theorem 3.8	Backward distinguishability and invariance, analyticity of T	Unknown (large enough)	$A = \begin{pmatrix} \lambda & & & \\ & \lambda & & \\ & & \ddots & \\ & & & 1 & \lambda \end{pmatrix}, B = (1 \ 0 \ \dots \ 0)^\top$
“Generic diagonal” KKL [4, Theorem 3]	Backward distinguishability	$2n + 2$	$A = \text{diag} \begin{pmatrix} \text{Re } \lambda_i & -\text{Im } \lambda_i \\ \text{Im } \lambda_i & \text{Re } \lambda_i \end{pmatrix}_{1 \leq i \leq n+1},$ $B = ((1, 0) \ \dots \ (1, 0))^\top$ for almost all (λ_i) 's in \mathbb{C}^{n+1}
“High-gain” KKL [4, Theorem 4]	Strong differential observability	Order of differential observability	For any (A_0, B) , $A = kA_0$ for k large enough

Acknowledgement

The authors are deeply thankful to Laurent Praly for many fruitful discussions on this topic.

A Results of numerical simulations

In this section, we provide some numerical results concerning Examples 1.2 and 4.5.

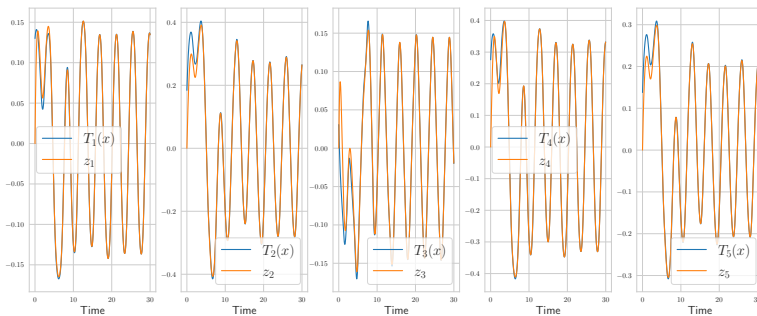


Fig. A.1. Convergence of z to $T_a(x)$, with x solution to (3), z solution to (1) with input $y = h(x)$ and T_a an approximate solution of (5) that has been learned numerically following the method of [35]. See Example 1.2 for details.

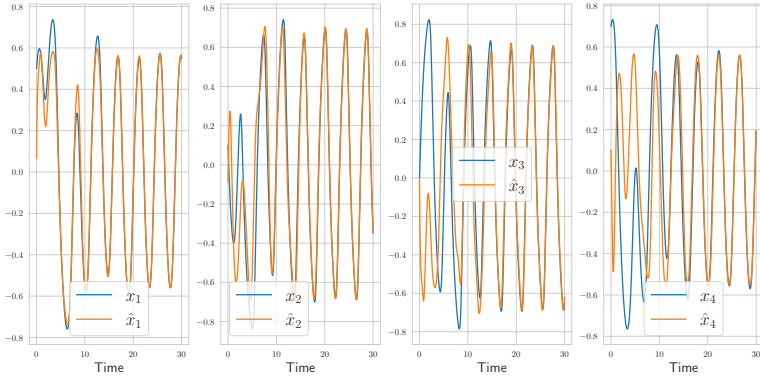


Fig. A.2. Convergence of $\hat{x} = T_a^*(z)$ to x , with x solution to (3) and z solution to (1) with input $y = h(x)$ and a left-inverse T_a^* of T_a has been learned numerically following the method of [35]. See Example 1.2 for details.

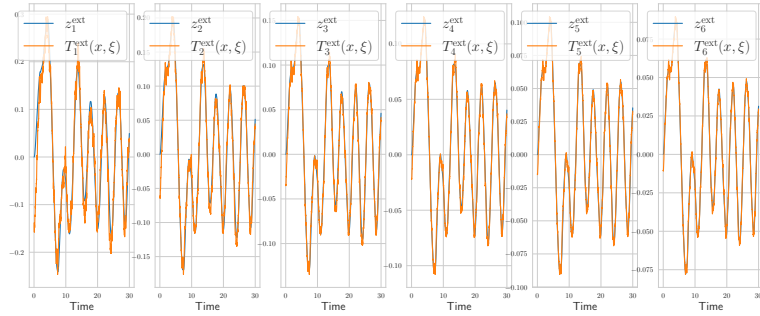


Fig. A.3. Convergence of z to $T_a^{\text{ext}}(x, \xi)$, with (x, ξ) solution to (22), but where ξ is fed with $y_m = h(x) + \nu$ instead of $y = h(x)$ (see Figure A.4), ν a colored white noise, z solution to (1) with input y_ξ and T_a^{ext} an approximate solution of (27). See Example 4.5 for details.

B Proof of Theorem 1.1

Since T is injective on the compact set $\text{cl}(\mathcal{X})$, there exists a continuous map $T^{\text{inv}} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ such that (6) holds (see [28] for instance). According to (5) holding on \mathcal{X} , and because A is Hurwitz,

$$\lim_{t \rightarrow +\infty} |Z(z, x, t) - T(X(x, t))| = 0.$$

Consider $\delta > 0$. Since $X(x, t) \in \mathcal{X}$ for all $t \geq 0$, there exists $\bar{t} \geq 0$ such that for all $t > \bar{t}$, $Z(z, x, t) \in T(\mathcal{X}) + \delta$. Besides, T^{inv} is continuous on the compact set $\text{cl}(T(\mathcal{X}) + \delta)$, so there exists a class- \mathcal{K} map ρ such that

$$|T^{\text{inv}}(z_a) - T^{\text{inv}}(z_b)| \leq \rho(|z_a - z_b|) \quad \forall (z_a, z_b) \in \text{cl}(T(\mathcal{X}) + \delta) \times \text{cl}(T(\mathcal{X}) + \delta).$$

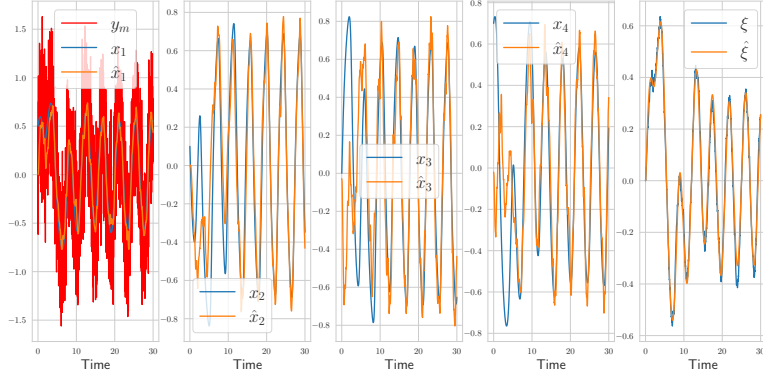


Fig. A.4. Convergence of $(\hat{x}, \hat{\xi}) = T_a^{\text{ext}*}(z)$ to (x, ξ) , with (x, ξ) solution to (22), but where $\hat{\xi}$ is fed with $y_m = h(x) + \nu$ instead of $y = h(x)$, ν a colored white noise, z solution to (1) with input y_ξ and an approximate left-inverse $T_a^{\text{ext}*}$ of T_a^{ext} . See Example 4.5 for details.

Applying this inequality with $z_a = T(X(x, t))$ and $z_b = Z(z, x, t)$ for $t > \bar{t}$ then gives the result using (6).

C Proof of Theorem 2.1

Since \mathcal{O} is backward invariant, we have for all x in \mathcal{O} and all $t \leq 0$,

$$\frac{d}{dt}T_a(X(x, t)) = AT_a(X(x, t)) + Bh(X(x, t)) ,$$

and,

$$\frac{d}{dt}T_b(X(x, t)) = AT_b(X(x, t)) + Bh(X(x, t)) ,$$

which implies

$$T_a(X(x, t)) - T_b(X(x, t)) = \exp(At)[T_a(x) - T_b(x)] , \forall t \leq 0 , \forall x \in \mathcal{O} .$$

Hence,

$$T_a(x) - T_b(x) = \exp(-At)[T_a(X(x, t)) - T_b(X(x, t))] , \forall t \leq 0 , \forall x \in \mathcal{O} .$$

Now make t go to $-\infty$.

In case (i), we get that $\exp(-At)[T_a(X(x, t)) - T_b(X(x, t))]$ tends towards 0 by boundedness of T_a and T_b since A is Hurwitz.

In case (ii), we get that $|X(x, t)| \leq M \exp(-\kappa_f t)$ for all $t \leq 0$ for some $M > 0$,

hence

$$|T_a(X(x, t)) - T_b(X(x, t))| \leq \kappa_a M^{q_a} \exp(-q_a \kappa_f t) + \rho_a + \kappa_b M^{q_b} \exp(-q_b \kappa_f t) + \rho_b .$$

Since $\kappa_f q_a < |\operatorname{Re} \lambda_m(A)|$ and $\kappa_f q_b < |\operatorname{Re} \lambda_m(A)|$, $\exp(-At)[T_a(X(x, t)) - T_b(X(x, t))]$ tends towards 0.

In case (iii), we get by Grönwall's inequality that

$$|T_a(X(x, t))| \leq (|T_a(x)| + \rho_a |t|) \exp(-\kappa_a t) \quad \text{and} \quad |T_b(X(x, t))| \leq (|T_b(x)| + \rho_b |t|) \exp(-\kappa_b t).$$

Hence

$$\begin{aligned} |T_a(X(x, t)) - T_b(X(x, t))| &\leq (|T_a(x)| + \rho_a |t|) \exp(-\kappa_a t) \\ &\quad + (|T_b(x)| + \rho_b |t|) \exp(-\kappa_b t). \end{aligned}$$

Since $\kappa_a < |\operatorname{Re} \lambda_m(A)|$ and $\kappa_b < |\operatorname{Re} \lambda_m(A)|$, $\exp(-At)[T_a(X(x, t)) - T_b(X(x, t))]$ tends towards 0.

Thus, in any case, $T_a(x) = T_b(x)$ for all x in \mathcal{O} .

D Proof of Theorem 3.4

D.1 A proof based on diagonalization

The proof relies on two main ideas:

- almost any matrix A of dimension $2n + 1$ is diagonalizable, with a spectrum decomposed into 2ℓ complex conjugate eigenvalues and $2(n - \ell) + 1$ real eigenvalues for some $\ell \in \{0, \dots, n\}$;
- a generic choice of ℓ complex eigenvalues and $2(n - \ell) + 1$ real eigenvalues for ℓ describing $\{0, \dots, n\}$ yields a generic choice of matrix A of dimension $2n + 1$.

The construction of the zero-measure set \mathcal{J} allowing to prove Theorem 3.4 is thus based on the following preliminary result which investigates the existence of an injective solution to (5) in the case where A is a diagonal Hurwitz matrix with ℓ complex eigenvalues with real part smaller than $-\rho$ and $2(n - \ell) + 1$ real eigenvalues smaller than $-\rho$. In that case, to define the observer state space, for ℓ in $\{0, \dots, n\}$, we introduce

$$\Omega_\ell = \mathbb{C}^\ell \times \mathbb{R}^{2(n-\ell)+1} . \tag{D.1}$$

Also, given a positive real number ρ and ℓ in $\{0, \dots, n\}$, we consider the set $\Omega_{\ell, \rho}$ defined as (see (9))

$$\Omega_{\ell, \rho} = \mathbb{C}_\rho^\ell \times \mathbb{R}_\rho^{2(n-\ell)+1} . \quad (\text{D.2})$$

The following result can be stated.

Proposition D.1 *Assume that system (3) is (\mathcal{O}, δ_d) -backward distinguishable for some open bounded set \mathcal{O} containing $\text{cl}(\mathcal{X})$ and some $\delta_d > 0$. Then there exist a positive real number ρ such that for each ℓ in $\{0, \dots, n\}$, there exists a zero Lebesgue measure subset \mathcal{I}_ℓ of $\Omega_{\ell, \rho}$ such that for each $(\lambda_1, \dots, \lambda_{2n-\ell+1})$ in $\Omega_{\ell, \rho} \setminus \mathcal{I}_\ell$, there exists an injective C^1 function $T_{\text{diag}} : \mathcal{O} \mapsto \Omega_\ell$ verifying (5) with*

$$A_{\text{diag}} = \text{diag}(\lambda_1 \dots, \lambda_{2n-\ell+1}) , \quad B_{\text{diag}} = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}^\top . \quad (\text{D.3})$$

Note that the first ℓ components of the map T_{diag} provided by Proposition D.1 are complex valued. Considering the real and imaginary parts of these components, we obtain a map $T_{\text{real}} : \mathcal{O} \rightarrow \mathbb{R}^{2n+1}$ which is an injective C^1 solution to (5) with

$$A_{\text{real}} = \text{diag}(\Lambda_1, \dots, \Lambda_\ell, \lambda_{\ell+1}, \dots, \lambda_{2n-\ell+1}) \quad , \quad B_{\text{real}} = \begin{bmatrix} B_1 \\ \vdots \\ B_{2n-\ell+1} \end{bmatrix} \quad (\text{D.4})$$

where Λ_i and B_i take the form

$$\Lambda_i = \begin{bmatrix} \text{Re}(\lambda_i) & -\text{Im}(\lambda_i) \\ \text{Im}(\lambda_i) & \text{Re}(\lambda_i) \end{bmatrix} \quad , \quad B_i = \begin{cases} \begin{bmatrix} 1 \\ 0 \end{bmatrix} & \text{for } i \in \{1, \dots, \ell\} \\ 1 & \text{for } i \in \{\ell+1, \dots, 2n-\ell+1\} . \end{cases}$$

The proof of Proposition D.1 is postponed to Section D.2. In the meantime, we prove Theorem 3.4 by translating the generic choice of eigenvalues in Proposition D.1 into a generic choice of the pair (A, B) of dimension $2n+1$.

Lemma D.2 *For ℓ in $\{1, \dots, n\}$, let \mathcal{I}_ℓ be a zero measure subset of Ω_ℓ . The set of matrices in $\mathbb{R}^{(2n+1) \times (2n+1)}$ with characteristic polynomial*

$$\prod_{j=1}^{\ell} (s^2 - 2 \text{Re}(\lambda_j)s + |\lambda_j|^2) \times \prod_{j=2\ell+1}^{2n-\ell+1} (s - \lambda_j)$$

for some $(\lambda_1, \dots, \lambda_{2n-\ell+1}) \in \mathcal{I}_\ell$ is of zero-measure in $\mathbb{R}^{(2n+1) \times (2n+1)}$.

Proof : Let $\mathbb{R}_{2n+1}[s]^{\text{monic}}$ be the set of monic real polynomials with indeterminate s and degree $2n + 1$ (i.e., real polynomials of degree $2n + 1$ in which the nonzero coefficient of highest degree is equal to 1). Consider $\phi : \Omega_{\ell, \rho} \rightarrow \mathbb{R}_{2n+1}[s]^{\text{monic}}$ such that

$$\phi(\lambda_1, \dots, \lambda_{2n-\ell+1}) = \prod_{j=1}^{\ell} (s^2 - 2 \operatorname{Re}(\lambda_j)s + |\lambda_j|^2) \times \prod_{j=2\ell+1}^{2n-\ell+1} (s - \lambda_j).$$

The map ϕ associates to a list of ℓ complex roots and $2(n - \ell) + 1$ real roots, the monic polynomial of degree $2n + 1$ with real coefficients having those roots. By identifying $\mathbb{R}_{2n+1}[s]^{\text{monic}}$ with a list of $2n + 1$ coefficients in \mathbb{R}^{2n+1} , ϕ is C^1 from \mathbb{R}^{2n+1} to \mathbb{R}^{2n+1} . From which, we concludes that $\phi(\mathcal{I}_{\ell})$ is a zero measure subset of $\mathbb{R}_{2n+1}[s]^{\text{monic}}$ assimilated to \mathbb{R}^{2n+1} (see for instance [17, Theorem 3 in §3]). Consider now $\Phi : \mathbb{R}^{(2n+1) \times (2n+1)} \rightarrow \mathbb{R}_{2n+1}[s]^{\text{monic}}$ defined as

$$\Phi(A) = \det(A - sI_{2n+1})$$

This map is C^{∞} (still identifying $\mathbb{R}_{2n+1}[s]^{\text{monic}}$ with \mathbb{R}^{2n+1}) Let us show that it is a submersion almost everywhere. All the coefficients of $\det(A - sI_{2n+1})$ are polynomials of the coefficients of A . It follows that $\frac{\partial \Phi}{\partial A}(A)$ is a rectangular matrix of dimension $(2n + 1) \times (2n + 1)^2$ whose coefficients are polynomials of the coefficients of A . The set of matrices A such that $\operatorname{rank} \frac{\partial \Phi}{\partial A}(A) < 2n + 1$ is characterized by the determinant of each minor being zero, which is thus an algebraic set of zero-measure. Hence Φ is a submersion. With ⁴ [33, Theorem 1], we therefore conclude that the set \mathcal{S}_{ℓ} defined as

$$\mathcal{S}_{\ell} = \{A, \Phi(A) \in \phi(\mathcal{I}_{\ell})\} \tag{D.5}$$

is a zero Lebesgue measure subset of $\mathbb{R}^{(2n+1) \times (2n+1)}$. □

Let ρ be a positive real number, and for each ℓ in $\{0, \dots, n\}$, let \mathcal{I}_{ℓ} be a zero measure subset of $\Omega_{\ell, \rho}$ as given by Proposition D.1 and consider the sets

$$\mathcal{J}_{NC} = \{(A, B) \in \mathbb{R}^{(2n+1) \times (2n+1)} \times \mathbb{R}^{2n+1}, (A, B) \text{ is not controllable}\}, \tag{D.6}$$

$$\mathcal{J}_{ND} = \{(A, B) \in \mathbb{R}^{(2n+1) \times (2n+1)} \times \mathbb{R}^{2n+1}, A \text{ is not diagonalizable in } \mathbb{C}\}, \tag{D.7}$$

$$\mathcal{J}_{\ell} = \{(A, B) \in \mathbb{R}^{(2n+1) \times (2n+1)} \times \mathbb{R}^{2n+1}, A \in \mathcal{S}_{\ell}\}. \tag{D.8}$$

⁴ Let $\Phi : U \subset \mathbb{R}^k \rightarrow \mathbb{R}^{k'}$ of class $C^{k-k'+1}$, where $k' \leq k$. Then, the pre-image of any zero-measure set is of zero-measure if and only if Φ is a submersion almost everywhere, i.e.,

$$\operatorname{rank} \frac{\partial \Phi}{\partial x}(x) = k' \quad \text{for almost all } x \in U$$

It is well-known that \mathcal{J}_{NC} and \mathcal{J}_{ND} are of zero-measure. Applying Lemma D.2, we conclude that the set $\mathcal{J} = \mathcal{J}_{NC} \cup \mathcal{J}_{ND} \cup (\bigcup_{\ell=0}^n \mathcal{J}_\ell)$ is of zero Lebesgue measure in $\mathbb{R}^{(2n+1) \times (2n+1)} \times \mathbb{R}^{2n+1}$.

Consider now (A, B) in $\mathbb{R}^{(2n+1) \times (2n+1)} \times \mathbb{R}^{2n+1} \setminus \mathcal{J}$, such that $A + \rho I_{2n+1}$ is Hurwitz. We wish to transform (A, B) into $(A_{\text{real}}, B_{\text{real}})$ defined in (D.4) in order to apply Proposition D.1. The spectrum of A can be decomposed into 2ℓ complex conjugate eigenvalues and $2(n - \ell) + 1$ real eigenvalues for some ℓ in $\{0, \dots, n\}$. By definition of \mathcal{J} , A is diagonalizable in \mathbb{C} , so there exist $(\lambda_1, \dots, \lambda_{2n-\ell+1})$ in $\Omega_{\ell, \rho}$ and an invertible matrix P in $\mathbb{R}^{(2n+1) \times (2n+1)}$ such that

$$A_{\text{real}} = P^{-1}AP$$

with A_{real} defined in (D.4). Let

$$\tilde{B} = P^{-1}B = \begin{bmatrix} \tilde{B}_1 \\ \vdots \\ \tilde{B}_{2n-\ell+1} \end{bmatrix}$$

with

$$\tilde{B}_i = \begin{cases} \begin{bmatrix} \tilde{b}_{i,1} \\ \tilde{b}_{i,2} \end{bmatrix} \in \mathbb{R}^2 & \text{for } i \in \{1, \dots, \ell\} \\ \tilde{b}_i \in \mathbb{R} & \text{for } i \in \{\ell + 1, \dots, 2n - \ell + 1\}. \end{cases}$$

and

$$M = \text{diag}(M_1, \dots, M_{2n-\ell+1})$$

with

$$M_i = \begin{cases} \begin{bmatrix} \tilde{b}_{i,1} & -\tilde{b}_{i,2} \\ \tilde{b}_{i,2} & \tilde{b}_{i,1} \end{bmatrix} & \text{for } i \in \{1, \dots, \ell\} \\ \tilde{b}_i & \text{for } i \in \{\ell + 1, \dots, 2n - \ell + 1\} \end{cases}$$

so that $MA_{\text{real}} = A_{\text{real}}M$ and $MB_{\text{real}} = \tilde{B}$.

Since $A \notin \mathcal{J}_\ell$, the vector $(\lambda_1, \dots, \lambda_{2n-\ell+1})$ with (D.5) is not in \mathcal{I}_ℓ . Hence, according to Proposition D.1, there exists an injective C^1 function $T_{\text{real}} : \mathcal{O} \mapsto \mathbb{R}^{2n+1}$ such that, for all x in \mathcal{X} ,

$$\frac{\partial T_{\text{real}}}{\partial x} f(x) = A_{\text{real}}T_{\text{real}}(x) + B_{\text{real}}h(x). \quad (\text{D.9})$$

with B_{real} as in (D.4).

Finally, let $T : \mathcal{O} \mapsto \mathbb{R}^{2n+1}$ be the mapping $T(x) = PMT_{\text{real}}(x)$. Since the pair (A, B) is controllable, and P invertible, the pair $(A_{\text{real}}, \tilde{B})$ is also controllable. Hence, this yields that for all i , $\tilde{B}_i \neq 0$. Consequently, the matrix M is

invertible. Thus T is injective on \mathcal{O} . Besides, for all x in \mathcal{X} ,

$$\begin{aligned} \frac{\partial T}{\partial x}(x)f(x) &= PMA_{\text{real}}T_{\text{real}}(x) + PMB_{\text{real}}h(x) , \\ &= PA_{\text{real}}P^{-1}PMT_{\text{real}}(x) + P\tilde{B}h(x) , \\ &= AT(x) + Bh(x) . \end{aligned}$$

D.2 Proof of Proposition D.1

D.2.1 Some variations on Coron's lemma

The proof of Proposition D.1 is based on the following lemma.

Lemma D.3 *Let Υ be an open subset of \mathbb{R}^{2n} , and Θ_i, g_i, p_i be such that for all $i \in \{1, \dots, m\}$,*

- *either Θ_i is an open subset of \mathbb{R} , $g_i : \Theta_i \times \Upsilon \rightarrow \mathbb{R}$ is in $C^\infty(\mathbb{R}; C^1(\mathbb{R}^{2n}; \mathbb{R}))$ and $p_i = 1$;*
- *or Θ_i is an open subset of \mathbb{C} , $g_i : \Theta_i \times \Upsilon \rightarrow \mathbb{C}$ is holomorphic with respect to λ and C^1 with respect to x , and $p_i = 2$.*

Then, if $\sum_i p_i \geq 2n + 1$, and if for all $i \in \{1, \dots, m\}$, for all $(\lambda, x) \in \Theta_i \times \Upsilon$, there exists $k_i \in \mathbb{N}$ such that

$$\frac{\partial^{k_i} g_i}{\partial \lambda^{k_i}}(\lambda, x) \neq 0 \tag{D.10}$$

then the following set has zero Lebesgue measure in $\prod_{i=1}^m \Theta_i$:

$$\mathcal{I} = \bigcup_{x \in \Upsilon} \left\{ (\lambda_i)_{i \in \{1, \dots, m\}} \in \prod_{i=1}^m \Theta_i : g_i(\lambda_i, x) = 0 \quad \forall i \in \{1, \dots, m\} \right\} . \tag{D.11}$$

This lemma is an extension of [16, Lemma 3.2] as well as the version given in [4, Lemma 3.2]:

- In those previous versions, the functions g_i were the same for each i but this does not make any significant difference in the proof.
- In [16, Lemma 3.2], the functions g_i are in $C^\infty(\mathbb{R} \times \mathbb{R}^{2n}; \mathbb{R})$ instead of $C^\infty(\mathbb{R}; C^1(\mathbb{R}^{2n}; \mathbb{R}))$ here. This loss of regularity is not a problem. Instead of the Malgrange theorem of [19], we employ the one obtained in [29].
- In [4, Lemma 3.2], the functions g_i are holomorphic with respect to λ and C^1 with respect to x .

Apart from these modifications, the proof follows readily and is based on the fact that \mathcal{I} is contained in the countable union of image sets through C^1

functions taking values in a real submanifold of dimension $2n$ of $\prod_i \Theta_i$ (which is a real manifold of dimension $\sum_i p_i \geq 2n + 1$). Hence the result is obtained from a variation of Sard theorem. The proof is provided in the long version of this paper, see [11, Appendix G] .

Now following [4], the idea to prove Proposition D.1 is first to exhibit a C^1 solution to (5) with $(A_{\text{diag}}, B_{\text{diag}})$ as in (D.3). This solution is parameterized by the (λ_i) 's. With the distinguishability assumption and the use of Lemma D.3, it is then shown that generically this function is injective on \mathcal{O} .

D.2.2 Construction of T_{diag}

Let $\delta_b > \delta_d$ be any positive real number where δ_d is given by backward-distinguishability. Let $\rho = \max_{x \in \mathcal{O} + \delta_b} \left| \frac{\partial \check{f}}{\partial x}(x) \right|$ where $\check{f} = \chi f$ and where $\chi : \mathbb{R}^n \rightarrow \mathbb{R}$ is a C^∞ function such that

$$\chi(x) = \begin{cases} 0, & x \notin \mathcal{O} + \delta_b \\ 1, & x \in \mathcal{O} + \delta_d. \end{cases} \quad (\text{D.12})$$

Fix ℓ in $\{1, \dots, n\}$. For each $(\lambda_1, \dots, \lambda_{2n-\ell+1})$ in $\Omega_{\ell, \rho}$, we can define the mapping $T_{\text{diag}} : \mathcal{O} \mapsto \Omega_{\ell, \rho}$ defined as

$$T_{\text{diag}}(x) = (T_0(\lambda_1, x), \dots, T_0(\lambda_{2n-\ell+1}, x)) \quad (\text{D.13})$$

with $T_0 : \mathbb{C}_\rho \times \mathcal{O} \rightarrow \mathbb{R}$ defined as

$$T_0(\lambda, x) = \int_{-\infty}^0 \exp(-\lambda s) h(\check{X}(x, s)) ds \quad (\text{D.14})$$

where $\check{X} : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}^n$ is the flow of the modified system

$$\dot{x} = \check{f}(x) = \chi(x)f(x). \quad (\text{D.15})$$

To prove Proposition D.1, we need to show that T_{diag} is solution to the PDE (5) and also that it has enough regularity to apply Lemma D.3 to obtain injectivity. First recall the following fact.

Proposition D.4 ([3, Proposition 3.3]) *The function T_{diag} is C^1 and satisfies (5) with $(A_{\text{diag}}, B_{\text{diag}})$ given in (D.3).*

Moreover,

- The map $T_0(\cdot, x)$ is holomorphic on \mathbb{C}_ρ for each $x \in \mathcal{O}$.

- The restriction of $T_0(\cdot, x)$ to \mathbb{R}_ρ is C^∞ (actually, analytic) for each $x \in \mathcal{O}$. Moreover, for all $k \in \mathbb{N}$, $\frac{\partial^k T_0}{\partial \lambda^k}(\lambda, \cdot)$ can be shown to be C^1 for any $\lambda \in \mathbb{R}_\rho$ by following readily the proof of [3, Proposition 3.3].

To prove Proposition D.1, it now remains to show injectivity by applying Lemma D.3.

D.2.3 Injectivity of T_{diag}

Let $\Upsilon = \{(x_a, x_b) \in \mathcal{O}^2, x_a \neq x_b\}$. Let also $\Theta_i = \mathbb{C}_\rho$ for $i \in \{1, \dots, \ell\}$ and $\Theta_i = \mathbb{R}_\rho$ for $i \in \{\ell + 1, \dots, 2n - \ell + 1\}$. Let $g_i : \Theta_i \times \Upsilon \mapsto \mathbb{C}$ for $i \in \{1, \dots, \ell\}$ and $g_i : \Theta_i \times \Upsilon \mapsto \mathbb{R}$ for $i \in \{\ell + 1, \dots, 2n - \ell + 1\}$ be defined by

$$g_i(\lambda, x_a, x_b) = T_0(\lambda, x_a) - T_0(\lambda, x_b) , \quad (\text{D.16})$$

$$= \int_{-\infty}^0 \exp(-(\lambda + \rho)s) \Delta(x_a, x_b, s) ds , \quad (\text{D.17})$$

for all $(x_a, x_b) \in \Upsilon$ and all $\lambda \in \Theta_i$, where

$$\Delta(x_a, x_b, s) = \exp(\rho s) \left[h(\check{X}(x_a, s)) - h(\check{X}(x_b, s)) \right] .$$

By backward-distinguishability, for all (x_a, x_b) in Υ there exists a negative time t in $(\max \{ \sigma_{\mathcal{O}+\delta_d}^-(x_a), \sigma_{\mathcal{O}+\delta_d}^-(x_b) \}, 0]$ such that $h(X(x_a, t)) \neq h(X(x_b, t))$. Moreover, by definition of χ in (D.15), $X(x, s) = \check{X}(x, s)$ for all $x \in \mathcal{O}$ and all $s \in (\sigma_{\mathcal{O}+\delta_d}^-(x), 0]$. It yields that for all (x_a, x_b) in Υ , there exists $s < 0$ such that $\Delta(x_a, x_b, s) \neq 0$.

From there, two cases may be distinguished.

- For $i \in \{1, \dots, \ell\}$, for each $(x_a, x_b) \in \Upsilon$, $g_i(\cdot, x_a, x_b)$ is holomorphic (since $T_0(\cdot, x)$ is holomorphic for each $x \in \mathcal{O}$) and consequently, there exists k_i such that (D.10) is satisfied.
- For $i \in \{\ell + 1, \dots, 2n - \ell + 1\}$, similarly to the proof of the injectivity of the Laplace transform, for all $\lambda \in \mathbb{R}_\rho$, with $u = \exp(s)$, yields

$$g_i(\lambda, x_a, x_b) = \int_0^1 u^{-(\lambda+\rho)-1} \bar{\Delta}(u) du$$

where $\bar{\Delta}$ is a continuous function defined by $\bar{\Delta}(u) = \Delta(x_a, x_b, \ln(u))$ for $u > 0$ and $\bar{\Delta}(0) = 0$. We deduce that $g_i(\cdot, x_a, x_b)$ is not identically zero on \mathbb{R}_ρ . Indeed, otherwise, picking $\lambda = -(j + \rho + 1)$ for each $j \in \mathbb{N}$, we get

$$\int_0^1 u^j \bar{\Delta}(u) du = 0 . \quad (\text{D.18})$$

By Stone-Weierstrass theorem, for each $\epsilon > 0$, there exists a polynomial P_ϵ such that

$$|\bar{\Delta}(u) - P_\epsilon(u)| \leq \epsilon, \quad \forall u \in [0, 1].$$

Moreover, since P_ϵ is a polynomial, (D.18) yields $\int_0^1 \bar{\Delta}(u)P_\epsilon(u)du = 0$. Hence,

$$\int_0^1 \bar{\Delta}(u)^2 du = \int_0^1 \bar{\Delta}(u)(\bar{\Delta}(u) - P_\epsilon(u))du \leq \max_{u \in [0,1]} |\bar{\Delta}(u)| \epsilon.$$

The former inequality being true for all ϵ , it yields that $\bar{\Delta}$ is identically zero on $[0, 1]$, which is a contradiction since Δ is not identically zero. Therefore, $g_i(\cdot, x_a, x_b)$ is not identically zero on \mathbb{R}_ρ . Since moreover g_i is analytic, it yields that there exists k_i such that (D.10) is satisfied.

We can finally apply Lemma D.3, to obtain the set \mathcal{I}_ℓ given in (D.11). By definition of g_i and of \mathcal{I}_ℓ , we conclude that the map T_{diag} defined in (D.13) is injective on \mathcal{O} for any $(\lambda_1, \dots, \lambda_{2n-\ell+1})$ in $\Omega_{\ell,\rho} \setminus \mathcal{I}_\ell$. This concludes the proof.

E Proof of Theorem 3.8

First of all, the set \mathcal{O} being backward invariant for the flow and the mapping h being bounded in \mathcal{O} , this implies that the function

$$S(\lambda, x) = \int_{-\infty}^0 \exp(-\lambda s) h(X(x, s)) ds \quad (\text{E.1})$$

is well defined on $\mathbb{R}_{<0} \times \mathcal{O}$ and such that for all $\lambda < 0$, $S(\lambda, \cdot)$ is bounded and solution to (14) on \mathcal{O} . Moreover, for all $x \in \mathcal{O}$, $S(\cdot, x)$ is analytic on $\mathbb{R}_{<0}$. With Theorem 2.1, it implies that $T_0 = S$ on $\Theta \times \mathcal{O}$ and therefore, for all $\lambda \in \Theta$, $S(\lambda, \cdot)$ is analytic on \mathcal{O} . For m in \mathbb{N} , let

$$T_{\lambda,m}(x) = (T_0(\lambda, x), \dots, T_{m-1}(\lambda, x)) \quad (\text{E.2a})$$

where

$$T_i(\lambda, x) = \frac{\partial^i T_0}{\partial \lambda^i}(\lambda, x), \quad i = \{0, \dots, m-1\}. \quad (\text{E.2b})$$

Since T_0 is C^∞ , for all $(\lambda, x) \in \Theta \times \mathcal{O}$,

$$\frac{\partial T_1}{\partial x}(\lambda, x)f(x) = \frac{\partial^2 T_0}{\partial \lambda \partial x}(\lambda, x)f(x) = \frac{\partial}{\partial \lambda}(\lambda T_0(\lambda, x) + h(x)) = \lambda T_1(\lambda, x) + T_0(\lambda, x)$$

and iteratively for all i

$$\frac{\partial T_i}{\partial x}(\lambda, x)f(x) = \lambda T_i(\lambda, x) + T_{i-1}(\lambda, x),$$

and consequently, $T_{\lambda,m}$ is the (unique) bounded solution of (5) with $(A, B) = (A_{\lambda,m}, B_m) \in \mathbb{R}^{m \times m} \times \mathbb{R}^m$ given in (13).

Let $g : \mathbb{R}_{<0} \times \mathcal{O} \times \mathcal{O}$ be given by

$$g(\lambda, x_a, x_b) = S(\lambda, x_a) - S(\lambda, x_b) . \quad (\text{E.3})$$

Let $\lambda \in \Theta$ and $\mathcal{C} \subset \mathcal{O}$. Let $\mathcal{D}_{\lambda,\ell}$ be the sequence of open sets defined as

$$\mathcal{D}_{\lambda,\ell} = \left\{ (x_a, x_b) \in \mathcal{O}^2, x_a \neq x_b, \frac{\partial^k g}{\partial \lambda^k}(\lambda, x_a, x_b) = 0, k = 0, \dots, \ell \right\} .$$

We will show that there exists m such that $\mathcal{D}_{\lambda,m} \cap (\mathcal{C} \times \mathcal{C}) = \emptyset$ which implies that $T_{\lambda,m}$ is injective in \mathcal{C} . Note that we have $\mathcal{D}_{\lambda,\ell+1} \subset \mathcal{D}_{\lambda,\ell}$. The map $g(\lambda, \cdot, \cdot)$ being analytic since $S = T_0$ on $\Theta \times \mathcal{O}$, $(\mathcal{D}_{\lambda,\ell})_{\ell \in \mathbb{N}}$ is a decreasing sequence of analytic subsets of $\mathcal{O}^2 \subset \mathbb{R}^{2m}$. The ring of analytic functions being Noetherian [30, Corollary 1, p.99], $(\mathcal{D}_{\lambda,\ell})_{\ell \in \mathbb{N}}$ is a stationary sequence in all compact subsets, i.e. there exists m^* in \mathbb{N} such that, for all $m \geq m^*$,

$$\mathcal{D}_{\lambda,m+\ell} \cap (\mathcal{C} \times \mathcal{C}) = \mathcal{D}_{\lambda,m} \cap (\mathcal{C} \times \mathcal{C}) , \quad \forall \ell \in \mathbb{N} .$$

Assume $\mathcal{D}_{\lambda,m} \cap (\mathcal{C} \times \mathcal{C})$ non-empty and take $(x_a, x_b) \in \mathcal{D}_{\lambda,m} \cap (\mathcal{C} \times \mathcal{C})$. We have $\frac{\partial^k g}{\partial \lambda^k}(\lambda, x_a, x_b) = 0$ for all k . Since, moreover $g(\cdot, x_a, x_b)$ is analytic, this implies that $g(\lambda, x_a, x_b) = 0$ for all $\lambda < 0$. On another hand, with (E.1) and by injectivity of the Laplace transform⁵, this implies that $s \mapsto h(X(x_a, s)) - h(X(x_b, s)) = 0$ for s in $(-\infty, 0]$ and $x_a \neq x_b$. This is a contradiction with the observability assumption. This implies that $\mathcal{D}_{\lambda,m} \cap (\mathcal{C} \times \mathcal{C}) = \emptyset$.

F Proof of Theorem 4.2

First, that for all (λ, x, ξ) in $\Theta_0^{\text{ext}} \times \mathcal{O} \times \mathbb{R}^{n_\xi}$

$$\begin{aligned} & \frac{\partial T_0^{\text{ext}}}{\partial x}(\lambda, x, \xi) f(x) + \frac{\partial T_0^{\text{ext}}}{\partial \xi}(\lambda, x, \xi) (F\xi + Gh(x)) \\ &= H(\lambda I - F)^{-1} G \frac{\partial T_0}{\partial x}(\lambda, x) f(x) - H(\lambda I - F)^{-1} (F\xi + Gh(x)) \\ &= \lambda \left(T_0^{\text{ext}}(\lambda, x, \xi) + H(\lambda I - F)^{-1} \xi \right) - H(\lambda I - F)^{-1} F\xi \\ &= \lambda T_0^{\text{ext}}(\lambda, x, \xi) + H\xi. \end{aligned} \quad (\text{F.1})$$

⁵ This fact is recalled and proved in Section D.2.3 by using Stone-Weierstrass theorem.

We next show injectivity of the mapping T built from T_0^{ext} by picking $2(n + n_\xi) + 1$ distinct λ . For that, our aim is to apply Lemma D.3. To do so, let $(\lambda, x_a, x_b, \xi_a, \xi_b)$ in $\Theta_0 \times \mathcal{O}^2 \times \mathbb{R}^{2n_\xi}$ verifying $(x_a, \xi_a) \neq (x_b, \xi_b)$. We have for each $\lambda \notin \sigma(F)$ $H(\lambda I - F)^{-1}G = \frac{\gamma}{d_{n_\xi}(\lambda)}$, $d_{n_\xi}(\lambda) = \lambda^{n_\xi} + \sum_{j=0}^{n_\xi-1} a_j \lambda^j$. Moreover, $H(\lambda I - F)^{-1} = \left(\frac{p_{n_\xi-1}(\lambda)}{d_{n_\xi}(\lambda)} \dots \dots \frac{p_1(\lambda)}{d_{n_\xi}(\lambda)} \frac{1}{d_{n_\xi}(\lambda)} \right)$ where p_j are polynomials of degree j and d_{n_ξ} a polynomial of degree n_ξ . Let us denote (forgetting the dependency in the variables (x_a, ξ_a, x_b, ξ_b))

$$\mathbf{g}(\lambda) = T_0^{\text{ext}}(\lambda, x_a, \xi_a) - T_0^{\text{ext}}(\lambda, x_b, \xi_b), \quad g_0(\lambda) = T_0(\lambda, x_a) - T_0(\lambda, x_b),$$

and $\tilde{\xi} = (\tilde{\xi}_a, \dots, \tilde{\xi}_{n_\xi}) = \xi_a - \xi_b$. Note that

$$\mathbf{g}(\lambda) = \frac{\sum_{j=0}^{n_\xi-1} \tilde{\xi}_j p_j(\lambda) + \gamma g_0(\lambda)}{d_{n_\xi}(\lambda)}. \quad (\text{F.2})$$

This gives⁶

$$\mathbf{g}^{(1)}(\lambda) = \frac{d_{n_\xi}^{(1)}(\lambda)}{d_{n_\xi}(\lambda)} \mathbf{g}(\lambda) + \frac{\sum_{j=1}^{n_\xi-1} \tilde{\xi}_j p_j^{(1)}(\lambda) + \gamma g_0^{(1)}(\lambda)}{d_{n_\xi}(\lambda)}.$$

which more generally gives for all $\ell \in \mathbb{N}$ and some integers $(c_{ir\ell})$

$$\mathbf{g}^{(\ell)}(\lambda) = \sum_{r=0}^{\ell-1} \sum_{i=1}^{\ell} c_{ir\ell} \frac{d_{n_\xi}^{(i)}(\lambda)}{d_{n_\xi}(\lambda)} \mathbf{g}^{(r)}(\lambda) + \frac{\sum_{j=\ell}^{n_\xi-1} \tilde{\xi}_j p_j^{(\ell)}(\lambda) + \gamma g_0^{(\ell)}(\lambda)}{d_{n_\xi}(\lambda)}. \quad (\text{F.3})$$

The former expression gives for $\ell \geq n_\xi$:

$$\mathbf{g}^{(\ell)}(\lambda) = \sum_{r=0}^{\ell-1} \sum_{i=1}^{\ell} c_{ir\ell} \frac{d_{n_\xi}^{(i)}(\lambda)}{d_{n_\xi}(\lambda)} \mathbf{g}^{(r)}(\lambda) + \frac{\gamma}{d_{n_\xi}(\lambda)} g_0^{(\ell)}(\lambda) \quad (\text{F.4})$$

If $x_a \neq x_b$, with Assumption 4.1, there exists k in \mathbb{N} such that

$$\forall i \in \{0, \dots, k-1\}, \quad \frac{\partial^{n_\xi+i} g_0}{\partial \lambda^{n_\xi+i}}(\lambda) = 0 \quad \text{and} \quad \frac{\partial^{n_\xi+k} g_0}{\partial \lambda^{n_\xi+k}}(\lambda) \neq 0. \quad (\text{F.5})$$

Combining (F.4) and (F.5), there exists k in \mathbb{N} such that

$$\frac{\partial^k T_0^{\text{ext}}}{\partial \lambda^k}(\lambda, x_a, \xi_a) - \frac{\partial^k T_0^{\text{ext}}}{\partial \lambda^k}(\lambda, x_b, \xi_b) \neq 0. \quad (\text{F.6})$$

Indeed, otherwise, (F.4) implies that $g_0^{(\ell)}(\lambda) = 0$ for all $\ell \geq n_\xi$ which contradicts (F.5).

⁶ With the notation $\mathbf{g}^{(\ell)}(\lambda) = \frac{\partial^\ell \mathbf{g}}{\partial \lambda^\ell}(\lambda)$.

Otherwise, $x_a = x_b$ and $\xi_a \neq \xi_b$. We thus have $g_0^{(\ell)}(\lambda) = 0$ for all ℓ , $\tilde{\xi} \neq 0$ and (by (F.3))

$$\mathfrak{g}^{(\ell)}(\lambda) = \sum_{r=0}^{\ell-1} \sum_{i=1}^{\ell} c_{ir\ell} \frac{d_{n_\xi}^{(i)}(\lambda)}{d_{n_\xi}(\lambda)} \mathfrak{g}^{(r)}(\lambda) + \frac{\sum_{j=\ell}^{n_\xi-1} \tilde{\xi}_j p_j^{(\ell)}(\lambda)}{d_{n_\xi}(\lambda)} .$$

Again, this implies (F.6) for some $k \in \mathbb{N}$. Indeed, otherwise, $\sum_{j=\ell}^{n_\xi-1} \tilde{\xi}_j p_j^{(\ell)}(\lambda) = 0$ for all ℓ which implies that $\tilde{\xi} = 0$ (since each p_j is of degree j), which contradicts $\xi_a \neq \xi_b$.

To conclude, this implies that for all $(\lambda, x_a, x_b, \xi_a, \xi_b)$ in $\Theta_1 \times \mathcal{O}^2 \times \mathbb{R}^{2n_\xi}$ verifying $(x_a, \xi_a) \neq (x_b, \xi_b)$, there exists k in \mathbb{N} such that (F.6) holds. Applying Lemma D.3 with $\Upsilon = (\mathcal{O} \times \mathbb{R}^{n_\xi})^2$, $\Theta_i = \Theta_0^{\text{ext}}$ and $g_i(\lambda, x_a, \xi_a, x_b, \xi_b) = T_0^{\text{ext}}(\lambda, x_a, \xi_a) - T_0^{\text{ext}}(\lambda, x_b, \xi_b)$ for $i \in \{1, \dots, 2(n + n_\xi) + 1\}$, we obtain the result.

G Proof of Lemma D.3

This part is a reproduction of the proof given in [4] with the small update related to the use of real or complex valued functions. The differences are in blue.

Let $\bar{\Theta} = \prod_i \Theta_i$. Assume that $\sum_i p_i \geq 2n + 1$. The idea of the proof is to show that the set \mathcal{I} is contained in a countable union of sets which have zero Lebesgue measure.

Given $(\epsilon, \underline{\Lambda}, \underline{x})$ in $\mathbb{R}_{>0} \times \bar{\Theta} \times \Upsilon$, we denote by $S_{\epsilon, \underline{\Lambda}, \underline{x}}$ the set :

$$S_{\epsilon, \underline{\Lambda}, \underline{x}} = \bigcup_{x \in \mathcal{B}_\epsilon(\underline{x})} \{ \Lambda \in \mathcal{B}_\epsilon(\underline{\Lambda}) : g_\ell(\lambda_\ell, x) = 0 \quad \forall \ell \in \{1, \dots, m\} \} . \quad (\text{G.1})$$

Assume for the time being that, for each pair $(\underline{\Lambda}, \underline{x})$ in $\Upsilon \times \bar{\Theta}$, we can find a positive real number ϵ and a countable family of C^1 functions $\sigma_i : \mathcal{B}_\epsilon(\underline{x}) \rightarrow \bar{\Theta}$, such that we have :

$$S_{\epsilon, \underline{\Lambda}, \underline{x}} \subset \bigcup_{i \in \mathbb{N}} \sigma_i(\mathcal{B}_\epsilon(\underline{x})) . \quad (\text{G.2})$$

The family $(\mathcal{B}_\epsilon(\underline{\Lambda}) \times \mathcal{B}_\epsilon(\underline{x}))_{(\underline{\Lambda}, \underline{x}) \in \bar{\Theta} \times \Upsilon}$ is a covering of $\bar{\Theta} \times \Upsilon$ by open subsets. From Lindelöf Theorem (see [10, Lemma 4.1] for instance), there exists a countable family $\{(\underline{\Lambda}_j, \underline{x}_j)\}_{j \in \mathbb{N}}$ such that we have :

$$\bar{\Theta} \times \Upsilon \subset \bigcup_{j \in \mathbb{N}} \mathcal{B}_{\epsilon_j}(\underline{\Lambda}_j) \times \mathcal{B}_{\epsilon_j}(\underline{x}_j) ,$$

where ϵ_j denotes the ϵ associated to the pair $(\underline{\Lambda}_j \times \underline{x}_j)$. With (G.2), it follows

that we have :

$$\mathcal{I} \subset \bigcup_{j \in \mathbb{N}} \bigcup_{i \in \mathbb{N}} \sigma_{i,j}(\mathcal{B}_{\epsilon_j}(\underline{x}_j)) ,$$

where $\sigma_{i,j}$ denotes the i th function σ associated with the pair $(\underline{\Lambda}_j, \underline{x}_j)$. The set $\sigma_{i,j}(\mathcal{B}_{\epsilon_j}(\underline{x}_j))$ is the image, contained in $\bar{\Theta}$, a real manifold of dimension $\sum_i p_i \geq 2n + 1$, by a C^1 function of $\mathcal{B}_{\epsilon_j}(\underline{x}_j)$, a real manifold of dimension $2n$. From a variation on Sard's Theorem (see [17, Theorem 3 in §3] for instance), this image $\sigma_{i,j}(\mathcal{B}_{\epsilon_j}(\underline{x}_j))$ has zero Lebesgue measure in $\bar{\Theta}$. So S , being a countable union of such zero Lebesgue measure subsets, has zero Lebesgue measure.

So all we have to do to establish Lemma D.3 is to prove the existence of ε and the functions σ_i satisfying (G.2) for each pair $(\underline{x}, \underline{\Lambda})$ in $\Upsilon \times \bar{\Theta}$. For ε , we consider two cases :

- (1) Consider a pair $(\underline{\Lambda}, \underline{x})$ such that $g_j(\underline{\lambda}_\ell, \underline{x})$ is non zero. By continuity of g_j , we can find a positive real number ϵ such that $g(\underline{\lambda}_\ell, x)$ is also non zero for all $\underline{\Lambda}$ in $\mathcal{B}_\epsilon(\underline{\lambda}) \subset \Theta_i$ and all x in $\mathcal{B}_\epsilon(\underline{x})$. In this case, the set $S_{\epsilon, \underline{\Lambda}, \underline{x}}$ is empty.
- (2) Consider a pair $(\underline{\Lambda}, \underline{x})$ such that $g(\underline{\lambda}_\ell, \underline{x})$ is zero. From the assumption (D.10), for each ℓ , there exists an integer k_ℓ satisfying :

$$\frac{\partial^i g_{j_\ell}}{\partial \lambda^i}(\underline{\lambda}_\ell, \underline{x}) = 0 \quad \forall i \in \{0, \dots, k_\ell - 1\} \quad , \quad \frac{\partial^{k_\ell} g_{j_\ell}}{\partial \lambda^{k_\ell}}(\underline{\lambda}_\ell, \underline{x}) \neq 0 .$$

For each ℓ in $\{1, \dots, m\}$, two cases may be distinguished.

- (a) If $\Theta_\ell = \mathbb{C}$, following the Weierstrass Preparation Theorem (see [19, Theorem IV.1.1]⁷ for instance), we know the existence of a positive real number ϵ_ℓ , a function $q_\ell : \mathcal{B}_{\epsilon_\ell}(\underline{\lambda}_\ell) \times \mathcal{B}_{\epsilon_\ell}(\underline{x}) \rightarrow \mathbb{C}$, and k_ℓ C^1 functions $a_j^\ell : \mathbb{R}^{2n} \rightarrow \mathbb{C}$ satisfying, for all (λ, x) in $B_{\epsilon_\ell}(\underline{\lambda}_\ell) \times \mathcal{B}_{\epsilon_\ell}(\underline{x})$,

$$q_\ell(\lambda, x) g_\ell(\lambda, x) = (\lambda - \underline{\lambda}_\ell)^{k_\ell} + \sum_{j=0}^{k_\ell-1} a_j^\ell(x) (\lambda - \underline{\lambda}_\ell)^j . \quad (\text{G.3})$$

- (b) If $\Theta_\ell = \mathbb{R}$, following the Malgrange Preparation Theorem [29], we know the existence of a positive real number ϵ_ℓ , a function $q_\ell : \mathcal{B}_{\epsilon_\ell}(\underline{\lambda}_\ell) \times \mathcal{B}_{\epsilon_\ell}(\underline{x}) \rightarrow \mathbb{R}$, and k_ℓ C^1 functions $a_j^\ell : \mathbb{R}^{2n} \rightarrow \mathbb{R}$ satisfying, for all (λ, x) in $B_{\epsilon_\ell}(\underline{\lambda}_\ell) \times \mathcal{B}_{\epsilon_\ell}(\underline{x})$,

$$q_\ell(\lambda, x) g_\ell(\lambda, x) = (\lambda - \underline{\lambda}_\ell)^{k_\ell} + \sum_{j=0}^{k_\ell-1} a_j^\ell(x) (\lambda - \underline{\lambda}_\ell)^j . \quad (\text{G.4})$$

⁷ In [19, Theorem IV.1.1], this theorem is stated with the assumption that g_ℓ is holomorphic in both λ and x . However, as far as x is concerned, it can be seen in the proof of this Theorem that we need only the implicit function theorem to apply. So continuous differentiability in x for each λ is enough.

We choose the real number ϵ , to be associated to $(\underline{\Lambda}, \underline{x})$ in the definition of $S_{\epsilon, \underline{\Lambda}, \underline{x}}$, as :

$$\epsilon = \inf_{\ell \in \{1, \dots, m\}} \epsilon_\ell .$$

In the following $P_\ell : \Theta_\ell \times \mathcal{B}_\epsilon(\underline{x}) \rightarrow \Theta_\ell$ and $a^\ell : \mathcal{B}_\epsilon(\underline{x}) \rightarrow \Theta^{k_\ell}$ denote :

$$P_\ell(\lambda, x) = (\lambda - \underline{\lambda}_\ell)^{k_\ell} + \sum_{j=0}^{k_\ell-1} a_j^\ell(x) (\lambda - \underline{\lambda}_\ell)^j \quad , \quad a^\ell(x) = (a_0^\ell(x), \dots, a_{k_\ell-1}^\ell(x)) .$$

With this definition of ϵ , we have the following implication, for Λ in $\mathcal{B}_\epsilon(\underline{\Lambda})$ and x in $\mathcal{B}_\epsilon(\underline{x})$,

$$g(\lambda_\ell, x) = 0 \quad \forall \ell \in \{1, \dots, m\} \quad \Rightarrow \quad (\lambda_\ell, a^\ell(x)) \in M^\ell \quad \forall \ell \in \{1, \dots, m\} \quad (\text{G.5})$$

where M^ℓ is the set :

$$M^\ell = \left\{ (\lambda, (b_0, \dots, b_{k_\ell-1})) \in \Theta_\ell \times \Theta_\ell^{k_\ell} : (\lambda - \underline{\lambda}_\ell)^{k_\ell} + \sum_{j=0}^{k_\ell-1} b_j (\lambda - \underline{\lambda}_\ell)^j = 0 \right\} \quad (\text{G.6})$$

Our interest in this set follows from the following Lemma, which follows directly from [15, Equations (3.51)-(3.52)] (for $\Theta = \mathbb{R}$) and [4, Lemma 2] (for $\Theta = \mathbb{C}$).

Lemma G.1 *Let M be the set defined as :*

$$M = \left\{ (\lambda, b_0, \dots, b_{k-1}) \in \Theta \times \Theta^k : \lambda^k + \sum_{j=0}^{k-1} b_j \lambda^j = 0 \right\} .$$

where $\Theta = \mathbb{C}$ or $\Theta = \mathbb{R}$. There exists a countable family $\{M_m\}_{m \in \mathbb{N}}$ of regular submanifolds of Θ^k and a countable family of C^1 functions $\rho_m : M_m \rightarrow \Theta$ such that we have the inclusion :

$$M \subset \bigcup_{m \in \mathbb{N}} \bigcup_{b \in M_m} \{(\rho_m(b), b)\} . \quad (\text{G.7})$$

So, for each ℓ in $\{1, \dots, n+1\}$ we have a countable family $\{M_{m_\ell}^\ell\}_{m_\ell \in \mathbb{N}}$ of regular submanifolds of \mathbb{C}^{k_ℓ} and a countable family of C^1 functions $\rho_{m_\ell}^\ell : M_{m_\ell}^\ell \rightarrow \mathbb{C}$ such that, for each x in $\mathcal{B}_\epsilon(\underline{x})$, if $P_\ell(\lambda_\ell, x)$ is zero, then there exists an integer m_ℓ such that we have :

$$a^\ell(x) \in M_{m_\ell}^\ell \quad , \quad \lambda_\ell = \rho_{m_\ell}^\ell(a^\ell(x)) . \quad (\text{G.8})$$

Hence, with (G.5), if :

$$g(\lambda_\ell, x) = 0 \quad \forall \ell \in \{1, \dots, m\} ,$$

then there exists an m -tuple $\mu = (m_1, \dots, m_m)$ of integers satisfying :

$$a^\ell(x) \in M_{m_\ell}^\ell, \quad \lambda_\ell = \rho_{m_\ell}^\ell(a^\ell(x)) \quad \forall \ell \in \{1, \dots, m\}.$$

So, by letting :

$$S_{\epsilon, \underline{\Delta}, \underline{x}}^\mu = \bigcup_{\{x \in \mathcal{B}_\epsilon(\underline{x}) : a^\ell(x) \in M_{m_\ell}^\ell \quad \forall \ell \in \{1, \dots, n+1\}\}} \{(\rho_{m_1}^1(a^1(x)), \dots, \rho_{m_m}^m(a^m(x)))\} \quad (\text{G.9})$$

we have established :

$$S_{\epsilon, \underline{\Delta}, \underline{x}} \subset \bigcup_{\mu \in \mathbb{N}^m} S_{\epsilon, \underline{\Delta}, \underline{x}}^\mu. \quad (\text{G.10})$$

Comparing (G.2) with (G.10) and using the definition (G.9), we see that a candidate for the function σ_i is :

$$\sigma_i(x) = \left(\rho_{m_\ell}^\ell \left(R_{M_{m_\ell}^\ell} (a^\ell(x)) \right) \right)_{\ell \in \{1, \dots, m\}}$$

where i happens to be the m -tuple μ and $R_{M_{m_\ell}^\ell} : \Theta_\ell^{k_\ell} \rightarrow M_{m_\ell}^\ell$ is a ‘‘restriction’’ to $M_{m_\ell}^\ell$ since we have to consider only those $a^\ell(x)$ which are in $M_{m_\ell}^\ell$. Finding such functions $R_{M_{m_\ell}^\ell}$ such that σ_i is C^1 may not be possible. But, following [16, Lemma 3.3], we know the existence, for each ℓ , of a countable family of C^1 functions $R_\nu^\ell : \Theta_\ell^{k_\ell} \rightarrow M_{m_\ell}^\ell$ such that we get :

$$S_{\epsilon, \underline{\Delta}, \underline{x}}^\mu \subset \bigcup_{\nu \in \mathbb{N}} \left\{ \left(\rho_{m_\ell}^\ell \left(R_\nu^\ell (a^\ell(\mathcal{B}_\epsilon(\underline{x}))) \right) \right)_{\ell \in \{1, \dots, n+1\}} \right\}.$$

In other words the family of functions σ_i is actually given by the family :

$$\sigma_{\mu, \nu} = \left(\rho_{m_\ell}^\ell \circ R_\nu^\ell \circ a_\ell \right)_{\ell \in \{1, \dots, n+1\}}$$

i.e. we have :

$$S_{\epsilon, \underline{\Delta}, \underline{x}} \subset \bigcup_{\mu \in \mathbb{N}^{n+1}} \bigcup_{\nu \in \mathbb{N}} \sigma_{\mu, \nu}(\mathcal{B}_\epsilon(\underline{x})).$$

H Proof that $T_{\lambda, 4}$ is injective

In this section we show that the mapping $T_{\lambda, 4}$ defined in (21) is injective. Indeed, consider x_a and x_b in $(\mathbb{R}^2 \setminus \{(0, 0)\}) \times \mathbb{R}_{\geq 0}$. Denote $w_a = \lambda^2 + x_{a,3}$ and $w_b = \lambda^2 + x_{b,3}$. We have

$$\begin{aligned} T_0(\lambda, x_a) = T_0(\lambda, x_b) =: z_1 &\iff \lambda(x_{b,1}w_a - x_{a,1}w_b) = x_{a,2}w_b - x_{b,2}w_a \\ \frac{\partial T_0}{\partial \lambda}(\lambda, x_a) = \frac{\partial T_0}{\partial \lambda}(\lambda, x_b) =: z_2 &\iff x_{b,1}w_a - x_{a,1}w_b = 2\lambda z_1(w_b - w_a) \end{aligned}$$

which thus gives

$$\begin{cases} T_0(\lambda, x_a) = T_0(\lambda, x_b) \\ \frac{\partial T_0}{\partial \lambda}(\lambda, x_a) = \frac{\partial T_0}{\partial \lambda}(\lambda, x_b) \end{cases} \iff \begin{cases} x_{b,1}w_a - x_{a,1}w_b = 2\lambda z_1(w_b - w_a) \\ x_{a,2}w_b - x_{b,2}w_a = 2\lambda^2 z_1(w_b - w_a) \end{cases}$$

Then, continuing differentiating,

$$\begin{aligned} \frac{\partial^2 T_0}{\partial \lambda^2}(\lambda, x_a) = \frac{\partial^2 T_0}{\partial \lambda^2}(\lambda, x_b) &\iff \frac{2}{w_a} \left[-2\lambda \frac{\partial T_0}{\partial \lambda}(\lambda, x_a) - T(x_a) \right] \\ &= \frac{2}{w_b} \left[-2\lambda \frac{\partial T_0}{\partial \lambda}(\lambda, x_b) - T(x_b) \right] \end{aligned}$$

and

$$\begin{aligned} \frac{\partial^3 T_0}{\partial \lambda^3}(\lambda, x_a) &= \frac{\partial^3 T_0}{\partial \lambda^3}(\lambda, x_b) \\ \iff \frac{-4\lambda}{w_a} \left[-2\lambda \frac{\partial T_0}{\partial \lambda}(\lambda, x_a) - T(x_a) \right] + \frac{2}{w_a} \left[-3 \frac{\partial T_0}{\partial \lambda}(\lambda, x_a) - 2\lambda \frac{\partial^2 T_0}{\partial \lambda^2}(\lambda, x_a) \right] \\ &= \frac{-4\lambda}{w_b} \left[-2\lambda \frac{\partial T_0}{\partial \lambda}(\lambda, x_b) - T(x_b) \right] + \frac{2}{w_b} \left[-3 \frac{\partial T_0}{\partial \lambda}(\lambda, x_b) - 2\lambda \frac{\partial^2 T_0}{\partial \lambda^2}(\lambda, x_b) \right] \end{aligned}$$

Assume therefore that $T_{\lambda,4}(x_a) = T_{\lambda,4}(x_b)$, namely

$$\begin{aligned} \left(T_0(\lambda, x_a), \frac{\partial T_0}{\partial \lambda}(\lambda, x_a), \frac{\partial^2 T_0}{\partial \lambda^2}(\lambda, x_a), \frac{\partial^3 T_0}{\partial \lambda^3}(\lambda, x_a) \right) \\ = \left(T_0(\lambda, x_b), \frac{\partial T_0}{\partial \lambda}(\lambda, x_b), \frac{\partial^2 T_0}{\partial \lambda^2}(\lambda, x_b), \frac{\partial^3 T_0}{\partial \lambda^3}(\lambda, x_b) \right) \end{aligned}$$

which we denote (z_1, z_2, z_3, z_4) . Then, we get two cases :

- either $-2\lambda z_2 - z_1 \neq 0$ and we get $w_a = w_b$ from the third equality, and then $x_a = x_b$ from the first two;
- or $-2\lambda z_2 - z_1 = 0$ and thus $z_3 = 0$, so that the fourth equality provides $z_2 \left(\frac{1}{w_a} - \frac{1}{w_b} \right) = 0$. So either $z_2 \neq 0$ and we recover $w_a = w_b$ and conclude as above; or $z_2 = 0$, but then also $z_1 = 0$, and necessarily $x_{a,1} = x_{a,2} = x_{b,1} = x_{b,2} = 0$, which is impossible.

We conclude that $T_{\lambda,4}$ is injective on $(\mathbb{R}^2 \setminus \{(0, 0)\}) \times \mathbb{R}_+$.

References

- [1] D. Aeyels. Generic observability of differentiable systems. *SIAM Journal on Control and Optimization*, 19(5):595–603, 1981.

- [2] C. Afri, V. Andrieu, L. Bako, and P. Dufour. State and parameter estimation: A nonlinear luenberger observer approach. *IEEE Transactions on Automatic Control*, 62(2):973–980, 2017.
- [3] V. Andrieu. Convergence speed of nonlinear luenberger observers. *SIAM Journal on Control and Optimization*, 52(5):2831–2856, 2014.
- [4] V. Andrieu and L. Praly. On the existence of a kazantzis–kravaris/luenberger observer. *SIAM Journal on Control and Optimization*, 45:432–456, 02 2006.
- [5] P. Bernard and V. Andrieu. Luenberger observers for nonautonomous nonlinear systems. *IEEE Transactions on Automatic Control*, 69:270–281, 2019.
- [6] P. Bernard, V. Andrieu, and D. Astolfi. Observer design for continuous-time dynamical systems. *Annual Reviews in Control*, 53:224–248, 2022.
- [7] P. Bernard, T. Devos, A. Jebai, P. Martin, and L. Praly. Kkl observer design for sensorless induction motors. *To be presented at the Conference on Decision and Control (CDC)*, 2022.
- [8] P. Bernard and L. Praly. Estimation of position and resistance of a sensorless PMSM : a nonlinear luenberger approach for a non-observable system. *IEEE Transactions on Automatic Control*, 66:481–496, 2021.
- [9] M. Bin, P. Bernard, and L. Marconi. Approximate nonlinear regulation via identification-based adaptive internal models. *IEEE Transactions on Automatic Control*, 2020.
- [10] W. M. Boothby and W. M. Boothby. *An introduction to differentiable manifolds and Riemannian geometry, Revised*, volume 120. Gulf Professional Publishing, 2003.
- [11] L. Brivadis, V. Andrieu, P. Bernard, and U. Serres. Further remarks on KKL observers. Submitted at SCL, available at <https://hal-mines-paristech.archives-ouvertes.fr/hal-03695863/>, 2022.
- [12] L. Brivadis, V. Andrieu, and U. Serres. Luenberger observers for discrete-time nonlinear systems. *IEEE Conference on Decision and Control*, 2019.
- [13] M. Buisson-Fenet, L. Bahr, and F. D. Meglio. Learning to observe : neural network-based KKL observers. Python toolbox available at https://github.com/Centre-automatique-et-systemes/learn_observe_KKL.git, 2022.
- [14] M. Buisson-Fenet, L. Bahr, and F. D. Meglio. Towards gain tuning for numerical kkl observers. Available at <https://arxiv.org/abs/2204.00318>, 2022.
- [15] J. Coron. *Control and Nonlinearity*. Mathematical surveys and monographs. American Mathematical Society, 2007.
- [16] J.-M. Coron. On the stabilization of controllable and observable systems by an output feedback law. *Mathematics of Control, Signals and Systems*, 7(3):187–216, 1994.

- [17] G. De Rham. *Differentiable Manifolds: Forms, Currents, Harmonic Forms*, volume 266. Springer Berlin, Heidelberg, 1984.
- [18] J.-P. Gauthier and I. Kupka. *Deterministic observation theory and applications*. Cambridge university press, 2001.
- [19] M. Golubitsky and V. Guillemin. *Stable mappings and their singularities*, volume 14. Springer Science & Business Media, 2012.
- [20] N. Henwood, J. Malaizé, and L. Praly. A robust nonlinear Luenberger observer for the sensorless control of SM-PMSM : Rotor position and magnets flux estimation. *IECON Conference on IEEE Industrial Electronics Society*, 2012.
- [21] S. Janny, V. Andrieu, M. Nadri, and C. Wolf. Deep kkl: Data-driven output prediction for non-linear systems. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pages 4376–4381. IEEE, 2021.
- [22] N. Kazantzis and C. Kravaris. Nonlinear observer design using lyapunov’s auxiliary theorem. *Systems & Control Letters*, 34(5):241–247, 1998.
- [23] G. Kreisselmeier and R. Engel. Nonlinear observers for autonomous Lipschitz continuous systems. *IEEE Transactions on Automatic Control*, 48(3), 2003.
- [24] A. J. Krener and M. Xiao. Nonlinear observer design in the siegel domain through coordinate changes. *IFAC Proceedings Volumes*, 34(6):519–524, 2001.
- [25] D. G. Luenberger. Observing the state of a linear system. *IEEE Transactions on Military Electronics*, 8(2):74–80, April 1964.
- [26] L. Marconi and L. Praly. Uniform practical nonlinear output regulation. *IEEE Transactions on Automatic Control*, 53(5):1184–1202, 2008.
- [27] L. Marconi, L. Praly, and A. Isidori. Output stabilization via nonlinear luenberger observers. *SIAM Journal on Control and Optimization*, 45(6):2277–2298, 2007.
- [28] E. J. McShane. Extension of range of functions. *Bull. Amer. Math. Soc.*, 40(12):837–842, 1934.
- [29] P. Michor. The division theorem on banach spaces. *Österrich. Akad. Wiss. Math.- Natur. Kl. Sitzungsber II*, 189:1–18, 1980.
- [30] R. Narasimhan. *Introduction to the Theory of Analytic Spaces*. Springer Berlin Heidelberg, 1966.
- [31] M. U. B. Niazi, J. Cao, X. Sun, A. Das, and K. H. Johansson. Learning-based design of luenberger observers for autonomous nonlinear systems, 2022.
- [32] J. Peralez and M. Nadri. Deep learning-based luenberger observer design for discrete-time nonlinear systems. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pages 4370–4375. IEEE, 2021.
- [33] S. P. Ponomarev. Submersions and preimages of sets of measure zero. *Siberian Mathematical Journal*, 28(1):153–163, 1987.

- [34] L. Praly, A. Isidori, and L. Marconi. A new observer for an unknown harmonic oscillator. In *17th international symposium on mathematical theory of networks and systems*, pages 24–28, 2006.
- [35] L. Ramos, F. Di Meglio, V. Morgenthaler, L. Silva, and P. Bernard. Numerical design of luenberger observers for nonlinear systems. *IEEE Conference on Decision and Control*, pages 5435–5442, 12 2020.
- [36] A. Shoshitaishvili. Singularities for projections of integral manifolds with applications to control and observation problems. *Theory of singularities and its applications*, 1:295, 1990.
- [37] E. D. Sontag. For differential equations with r parameters, $2r+1$ experiments are enough for identification. *Journal of Nonlinear Science*, 12(6):553–583, 2002.
- [38] M. Spirito, P. Bernard, and L. Marconi. On the existence of KKL functional observers. *American Control Conference*, 2022.
- [39] F. Takens. Detecting strange attractors in turbulence. In *Dynamical systems and turbulence, Warwick 1980*, pages 366–381. Springer, 1981.
- [40] H. R. Wilson and J. D. Cowan. Excitatory and inhibitory interactions in localized populations of model neurons. *Biophysical Journal*, 12(1):1–24, 1972.