



**HAL**  
open science

# Optimal control for neural ODE in a long time horizon and applications to the classification and simultaneous controllability problems

Jon Asier Bárcena-Petisco

► **To cite this version:**

Jon Asier Bárcena-Petisco. Optimal control for neural ODE in a long time horizon and applications to the classification and simultaneous controllability problems. 2021. hal-03299270v1

**HAL Id: hal-03299270**

**<https://hal.science/hal-03299270v1>**

Preprint submitted on 26 Jul 2021 (v1), last revised 4 Jan 2024 (v4)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 Optimal control for neural ODE in a long time horizon  
2 and applications to the classification and simultaneous  
3 controllability problems

4 Jon Asier Bárcena-Petisco\*

5 July 26, 2021

6 **Abstract:** We study the optimal control in a long time horizon of neural ordinary differential  
7 equations which are affine or whose activation function is homogeneous. When considering the  
8 classical regularized empirical risk minimization problem we show that, in long time and under  
9 suitable assumptions, the final state of the optimal trajectories has zero training error. We  
10 assume that the data can be interpolated and that the error can be taken to zero with a cost  
11 proportional to the error. These hypotheses are fulfilled in the classification and simultaneous  
12 controllability problems for some relevant activation and loss functions. Our proofs are mainly  
13 constructive combined with reductio ad absurdum: We find that in long time horizon if the  
14 final error is not zero, we can construct a less expensive control which takes the error to zero.  
15 Moreover, we prove that the norm of the optimal control is constant. Finally, we show the  
16 sharpness of our hypotheses by giving an example for which the error of the optimal state, even  
17 if it decays to 0, is strictly positive for any time.

18 **Key words:** data classification, exact controllability, neural ODE, nonlinear systems, optimal  
19 control, simultaneous controllability

20 **AMS subject classification:** 34H05, 49N10, 93B05

21 **Abbreviated title:** Optimal control for neural ODE

22 **Acknowledgements:** This project has received funding from the European Research Council  
23 (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant  
24 agreement NO: 694126-DyCon). I would like to thank Carlos Esteve, Borjan Geshkovski and  
25 Domènec Ruiz i Balet for fruitful discussion. In particular, I would like to thank Carlos Esteve  
26 for proposing Proposition A.1.

---

\*[1] Departamento de Matemáticas, Universidad Autónoma de Madrid, 28049 Madrid, Spain, [2] Chair of Computational Mathematics, Fundación Deusto, Avenida de las Universidades 24, 48007 Bilbao, Basque Country, Spain. <https://orcid.org/0000-0002-6583-866X>, [jon.barcena@uam.es](mailto:jon.barcena@uam.es)

# 1 Introduction

In this paper we study the optimal control of neural ordinary differential equations for a long time horizon. Neural ODE have been used in machine learning in the last five years, a trend started with [Wei17, HR17]. However, they date back to the 90s, when they were already used for the construction of controls (see the survey [Son92]) and when their controllability properties were first studied (see, for example, [Zbi93] and [SS97]). The control systems ruled by neural ODE have considerably better controllability properties than linear control systems. In fact, as pointed out in [RBZ21], for a fixed  $d \in \mathbb{N}$ , if chosen the right neural ODE we can control an arbitrarily large amount of data in  $\mathbb{R}^d$ , whereas in linear systems we can at most control an amount of data equal to the dimension of the control.

The problem under study is the following: Given a dataset  $\mathbf{x} = (x^1, \dots, x^N) \in (\mathbb{R}^d)_*^N$ , for

$$(\mathbb{R}^d)_*^N := \{(x^1, \dots, x^N) \in (\mathbb{R}^d)^N : x^i \neq x^j \quad \forall i, j \in \{1, \dots, N\} : i \neq j\},$$

we seek to take simultaneously the data set to some target points or regions in  $\mathbb{R}^d$  in a given time  $T > 0$ . The distance to those targets is measured with an error function (also known as *loss function*). The control is the minimizer of the risk minimization functional, which provides a balance between a small cost for the control and a small value for the loss function at the final state of the optimal trajectory.

We study the controllability on affine neural networks, which are given by the following equations:

$$\begin{cases} \dot{y}(t) = w(t)\sigma(y(t)) + b(t), \\ y(0) = x, \end{cases} \quad (1.1)$$

for  $x \in \mathbb{R}^d$  the initial value, and  $\sigma : \mathbb{R}^d \mapsto \mathbb{R}^d$  a Lipschitz function, which is called the *activation function*. The functions  $(w, b)$  are the controls and they belong to  $L^2(0, T; \mathcal{U})$ , for  $\mathcal{U}$  defined by:

$$\mathcal{U} := \mathcal{L}(\mathbb{R}^d, \mathbb{R}^d) \times \mathbb{R}^d.$$

If we want to emphasize the dependence of (1.1) to the initial value and the control, we write  $y(\cdot; x, w, b)$ . Similarly, we denote the sequence of solutions of (1.1) for some fixed control  $(w, b)$  and a data set  $\mathbf{x}$  as:

$$y(\cdot; \mathbf{x}, w, b) := (y(\cdot; x^1, w, b), \dots, y(\cdot; x^N, w, b)).$$

Since  $\sigma$  is Lipschitz, (1.1) is well-posed by Cauchy-Lipschitz Theorem.

1 In addition, we also study more complex neural networks, which are given by the equations:

$$\begin{cases} \dot{y}(t) = r(t)\sigma(w(t)y(t) + b(t)), \\ y(0) = x. \end{cases} \quad (1.2)$$

2 Here  $x$  is the initial value and  $(r, w, b)$  is the control, which belongs to  $L^2(0, T; \tilde{\mathcal{U}})$ , for:

$$\tilde{\mathcal{U}} := \mathcal{M} \times \mathcal{L}(\mathbb{R}^d, \mathbb{R}^d) \times \mathbb{R}^d,$$

3 for:

$$\mathcal{M} := \{L \in \mathcal{L}(\mathbb{R}^d, \mathbb{R}^d) : Le_i \in \{1, -1\}, \forall i = 1, \dots, d\}. \quad (1.3)$$

4 We remark that  $r$  models the direction of the flow, which is necessary if  $\sigma$  is positive. The  
 5 intensity of the flow, on the other hand, is modelled by  $(w, b)$ . We assume that the activation  
 6 function  $\sigma$  is Lipschitz and homogeneous in the sense that:

$$\sigma(\lambda x) = \lambda \sigma(x), \quad \forall \lambda > 0, \quad \forall x \in \mathbb{R}^d. \quad (1.4)$$

This includes important *activation functions* such as rectified linear units, which are given by:

$$\sigma(x) = (\max\{x_1, 0\}, \dots, \max\{x_d, 0\}),$$

see [NH10]; parametric rectified units, given by:

$$\sigma(x) = (\alpha x_1 1_{x_1 < 0} + x_1 1_{x_1 > 0}, \dots, \alpha x_d 1_{x_d < 0} + x_d 1_{x_d > 0}),$$

7 see [HZRS15]; and, of course, the identity,  $\sigma(x) = x$ . As in the previous system,  $y(\cdot; x, r, w, b)$   
 8 and  $y(\cdot; \mathbf{x}, r, w, b)$  denote the solutions of (1.2), and (1.2) is well-posed by Cauchy-Lipschitz  
 9 Theorem.

10 As stated in the first paragraph, we study the properties of the optimal control in a long  
 11 time horizon. The main contribution of our paper is that, if the data can be interpolated  
 12 and the error can be corrected with a cost proportional to the current error, we improve the  
 13 asymptotic bound  $\mathcal{O}(1/T)$  for the error of the final state of the optimal trajectory obtained in  
 14 [EGPZ20a] and prove that it is exactly 0 for a sufficiently large time. In fact, in the simulations  
 15 presented in [EGPZ20a, Examples 4.2 and 4.4] the final errors seem to be 0, so we want to  
 16 determine theoretically if, as their simulation suggests, the error is taken exactly to 0. We  
 17 work in an abstract setting, though we give concrete examples of problems that satisfy our  
 18 assumptions, notably the simultaneous controllability and classification problems. As the name  
 19 suggests, in simultaneous controllability we aim to control two or more independent equations  
 20 by applying the same control. The study of simultaneous controllability dates back to [Rus86]  
 21 and [Lio88, Chapter 5], and relevant papers on this topic include [TW00, Mor14, LZ16, WZL17,  
 22 AB19, RBZ21]. As for the classification problem, it is a simplified version of the simultaneous

1 controllability problem, where the objective is to split the data into two sets, for instance,  
 2  $\{x_1 \leq 1\}$  and  $\{x_1 \geq 1\}$ . An additional contribution of our paper is an example of neural ODE  
 3 and loss functions where the error can be taken to 0, but for all time  $T > 0$  the error at time  
 4  $T$  of the optimal trajectories is strictly positive. This illustrates that the results are far from  
 5 trivial.

6 This paper follows a well-established research line that studies the properties of the optimal  
 7 control and trajectories in a long time horizon. This allows, for instance, that when doing  
 8 numerical simulations one might discard local minimum that are not optimal control. The  
 9 pioneer work is [PZ13], where the authors introduce the *turnpike property*, which means that  
 10 when minimizing certain functionals all the optimal trajectories are most of the time near some  
 11 specific state (the *turnpike*) independently of the initial value and the target. Interesting papers  
 12 on this subject include [DGSW14, TZ15, GTZ16, GG18, TZZ18, Zam18, GH19, SPZ19, BP20,  
 13 EKPZ20, EGPZ20b, GSZ21, MRB20, SN20, Tré20, FHS21, WZ21]. Finally, the optimal control  
 14 for neural ODE is also studied in [EYG21], where the authors consider the cost of the  $L^1$ -norm  
 15 of the control instead of the  $L^2$ -norm and obtain that for that norm the optimal control satisfies  
 16 some sparsity properties and that the error of the final state belongs to  $\mathcal{O}(1/T)$ .

## 17 1.1 Optimal trajectories for affine neural ODE

18 As stated in the first part of the introduction, we study the optimal control of a data set ruled  
 19 by a neural ODE. To measure how far the data is from the objective we introduce the *error*  
 20 *function* (also referred in the literature of Data Science as *loss function*)  $E : \mathbb{R}^d \times X \mapsto \mathbb{R}^+$ , for  
 21  $X \subset \mathbb{R}^d$ . We assume that  $X$  contains the data set  $\mathbf{x}$  and that for all  $x \in X$  the function  $E(\cdot, x)$   
 22 is continuous. This allows to define the *empirical risk minimization functional*:

$$J_T(w, b) := \mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}) + \int_0^T |(w(t), b(t))|^2 dt, \quad (1.5)$$

for:

$$\mathcal{E}(\tilde{\mathbf{x}}, \mathbf{x}) := \frac{1}{N} \sum_{i=1}^N E(\tilde{x}^i, x^i),$$

23 and  $y$  a solution of (1.1). In this paper we denote any minimizer of  $J_T$  by  $(w_T, b_T)$ . Moreover, the  
 24 trajectories induced by such minimizers, called *optimal trajectories*, are denoted by  $y_T(t; \mathbf{x}) :=$   
 25  $y(t; \mathbf{x}, w_T, b_T)$ .

26 First of all, we recall that the functional  $J_T$  has at least a minimizer:

27 **Proposition 1.1** (Existence of minimizers). *The functional  $J_T$  given in (1.5) has at least a*  
 28 *minimizer in  $L^2(0, T; \mathcal{U})$ .*

1 Proposition 1.1 is classical, and the proof can be found, for instance, in [Tré05, Proposition  
2 6.2.3]. The main idea of the proof is that  $J_T$  is a sum of a positive weakly continuous functional  
3 and a positive continuous convex functional. For this result the continuity of  $E(\cdot, x)$  is essential.

4 *Example 1.2* (On the non-uniqueness of minimizers). We might have more than one minimizer.  
5 If  $d = 1$ ,  $\mathbf{x} = x_1 = 0$ ,  $E(x, 0) = \min\{|1-x|, |1+x|\}$  and  $\sigma(s) = s$ ,  $J_T$  has at least two minimizers  
6 for any  $T > 0$ . Indeed,  $J_T(0, 1/T) = 1/T$  for  $T \geq 1$  and  $J_T(0, 1/2) = 1 - T/4$  for  $T \in (0, 1]$ ,  
7 so  $(0, 0)$  is not a minimizer. Moreover, if  $(w, b)$  is a minimizer, by symmetry  $(-w, -b)$  is also a  
8 minimizer.

9 For having  $J_T$  minimizers which take the error to 0 the first thing that we need, of course, is  
10 that the error can be taken to 0, a property known in Data Science as *interpolation*:

*Hypothesis 1* (Interpolation). For the data set  $\mathbf{x}$  there are controls:

$$(w_*, b_*) \in L^2(0, 1; \mathcal{U}),$$

11 such that  $\mathcal{E}(y(1; \mathbf{x}, w_*, b_*), \mathbf{x}) = 0$ .

12 Hypothesis 1 is used to show that the error of the final state of the optimal trajectories decays  
13 as  $T \rightarrow \infty$  (see Lemma 2.3 below).

14 In addition, we assume that the error can be taken to 0 with a cost proportional to the error:

*Hypothesis 2* (Local controllability of the system). Let the data set be  $\mathbf{x} \in (\mathbb{R}^d)_*^N$ . Then, there  
are  $C, \tilde{\varepsilon} > 0$  such that for all  $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_N) \in (\mathbb{R}^d)_*^N$  satisfying  $\mathcal{E}(\bar{\mathbf{x}}, \mathbf{x}) < \tilde{\varepsilon}$ , there are some  
controls  $(w, b)$  satisfying:

$$\|(w, b)\|_{L^\infty(0, 1; \mathcal{U})} < C\mathcal{E}(\bar{\mathbf{x}}, \mathbf{x}),$$

such that:

$$\mathcal{E}(y(1; \bar{\mathbf{x}}, w, b), \mathbf{x}) = 0.$$

15 As shown in Appendix A, an additional assumption to Hypothesis 1 is necessary to take the  
16 error to 0. Without Hypothesis 2 the cost to take the error to 0 may be considerably higher  
17 than obtaining some small error, so it might not compensate to take the error exactly to 0 (see  
18 Proposition A.1).

19 Now we have all the tools to state the first main result of this paper:

20 **Theorem 1.3** (Annihilation of the error in a long time horizon). *Let  $\mathbf{x} \in (\mathbb{R}^d)_*^N$ ,  $\sigma$  be a  
21 Lipschitz activation function and  $E$  be an error function such that Hypotheses 1 and 2 are  
22 satisfied. Then, for  $T > 0$  large enough depending on  $\sigma$ ,  $\mathbf{x}$  and  $E$ , if  $\varepsilon > 0$  there is  $\delta > 0$  such  
23 that  $J_T(w, b) < \inf J_T + \delta$  implies:*

$$\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}) < \varepsilon. \tag{1.6}$$

1 Moreover, for  $T > 0$  large enough the following equality holds for any optimal trajectory:

$$\mathcal{E}(y_T(T; \mathbf{x}), \mathbf{x}) = 0. \quad (1.7)$$

Theorem 1.3 is proved by showing that if  $T$  is sufficiently large and if  $\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x})$  is small and strictly positive, we can construct with Hypothesis 2 a control  $(\tilde{w}, \tilde{b})$  such that:

$$J(\tilde{w}, \tilde{b}) \leq J(w, b) - \frac{1}{2}\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}).$$

2 The construction of such control is far from trivial and, as illustrated in Appendix A, the  
 3 hypotheses are rather sharp. As explained in the first part of the introduction, Theorem 1.3  
 4 improves the results presented in [EGPZ20a], where the authors prove that the error of the  
 5 final state of the optimal trajectory is of size  $\mathcal{O}(1/T)$ .

*Example 1.4* (Application of Theorem 1.3 to the classification problem). Let us consider:

$$\mathbf{x} = (x^1, \dots, x^M, x^{M+1}, \dots, x^N) \in (\mathbb{R}^d)_*^N,$$

6 the error function  $E(x, x^i) = (x_1 + 1)1_{x_1 > -1}(x_1)$  for  $i = 1, \dots, M$ , and  $E(x, x^i) = (x_1 -$   
 7  $1)1_{x_1 > 1}(x_1)$  for  $i = M+1, \dots, N$ ; and any neural function  $\sigma$  of the type  $\sigma(x) = (\tilde{\sigma}(x_1), \dots, \tilde{\sigma}(x_d))$   
 8 such that there is  $c > 0$  such that  $cs \leq \tilde{\sigma}(s)$  for all  $s \geq 0$  and  $\tilde{\sigma}(s) \leq cs$  for all  $s \leq 0$ .  
 9 Hypothesis 2 is clearly satisfied, as it suffices to consider  $\tilde{\varepsilon} = 1/(2N)$ ,  $b = 0$  and  $w(t)x =$   
 10  $(2Nc^{-1}\mathcal{E}(\bar{\mathbf{x}}, \mathbf{x})x_1, 0, \dots, 0)$ . Thus, Theorem 1.3 implies that if the data can be classified (i.e. if  
 11 Hypothesis 1 is satisfied), then by computing the optimal control for a sufficiently large time,  
 12 the data is sent to the sets  $\{x_1 \leq -1\}$  and  $\{x_1 \geq 1\}$ .

## 13 1.2 Optimal trajectories for neural ODE with an homogeneous ac- 14 tivation function

15 In this section we present the analogous results to those in Section 1.1 for the neural ODE  
 16 (1.2) with activation functions which satisfy (1.4). Let us reformulate Hypotheses 1 and 2 in  
 17 the context of (1.2):

*Hypothesis 3* (Interpolation). For the data set  $\mathbf{x}$  there are controls:

$$(r_*, w_*, b_*) \in L^2(0, 1; \tilde{\mathcal{U}}),$$

18 such that  $\mathcal{E}(y(1; \mathbf{x}, r_*, w_*, b_*), \mathbf{x}) = 0$ .

*Hypothesis 4* (Local controllability of the system). Let the data set be  $\mathbf{x} \in (\mathbb{R}^d)_*^N$ . Then, there  
 are  $C, \tilde{\varepsilon} > 0$  such that for all  $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_N)$  satisfying  $\mathcal{E}(\bar{\mathbf{x}}, \mathbf{x}) < \tilde{\varepsilon}$ , there are some controls  
 $(r, w, b)$  satisfying:

$$\|(w, b)\|_{L^\infty(0, 1; \tilde{\mathcal{U}})} < C\mathcal{E}(\bar{\mathbf{x}}, \mathbf{x}),$$

such that:

$$\mathcal{E}(y(1; \bar{\mathbf{x}}, r, w, b), \mathbf{x}) = 0.$$

Again, we seek to get sufficient conditions so that the optimal trajectories induced by:

$$\tilde{J}_T(r, w, b) := \mathcal{E}(y(T; \mathbf{x}, r, w, b), \mathbf{x}) + \int_0^T |(w(t), b(t))|^2 dt, \quad (1.8)$$

satisfy  $\mathcal{E}(y_T(T; \mathbf{x}), \mathbf{x}) = 0$ . Since  $|r|$  is constant (see (1.3)), it makes no sense to include it in the definition of  $\tilde{J}_T$ . For the functional  $\tilde{J}_T$  the following result holds:

**Theorem 1.5** (Annihilation of the error for a sufficiently large time). *Let  $\sigma$  be a Lipschitz activation function satisfying (1.4) and  $E$  an error function satisfying Hypothesis 3 and 4. Then, for  $T > 0$  large enough depending on  $\sigma$ ,  $\mathbf{x}$  and  $E$ , and all  $\varepsilon > 0$  there is  $\delta > 0$  such that if  $J_T(r, w, b) < \inf J_T + \delta$ , then:*

$$\mathcal{E}(y(T; \mathbf{x}, r, w, b), \mathbf{x}) < \varepsilon. \quad (1.9)$$

Moreover, if  $T$  is large enough and if  $\tilde{J}_T$  has an optimal trajectory:

$$\mathcal{E}(y_T(T; \mathbf{x}), \mathbf{x}) = 0. \quad (1.10)$$

The proof of Theorem 1.5 is analogous to that of Theorem 1.3, so we just give some brief explanations in the first comment of Section 3. As with Theorem 1.3, Theorem 1.5 improves the results presented in [EGPZ20a], where the authors prove that the error of the optimal trajectory at a final time  $T$  is of size  $\mathcal{O}(1/T)$  also for the solutions of (1.2) with an activation functions satisfying (1.4).

*Remark 1.6* (Existence of minimizers of  $\tilde{J}_T$ ). We have stated “if  $\tilde{J}_T$  has an optimal trajectory” in Theorem 1.5 because, as far as we know, it is an open question to see if  $\tilde{J}_T$  admits a minimizer. The main obstacle to adapt the proof of Proposition 1.1 is that nonlinear functions and weak limits may not commute. However, as we see in the first comment of Section 3, we can improve Theorem 1.5 and obtain that for  $T$  large enough and all  $\varepsilon > 0$  there are controls  $(r, w, b)$  such that  $J_T(r, w, b) < \inf J_T + \varepsilon$  and  $\mathcal{E}(y(t; \mathbf{x}, r, w, b), \mathbf{x}) = 0$ .

*Example 1.7* (Application of Theorem 1.5 to simultaneous controllability). Theorem 1.5 can be applied to the simultaneous controllability problem. Let  $\mathbf{x} \in (\mathbb{R}^d)_*^N$  for  $d \geq 2$ , the activation function:

$$\sigma(x) = (\max\{x_1, 0\}, \dots, \max\{x_d, 0\}), \quad (1.11)$$

the targets  $\mathbf{z} = (z^1, \dots, z^N) \in (\mathbb{R}^d)_*^N$ , and the error function  $E(x, x^i) = |x - z^i|$ . Then, it is proved in [RBZ21, Theorem 2] that Hypothesis 3 is satisfied. Moreover, as we prove in Appendix C, Hypothesis 4 also holds. We present the proof because the bounds for the cost of the control is not a straight consequence of the computations in [RBZ21]. Consequently, Theorem 1.5 (and all the auxiliary results and corollaries) can be applied to this neural problem.

### 1.3 Outline of the paper

The rest of the paper is organized as follows: in Section 2 we present the proof of Theorem 1.3, in Section 3 we comment some additional facts about neural ODE and present some open problems, in Appendix A we present a pathological case that motivate Hypotheses 2 and 4, in Appendix B we present the proof of a technical result involving measure theory, and in Appendix C we show that the simultaneous controllability problem satisfies Hypothesis 4.

## 2 Optimal control for affine neural ODE

In this section we work in the control problem described by (1.1) and the risk minimization functional  $J_T$  given by (1.5). In this section  $C > 0$  denotes an arbitrary constant that may change from line to line and depends only on  $\sigma$ ,  $E$  and  $\mathbf{x}$ . Similarly, when we assume that  $T$  is large enough we mean with respect to  $\sigma$ ,  $E$  and  $\mathbf{x}$ . We first present some technical results in Section 2.1, then conclude the proof of Theorem 1.3 in Section 2.2 by reductio ad absurdum, and finally provide additional properties of the optimal controls in Section 2.3.

### 2.1 Preliminaries

We first remark that the trajectories of (1.1) are invariant with the time variable:

**Lemma 2.1** (Invariance of trajectories with the time variable). *Let  $\mathbf{x} \in (\mathbb{R}^d)^N$  and*

$$\phi \in L^1_{loc}(0, \infty; \mathbb{R}^+).$$

*Then:*

$$y(\Phi(t); \mathbf{x}, w, b) = y(t; \mathbf{x}, \phi(\Phi(s))w(\Phi(s)), \phi(\Phi(s))b(\Phi(s))), \quad \forall t \in [0, T^*], \quad (2.1)$$

*for  $T^* > 0$  and  $\Phi$  any solution of:*

$$\begin{cases} \dot{\Phi}(s) = \phi(\Phi(s)), & s \in [0, T^*), \\ \Phi(0) = 0. \end{cases} \quad (2.2)$$

**Remark 2.2** (Invariance of trajectories when  $\phi$  is constant). An important application of Lemma 2.1 is the case  $\phi(t) = \lambda \in \mathbb{R}^+$ ; that is, when  $\phi$  is constant. Then, (2.1) becomes:

$$y(\lambda t; \mathbf{x}, w, b) = y(t; \mathbf{x}, \lambda w(\lambda s), \lambda b(\lambda s)). \quad (2.3)$$

*Proof of Lemma 2.1.* It suffices to see that for all  $i$  the function  $t \mapsto y(\Phi(t); x^i, w, b)$  is a solution of (1.1) with initial value  $x^i$  and controls  $\phi(\Phi(t))w(\Phi(t))$  and  $\phi(\Phi(t))b(\Phi(t))$ , since (1.1) has a unique solution by Cauchy-Lipschitz Theorem. From (2.2)<sub>2</sub> we obtain that:

$$y(\Phi(0); x^i, w, b) = y(0; x^i, w, b) = x^i.$$

1 Moreover, from (2.2)<sub>1</sub> and the chain rule:

$$\begin{aligned}
\frac{d}{dt} \left( y(\Phi(t); x^i, w, b) \right) &= \phi(\Phi(t)) \dot{y}(\Phi(t); x^i, w, b) \\
&= \phi(\Phi(t)) \left( w(\Phi(t)) \sigma(y(\Phi(t); x^i, w, b)) + b(\Phi(t)) \right) \\
&= (\phi(\Phi(t)) w(\Phi(t))) \sigma(y(\Phi(t); x^i, w, b)) + (\phi(\Phi(t)) b(\Phi(t))).
\end{aligned} \tag{2.4}$$

2 □

3 Next, we recall that Hypothesis 1 implies that the error is at most of size  $\mathcal{O}(1/T)$ :

4 **Lemma 2.3** (Boundedness of the error with respect to  $T$ ). *Let  $T > 0$ ,  $\sigma$  be an activation*  
5 *function and  $E$  an error function satisfying Hypothesis 1. Then:*

$$\mathcal{E}(y_T(T; \mathbf{x}), \mathbf{x}) \leq \frac{C}{T}. \tag{2.5}$$

6 Lemma 2.3 is proved in [EGPZ20a]. Briefly, it is a consequence of the definition of  $y_T$  as the  
7 optimal trajectory and that, by Remark 2.2,  $\frac{1}{T} w_*(\frac{\cdot}{T})$  and  $\frac{1}{T} b_*(\frac{\cdot}{T})$  are controls that take the  
8 error to 0 (see Hypothesis 1 for the definition of  $(w_*, b_*)$ ). In fact:

$$\begin{aligned}
\mathcal{E}(y_T(T; \mathbf{x}), \mathbf{x}) &\leq J_T(w_T, b_T) \leq J_T \left( \frac{1}{T} w_* \left( \frac{\cdot}{T} \right), \frac{1}{T} b_* \left( \frac{\cdot}{T} \right) \right) \\
&= \frac{1}{T} \int_0^1 |(w_*(t), b_*(t))|^2 dt.
\end{aligned} \tag{2.6}$$

9 Finally, we prove that we may assume that the norm of the optimal control is constant if it  
10 does not get null or explode:

**Lemma 2.4** (More efficient controls). *Let:*

$$(w, b) \in C^1([0, T]; \mathcal{U} \setminus \{(0, 0)\}),$$

11 *be such that  $t \mapsto |(w(t), b(t))|$  is not constant. Then, there is a control  $(\tilde{w}, \tilde{b})$  such that  $t \mapsto$*   
12  *$|(\tilde{w}(t), \tilde{b}(t))|$  is constant, such that:*

$$|(\tilde{w}(t), \tilde{b}(t))| \in \left( \min_{[0, T]} |(w, b)|, \max_{[0, T]} |(w, b)| \right) \text{ in } [0, T], \tag{2.7}$$

13 *and such that:*

$$y(T; \mathbf{x}, \tilde{w}, \tilde{b}) = y(T; \mathbf{x}, w, b), \quad \int_0^T |(\tilde{w}(t), \tilde{b}(t))|^2 dt < \int_0^T |(w(t), b(t))|^2 dt. \tag{2.8}$$

1 *Proof of Lemma 2.4.* Let us consider the auxiliary function:

$$\phi_\gamma(t) = \frac{\gamma}{|(w(t), b(t))|} 1_{[0, T]}(t), \quad (2.9)$$

2 for  $\gamma > 0$  to be fixed later, and  $\Phi_\gamma$  given by:

$$\begin{cases} \dot{\Phi}_\gamma(s) = \phi_\gamma(\Phi_\gamma(s)), & s \in [0, T_\gamma], \\ \Phi_\gamma(0) = 0, \end{cases} \quad (2.10)$$

3 for:

$$T_\gamma := \sup\{t : \Phi_\gamma(t) < T\}.$$

4 Note that from the definition of  $T_\gamma$  it follows that:

$$\Phi_\gamma(T_\gamma) = T. \quad (2.11)$$

5 Since  $\phi_\gamma$  is  $C^1$  (as  $(w, b) \neq 0$ , by compactness  $\min_{[0, T]} |(w, b)| > 0$ ), (2.10) has a unique solution  
6 by Cauchy-Lipschitz Theorem. Moreover,  $\gamma \mapsto T_\gamma$  is continuous and decreasing,  $\lim_{\gamma \rightarrow 0} T_\gamma = \infty$ ,  
7 and  $\lim_{\gamma \rightarrow \infty} T_\gamma = 0$ , so there is  $\gamma^* > 0$  such that:

$$T_{\gamma^*} = T. \quad (2.12)$$

8 For that value  $\gamma^*$  it holds:

$$\gamma^* = \frac{1}{\frac{1}{T} \int_0^T \frac{dt}{|(w(\Phi_{\gamma^*}(t)), b(\Phi_{\gamma^*}(t)))|}}, \quad (2.13)$$

9 because by (2.11) and (2.12):

$$\int_0^T \frac{\gamma^* dt}{|(w(\Phi_{\gamma^*}(t)), b(\Phi_{\gamma^*}(t)))|} = \int_0^T \dot{\Phi}_{\gamma^*}(t) dt = \Phi_{\gamma^*}(T) = \Phi_{\gamma^*}(T_{\gamma^*}) = T.$$

10 Considering (2.1) we obtain that:

$$\begin{aligned} y(T; \mathbf{x}, \phi_{\gamma^*}(\Phi_{\gamma^*}(s))w(\Phi_{\gamma^*}(s)), \phi_{\gamma^*}(\Phi_{\gamma^*}(s))b(\Phi_{\gamma^*}(s))) &= y(\Phi_{\gamma^*}(T); \mathbf{x}, w, b) \\ &= y(T; \mathbf{x}, w, b). \end{aligned} \quad (2.14)$$

11 Moreover, considering the strict inequality between the harmonic and arithmetic means (see,  
12 for instance, [Kom11]), (2.10)-(2.13), and the Change of Variables Theorem we get that:

$$\begin{aligned} \int_0^T \phi_{\gamma^*}^2(\Phi_{\gamma^*}(t)) |(w(\Phi_{\gamma^*}(t)), b(\Phi_{\gamma^*}(t)))|^2 dt &= \int_0^T (\gamma^*)^2 dt = (\gamma^*)^2 T \\ &= \frac{\gamma^* T}{\frac{1}{T} \int_0^T \frac{dt}{|(w(\Phi_{\gamma^*}(t)), b(\Phi_{\gamma^*}(t)))|}} \\ &< \frac{1}{T} \int_0^T \gamma^* T |(w(\Phi_{\gamma^*}(t)), b(\Phi_{\gamma^*}(t)))| dt \\ &= \int_0^T \phi_{\gamma^*}(\Phi_{\gamma^*}(t)) |(w(\Phi_{\gamma^*}(t)), b(\Phi_{\gamma^*}(t)))|^2 dt \\ &= \int_0^T |(w(t), b(t))|^2 dt. \end{aligned} \quad (2.15)$$

Therefore, combining (2.14) and (2.15) we obtain (2.8) for:

$$\tilde{w}(t) = \phi_{\gamma^*}(\Phi_{\gamma^*}(t))w(\Phi_{\gamma^*}(t)), \quad \tilde{b}(t) = \phi_{\gamma^*}(\Phi_{\gamma^*}(t))b(\Phi_{\gamma^*}(t)).$$

Finally, since  $\gamma^*$  is the harmonic mean of values in:

$$\left[ \min_{[0,T]} |(w, b)|, \max_{[0,T]} |(w, b)| \right],$$

1 we obtain (2.7). □

## 2 2.2 Construction of controls which take the error to zero

3 Let us state the properties of the controls that we construct in this section:

4 **Proposition 2.5.** *Let  $\sigma$  be an activation function and  $E$  an error function that satisfy Hy-*  
 5 *potheses 1 and 2. Let  $T > 0$  be large enough and  $(w, b)$  be such that:*

$$J_T(w, b) \leq 2 \inf J_T. \tag{2.16}$$

6 *Then, there is a control  $(\hat{w}, \hat{b})$  such that:*

$$\mathcal{E}(y(T; \mathbf{x}, \hat{w}, \hat{b}), \mathbf{x}) = 0, \tag{2.17}$$

7 *and:*

$$J_T(\hat{w}, \hat{b}) \leq J_T(w, b) - \frac{1}{2} \mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}). \tag{2.18}$$

8 The first step is to remark that  $\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x})$  is small for large  $T$ . The second step is to  
 9 approximate  $(w, b)$  by some control  $(\tilde{w}, \tilde{b})$  satisfying the hypothesis of Lemma 2.4. The third  
 10 step is to show that if (1.7) is false, we may prolong for some  $\tau > 0$  the controls  $\tilde{w}$  and  $\tilde{b}$  in  
 11  $[T, T + \tau]$  so that:

$$\tilde{y}(T + \tau; \mathbf{x}, \tilde{w}, \tilde{b}) = 0. \tag{2.19}$$

12 The fourth step is to take those trajectories to  $[0, T]$  with (2.3). The fifth and last step is to  
 13 check that the new control satisfies (2.18).

14 *Proof of Proposition 2.5. Step 1: estimate of  $\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x})$ .* If  $\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}) = 0$ ,  
 15 then it suffices to consider  $(\hat{w}, \hat{b}) = 0$ , so we suppose from now on that  $\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}) > 0$ .  
 16 Moreover, from (2.6) and (2.16) we obtain for  $T$  large enough that:

$$\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}) \in (0, \tilde{\varepsilon}/2), \tag{2.20}$$

17 for  $\tilde{\varepsilon}$  the value in Hypothesis 2.

**Step 2: approximating the control.** Clearly,  $C^1([0, T]; \mathcal{U} \setminus \{0\})$  is dense in  $L^2(0, T; \mathcal{U})$ . Moreover,

$$(w, b) \mapsto \mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}),$$

1 is continuous from  $L^2(0, T; \mathcal{U})$  to  $\mathbb{R}$ . Thus, there is  $(\tilde{w}, \tilde{b}) \in C^1([0, T]; \mathcal{U} \setminus \{0\})$  such that:

$$\|(\tilde{w}, \tilde{b})\|_{L^2(0, T; \mathcal{U})} \leq \|(w, b)\|_{L^2(0, T; \mathcal{U})}, \quad (2.21)$$

2 and:

$$\mathcal{E}(y(T; \mathbf{x}, \tilde{w}, \tilde{b}), \mathbf{x}) \leq 2\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}). \quad (2.22)$$

3 Moreover, by Lemma 2.4 we can suppose that  $t \mapsto |(\tilde{w}(t), \tilde{b}(t))|$  is constant. In addition, from  
4 (2.6), (2.16) and (2.21) it follows that:

$$\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})} \leq \frac{C}{\sqrt{T}}. \quad (2.23)$$

**Step 3: taking the error to 0.** From Hypothesis 2, (2.20) and (2.22) we obtain a control  $(\bar{w}, \bar{b}) \in L^\infty(0, T; \mathcal{U})$  that takes the solution from  $y(T; \mathbf{x}, \tilde{w}, \tilde{b})$  to a state  $\tilde{\mathbf{x}}$  such that  $\mathcal{E}(\tilde{\mathbf{x}}, \mathbf{x}) = 0$ . Moreover,

$$\|(\bar{w}, \bar{b})\|_{L^\infty(0, 1; \mathcal{U})} \leq C\mathcal{E}(y(T; \mathbf{x}, \tilde{w}, \tilde{b}), \mathbf{x}) \leq C\mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}).$$

5 Consequently, by Remark 2.2, for some:

$$\tau \leq C \frac{\mathcal{E}(y(T; \mathbf{x}, \tilde{w}, \tilde{b}), \mathbf{x})}{\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})}}, \quad (2.24)$$

6 the control  $(\tilde{w}, \tilde{b})$  can be prolonged to  $[0, T + \tau]$  so that both:

$$\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T + \tau; \mathcal{U})} = \|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})}, \quad (2.25)$$

7 and (2.19) are satisfied.

8 **Step 4: taking the trajectory to  $[0, T]$ .** We consider:

$$\begin{aligned} \hat{w}(t) &:= \frac{T + \tau}{T} \tilde{w} \left( \frac{T + \tau}{T} t \right), \\ \hat{b}(t) &:= \frac{T + \tau}{T} \tilde{b} \left( \frac{T + \tau}{T} t \right). \end{aligned} \quad (2.26)$$

9 Then, (2.17) is true. In fact, the equation (2.3) with  $\lambda = \frac{T + \tau}{T}$  implies:

$$y(T; \mathbf{x}, \hat{w}, \hat{b}) = y(T + \tau; \mathbf{x}, \tilde{w}, \tilde{b}).$$

**Step 5: efficiency of the new control.** First, we realize that:

$$J_T(w, b) - J_T(\hat{w}, \hat{b}) = \mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}) + \int_0^T |(w(t), b(t))|^2 dt - \left(\frac{T+\tau}{T}\right)^2 \int_0^T \left| \left( \tilde{w}\left(\frac{T+\tau}{T}t\right), \tilde{b}\left(\frac{T+\tau}{T}t\right) \right) \right|^2 dt. \quad (2.27)$$

- 1 Considering that  $t \mapsto (\tilde{w}(t), \tilde{b}(t))$  is constant in  $[0, T]$ , and that (2.21) and (2.25) are satisfied  
 2 we deduce that:

$$\int_0^T |(w(t), b(t))|^2 dt - \int_0^T \left| \left( \tilde{w}\left(\frac{T+\tau}{T}t\right), \tilde{b}\left(\frac{T+\tau}{T}t\right) \right) \right|^2 dt \geq 0. \quad (2.28)$$

Consequently, we obtain from (2.20), (2.23)-(2.25) and (2.27)-(2.28) that:

$$\begin{aligned} J_T(w, b) - J_T(\hat{w}, \hat{b}) &\geq \mathcal{E}(y_T(T; \mathbf{x}, w, b), \mathbf{x}) - \left(\frac{2\tau}{T} + \frac{\tau^2}{T^2}\right) \int_0^T \left| \left( \tilde{w}\left(\frac{T+\tau}{T}t\right), \tilde{b}\left(\frac{T+\tau}{T}t\right) \right) \right|^2 dt \\ &\geq \left(1 - C\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})} - C\varepsilon T^{-1}\right) \mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}) \\ &\geq (1 - CT^{-1/2}) \mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}), \end{aligned}$$

- 3 which implies (2.18) for  $T$  large enough. □

- 4 Now we may conclude the proof of Theorem 1.3 by reductio ad absurdum:

*Conclusion of the proof of Theorem 1.3.* Let  $\varepsilon > 0$ . It suffices to consider  $\delta = \varepsilon/3$ . If  $(w, b)$  are such that  $J_T(w, b) \leq \inf J_T + \varepsilon/3$ , then  $\mathcal{E}(y(T; \mathbf{x}, w, b)) < \varepsilon$ . Otherwise, by Proposition 2.5 there are  $(\hat{w}, \hat{b})$  such that:

$$J_T(\hat{w}, \hat{b}) \leq J_T(w, b) - \frac{\varepsilon}{2} \leq \inf J_T - \frac{\varepsilon}{6},$$

- 5 which is absurd. Similarly, if  $(w_T, b_T)$  is a minimizer of  $J_T$  and (1.7) is not satisfied, then the  
 6 control  $(\hat{w}, \hat{b})$  of Proposition 2.5 satisfies  $J_T(\hat{w}, \hat{b}) < J_T(w_T, b_T)$ , contradicting the definition of  
 7 minimizer. □

### 8 **2.3 Additional properties of the optimal control**

As a consequence of Remark 2.2, we can easily prove that, assuming Hypothesis 1 and 2, for a sufficiently large time the optimal controls are of the form:

$$\left( \frac{1}{T} w_* \left( \frac{t}{T} \right), \frac{1}{T} b_* \left( \frac{t}{T} \right) \right),$$

1 for  $(w_*, b_*)$  the minimizers of the functional:

$$t \mapsto \int_0^T |(w(t), b(t))|^2 dt,$$

2 considered in the domain:

$$\{(w, b) : \mathcal{E}(y(1; \mathbf{x}, w, b), \mathbf{x}) = 0\}.$$

3 In addition, we can prove that such minimizers belong to  $L^\infty(0, T)$  and satisfy that  $t \mapsto$   
 4  $|(w(t), b(t))|$  is constant, which follows from:

**Lemma 2.6** (A more efficient control). *Let  $(w, b)$  a control in  $L^2(0, 1)$  such that  $t \mapsto |(w(t), b(t))|$  is not constant. Then, there is a control  $(\tilde{w}, \tilde{b})$  such that:*

$$y(1; \mathbf{x}, \tilde{w}, \tilde{b}) = y(1; \mathbf{x}, w, b),$$

$$\|(\tilde{w}, \tilde{b})\|_{L^2(0, T; \mathcal{U})} < \|(w, b)\|_{L^2(0, T; \mathcal{U})},$$

6 and, if  $(w, b) \in L^\infty(0, T; \mathcal{U})$ ,

$$\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})} \leq 2\|(w, b)\|_{L^\infty(0, T; \mathcal{U})}. \quad (2.29)$$

7 The proof of Lemma 2.6 is based on classical results from Measure Theory and is postponed to  
 8 Appendix B. Lemma 2.6, compared to Lemma 2.4, has the advantage of having less restrictive  
 9 hypothesis. However, it has the disadvantage that we do not obtain neither a contraction for  
 10 the  $L^\infty$  norm (see Remark B.2) nor a control with constant norm, which is needed for proving  
 11 Proposition 2.5.

### 12 3 Further comments and open problems

- **Analogous results for neural ODE whose dynamics are described by (1.2).** Clearly Lemmas 2.1, 2.3, 2.4, and 2.6 and Proposition 2.5 can be proved for system (1.2) with  $\sigma$  satisfying (1.4) as in Section 2. The key lemma is Lemma 2.1, since the other results use the homogeneity of the system via Lemma 2.1. The analogous of Lemma 2.1 can be proved by replacing (2.4) by:

$$\begin{aligned} \frac{d}{dt}(y(\Phi(t); x^i, w, b, r)) &= \phi(\Phi(t))r(\Phi(t))\dot{y}(\Phi(t); x^i, w, b) \\ &= \phi(\Phi(t))r(\Phi(t))\sigma\left(w(\Phi(t))y(\Phi(t); x^i, w, b) + b(\Phi(t))\right) \\ &= r(\Phi(t))\sigma\left(\phi(\Phi(t))w(\Phi(t))y(\Phi(t); x^i, w, b) + \phi(\Phi(t))b(\Phi(t))\right). \end{aligned}$$

13 The last equality follows from (1.4). Finally, Theorem 1.5 and the analogous of Proposi-  
 14 tion 2.5 imply that for all  $\delta > 0$  there is a control  $(r, w, b)$  such that  $\tilde{J}_T(r, w, b) < \inf \tilde{J}_T - \delta$   
 15 and  $\mathcal{E}(y(T; \mathbf{x}, r, w, b)) = 0$ .

1 • **Direction of the flow in neural ODE described by (1.2).** The same results are valid  
 2 if we replace  $\mathcal{M}$  (see (1.3)) by any closed subset of orthogonal matrices. In particular,  
 3 this is true if  $\mathcal{M} = \{I\}$ ; that is, if  $r(t) = I$ , which is the case studied in [EGPZ20a].

4 • **Functionals allowing expensive controls.** As in [EGPZ20a], we can consider the  
 5 functional:

$$J_{T,\delta}(w, b) := \mathcal{E}(y(T; \mathbf{x}, w, b), \mathbf{x}) + \delta \int_0^T |(w(t), b(t))|^2 dt,$$

6 instead of  $J_T$  for (1.1), and:

$$J_{T,\delta}(r, w, b) := \mathcal{E}(y(T; \mathbf{x}, r, w, b), \mathbf{x}) + \delta \int_0^T |(w(t), b(t))|^2 dt,$$

7 instead of  $J_T$  for (1.2)-(1.4). By linearity (see Remark 2.2) it holds that:

$$J_{T,\delta}(w, b) = J_{T\delta^{-1},1}(\delta w(t\delta), \delta b(t\delta)),$$

8 and:

$$\tilde{J}_{T,\delta}(r, w, b) = \tilde{J}_{T\delta^{-1},1}(r(t\delta), \delta w(t\delta), \delta b(t\delta)),$$

9 respectively. A straight consequence is that  $(w, b)$  is a minimizer of  $J_{T,\delta}$  if and only  
 10 if  $(\delta w(t\delta), \delta b(t\delta))$  is a minimizer of  $J_{T\delta^{-1},1}$ . Similarly,  $(r, w, b)$  is a minimizer of  $\tilde{J}_{T,\delta}$  if  
 11 and only if  $(r(t\delta), \delta w(t\delta), \delta b(t\delta))$  is a minimizer of  $J_{T\delta^{-1},1}$ . Thus, analogous results to  
 12 Theorems 1.3 and 1.5 and all the auxiliary results hold true for  $J_{T,\delta}$  and  $\tilde{J}_{T,\delta}$  when  $T$  is  
 13 fixed and  $\delta > 0$  is small enough depending on  $\sigma$ ,  $E$ ,  $\mathbf{x}$  and  $T$ .

• **Optimal control for non-homogenous activation functions.** It remains an open  
 problem to determine if similar results to Theorem 1.5 hold for non-homogeneous activa-  
 tion functions satisfying  $\sigma(0) = 0$  such as the hyperbolic tangent,

$$\sigma(x) = (\tanh(x_1), \dots, \tanh(x_d)),$$

see [FS18]. We may wonder whether similar results hold with more general activation  
 functions if we replace  $\mathcal{M}$  (see (1.3)) by the unitary matrices or by  $\mathcal{L}(\mathbb{R}^d, \mathbb{R}^d)$  (of course,  
 the cost of  $r$  must also be included in the risk minimization functional). This would  
 include, for instance, sigmoid,

$$\sigma(x) = ((1 + e^{-x_1})^{-1}, \dots, (1 + e^{-x_d})^{-1}),$$

see [MS95]; and softplus,

$$\sigma(x) = (\log(1 + e^{x_1}), \dots, \log(1 + e^{x_d})),$$

14 see [GBB11]. The main difficulty is that the analogue of Lemma 2.1 cease to be true, so  
 15 another tool is needed to prove the main result, probably a local inverse theorem result.

- **Optimal control with other norms.** It is a relevant problem to determine if similar results to Theorems 1.3 and 1.5 hold for any other Lebesgue or Sobolev norms. In particular, the most interesting scenario is to replace both in  $J_T$  and  $\tilde{J}_T$  the terms  $\|(w, b)\|_{L^2(0, T)}^2$  by  $\|(w, b)\|_{H^1(0, T)}^2$  and adding the restriction that the component of  $r$  can only change signs if  $(w, b) = 0$  or to measure the  $H^1$  norm of  $r$  if the space  $\mathcal{M}$  is connected. The interest of this is double: thinking in potential applications it makes sense to also try to bound the variations in the time variable, which can be obtained by minimizing the time derivative. Moreover, if we consider the  $H^1$ -norm we can prove as in Proposition 1.1 that  $\tilde{J}_T$  admits a minimizer. The main difficulties when studying these norms are that Lemmas 2.4 and 2.6 may not be proved as easily (if they are true) because we need to keep track of the time derivative and because we cannot define the control on  $[T, T + \tau]$  independently to the controls on  $[0, T]$  due to the time derivative.

## A A pathological case

In this section we prove that without Hypothesis 2 the error may not be taken exactly to 0 if the ratio between the cost of correcting the error and the error explodes as the error vanishes. We present an example for the sake of simplicity, though the proof can be replicated whenever the gradient of the error is null on all the points where the error is null, which is the key impediment for taking the error exactly to zero.

**Proposition A.1** (Necessity of local controllability). *Let us consider  $d = 1$ ,  $\mathbf{x} = x_1 = 1$ ,  $E(x, 1) = x^2$ ,  $\sigma(s) = s$  and  $J_T$  given by (1.5). Then,  $y_T(T) > 0$  for all  $T > 0$ .*

*Proof of Proposition A.1.* Let  $(w_T, b_T)$  be a minimizer of  $J_T$ . Clearly  $w_T, b_T \leq 0$ . Let us prove by contradiction that  $y_T(T) > 0$ . For that, we suppose that  $y_T(T) = 0$ . By Lemma 2.6,  $t \mapsto |(w_T(t), b_T(t))|$  is a constant function equal to some constant  $\mathfrak{c}$ . In particular, for  $\delta > 0$  small enough the following inequality is satisfied:

$$(y_T(T - \delta))^2 - \int_{T-\delta}^T |(w_T(t), b_T(t))|^2 = (y_T(T - \delta))^2 - \mathfrak{c}^2 \delta \leq (C\delta)^2 - \mathfrak{c}^2 \delta < 0. \quad (\text{A.1})$$

The estimate  $|y_T(T - \delta)| \leq C\delta$  follows from the formula:

$$y_T(T - \delta) = - \int_{T-\delta}^T b_T(s) \exp\left(- \int_{T-\delta}^s w_T(z) dz\right) ds,$$

which follows from  $y_T(T) = 0$ . Consequently, we obtain from (A.1) that:

$$J_T(w_T 1_{(0, T-\delta)}, b_T 1_{(0, T-\delta)}) - J_T(w_T, b_T) < 0, \quad (\text{A.2})$$

which contradicts that  $(w_T, b_T)$  is a minimizer of  $J_T$ .  $\square$

*Remark A.2* (On Hypothesis 1). It is trivial that Hypothesis 1 is satisfied by the activation and error function introduced in Propositions A.1.

## 1 B Proof of Lemma 2.6

2 In this section we prove Lemma 2.6. Here  $\mu$  denotes the Lebesgue measure. In order to prove  
 3 Lemma 2.6 we need the following classical result of measure theory, whose proof can be found  
 4 in [Yeh06, Thm. 3.25]:

5 **Lemma B.1** (Comparison between sets of positive measure and open sets). *Let  $S \subset [0, T]$  be  
 6 a measurable set such that  $\mu(S) > 0$ . Then, for all  $\varepsilon > 0$  there is an open set  $\mathcal{O}^\varepsilon = \bigcup_{i=1}^{n^\varepsilon} (a_i^\varepsilon, b_i^\varepsilon)$   
 7 such that  $\mu(\mathcal{O}^\varepsilon \Delta S) < \varepsilon$ .*

8 *Proof of Lemma 2.6.* Since  $|(w, b)|$  is not constant, there are some sets  $S_1$  and  $S_2$  and some  
 9 constants  $C_1, C_2 > 0$  such that  $C_1 < C_2$ ,  $|(w, b)| < C_1$  on  $S_1$ ,  $|(w, b)| > C_2$  on  $S_2$  and:

$$\inf\{|x^2 - x^1| : x^1 \in S_1, x^2 \in S_2\} > 0. \quad (\text{B.1})$$

10 From Lemma B.1 we get that for  $\varepsilon > 0$  small enough there are two sets  $\mathcal{O}_1^\varepsilon = \bigcup_{i=1}^{n_1^\varepsilon} (a_{1,i}^\varepsilon, b_{1,i}^\varepsilon)$   
 11 and  $\mathcal{O}_2^\varepsilon = \bigcup_{i=1}^{n_2^\varepsilon} (a_{2,i}^\varepsilon, b_{2,i}^\varepsilon)$  satisfying:

$$\mu(\mathcal{O}_1^\varepsilon \setminus S_1) < \varepsilon, \quad \mu(\mathcal{O}_2^\varepsilon \setminus S_2) < \varepsilon, \quad (\text{B.2})$$

12 and:

$$\mu(\mathcal{O}_1^\varepsilon) = \mu(\mathcal{O}_2^\varepsilon) = \frac{\min\{\mu(S_1), \mu(S_2)\}}{2}. \quad (\text{B.3})$$

13 If  $\varepsilon$  is small enough, because of (B.1) we may also assume that:

$$\mathcal{O}_1^\varepsilon \cap \mathcal{O}_2^\varepsilon = \emptyset. \quad (\text{B.4})$$

14 Let us consider the auxiliary function:

$$\phi_\gamma(t) = \begin{cases} 1 & t \in [0, T] \setminus (\mathcal{O}_1^\varepsilon \cup \mathcal{O}_2^\varepsilon), \\ 1 + \gamma & t \in \mathcal{O}_1^\varepsilon, \\ \frac{1+\gamma}{1+2\gamma} & t \in \mathcal{O}_2^\varepsilon, \\ 0 & t \geq T, \end{cases} \quad (\text{B.5})$$

15 for  $\gamma > 0$  to be fixed later, and  $\Phi_\gamma$  given by:

$$\begin{cases} \dot{\Phi}_\gamma(s) = \phi_\gamma(\Phi_\gamma(s)), & \forall s \geq 0, \\ \Phi_\gamma(0) = 0. \end{cases} \quad (\text{B.6})$$

16 We remark that:

$$\Phi_\gamma(T) = T. \quad (\text{B.7})$$

Indeed, it can be proved that if  $\Phi_\gamma(T_*) = a$  and  $\phi_\gamma(t) = c$  on  $[a, b]$ , then  $\Phi_\gamma(T_* + \frac{b-a}{c}) = b$ .  
 Hence:

$$\Phi_\gamma \left( \mu([0, T] \setminus (\mathcal{O}_1^\varepsilon \cup \mathcal{O}_2^\varepsilon)) + \frac{1}{1+\gamma} \mu(\mathcal{O}_1^\varepsilon) + \frac{1+2\gamma}{1+\gamma} \mu(\mathcal{O}_2^\varepsilon) \right) = T,$$

17 which considering (B.3), (B.4) and (B.5) implies (B.7).

1 Consequently, the following controls satisfy the conclusions of Lemma 2.6:

$$( \tilde{w}, \tilde{b} ) = \phi_\gamma(\Phi_\gamma(t))(w(\Phi_\gamma(t)), b(\Phi_\gamma(t))). \quad (\text{B.8})$$

Indeed, from (2.1) and (B.1) it holds that:

$$\begin{aligned} y(T; \mathbf{x}, \phi_\gamma(s)w(\Phi_\gamma(s)), \phi_\gamma(s)b(\Phi_\gamma(s))) &= y(\Phi_\gamma(T); \mathbf{x}, w, b) \\ &= y(T; \mathbf{x}, w, b). \end{aligned}$$

In addition, if  $\gamma$  and  $\varepsilon$  are small enough:

$$\begin{aligned} &\int_0^T |(w(t), b(t))|^2 dt - \int_0^T |\phi_\gamma(\Phi_\gamma(t))(w(\Phi_\gamma(t)), b(\Phi_\gamma(t)))|^2 dt \\ &= \int_{\mathcal{O}_1^\varepsilon \cup \mathcal{O}_2^\varepsilon} |(w(t), b(t))|^2 dt \\ &\quad - \int_{\Phi_\gamma^{-1}(\mathcal{O}_1^\varepsilon) \cup \Phi_\gamma^{-1}(\mathcal{O}_2^\varepsilon)} \phi_\gamma^2(\Phi_\gamma(t)) |(w(\Phi_\gamma(t)), b(\Phi_\gamma(t)))|^2 dt \\ &= -\gamma \int_{\mathcal{O}_1^\varepsilon} |(w(t), b(t))|^2 dt + \frac{\gamma}{1+2\gamma} \int_{\mathcal{O}_2^\varepsilon} |(w(t), b(t))|^2 dt \\ &\geq \frac{\gamma}{1+2\gamma} C_2 \left( \frac{\min\{\mu(S_1), \mu(S_2)\}}{2} - \varepsilon \right) \\ &\quad - \gamma C_1 \frac{\min\{\mu(S_1), \mu(S_2)\}}{2} - \|(w, b)\|_{L^2(\mathcal{O}_1^\varepsilon \setminus S_1)}^2 > 0. \end{aligned}$$

2 The second equality follows from the change of variable  $s = \Phi_\gamma(t)$ , the first inequality from the  
 3 definitions of  $S_1, S_2, \mathcal{O}_1^\varepsilon, \mathcal{O}_2^\varepsilon$  and (B.3), and the last inequality from  $C_2 > C_1$ , (B.2), being  $\gamma$   
 4 and  $\varepsilon$  small enough, and the well known identity:

$$\lim_{c \rightarrow 0} \sup_{\mu(A)=c} \|g\|_{L^2(A, dx)} = 0, \quad \forall g \in L^2(0, T).$$

5 Finally, if  $(w, b) \in L^\infty(0, T; \mathcal{U})$  the estimate (2.29) follows from (B.5) and (B.8) by taking  
 6  $\gamma \leq 1$ . □

7 *Remark B.2* (Sharpness of the estimate (2.29)). The construction provided in the previous  
 8 proof may not ensure us that:

$$\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})} \leq \|(w, b)\|_{L^\infty(0, T; \mathcal{U})};$$

9 for instance if  $|(w, b)| = 1_\Omega$ , for  $\Omega \subset [0, T]$  a set such that  $\mu(\Omega) \in (0, T)$  and which contains an  
 10 open neighbourhood of every rational number in  $[0, T]$ . However, we can replace in the estimate  
 11 (2.29) the constant 2 by any constant strictly greater than 1.

# C Local simultaneous controllability

In this section we prove the following result:

**Lemma C.1** (Local simultaneous controllability result). *Let  $\sigma$  be the activation function defined by (1.11) and  $E$  defined in Example 1.7. Then  $\sigma$  and  $E$  satisfy Hypothesis 4.*

The main contribution with respect to the result on [RBZ21] is that we keep track of the cost and continuity of the control. The controls that we construct are different to those in [RBZ21], in which  $w$  and  $b$  have a single non-zero component at any time, since we do not search for a sparse property, but to obtain the continuity of the controls with respect to the initial data. We recall that  $C$  is a positive constant sufficiently large changing from line to line which depends on the target set  $\mathbf{z}$ . Moreover, for  $d = 2$  we denote:

$$r = \begin{bmatrix} r_1 & 0 \\ 0 & r_2 \end{bmatrix}, \quad w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}, \quad b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}.$$

*Proof of Lemma C.1.* In order to simplify the notation we prove Lemma C.1 for the case  $d = 2$ , though the proof is analogous for any  $d \geq 2$ . We prove Lemma C.1 by induction on  $N$ .

**Step 1: the base case.** Let us begin with the case  $N = 1$ . We may take  $x = (x_1, x_2)$  to  $z = (z_1, z_2)$  with a force proportional to  $|z - x|$  by applying the controls  $w = 0$ ,  $b_1 = |z_1 - x_1|$ ,  $b_2 = |z_2 - x_2|$ ,  $r_1 = \text{sign}(z_1 - x_1)$  and  $r_2 = \text{sign}(z_2 - x_2)$ .

**Step 2: the inductive case. Step 2.1: rearranging the points.** We may suppose by rearranging the indexes that  $|z^N| = \max_{i=1, \dots, N} |z^i|$ . For the rest of the proof we define  $e := \frac{z^N}{|z^N|}$  and:

$$\delta := \min\{|z^N| - \max_{i=1, \dots, N-1} |z^i|, 1\}. \quad (\text{C.1})$$

Then,  $\delta > 0$  since, for  $i = 1, \dots, N - 1$ , either  $|z^i| < |z^N|$  or  $|z^i| = |z^N|$  but  $z^i \neq z^N$ , so  $\cos(z^i, e) < 1$ .

**Step 2.2: controlling  $(x^1, \dots, x^{N-1})$  in  $[0, 1/2]$ .** By the induction hypothesis and linearity we know that for  $\hat{\varepsilon}$  small enough, if  $\sum_{i=1}^{N-1} |z^i - x^i| < \hat{\varepsilon}$  there are some controls  $(r, w, b)$  satisfying:

$$\|(w, b)\|_{L^\infty(0, 1/2)} < C \sum_{i=1}^{N-1} |z^i - x^i|, \quad (\text{C.2})$$

and such that:

$$y(1/2; (x^1, \dots, x^{N-1}), r, w, b) = (z^1, \dots, z^{N-1}). \quad (\text{C.3})$$

We fix:

$$\tilde{\varepsilon} = \min \left\{ \hat{\varepsilon}, \frac{\delta}{C(|z^N| + 1)} \right\}.$$

1 If:

$$\sum_{i=1}^N |z^i - x^i| < \tilde{\varepsilon}, \quad (\text{C.4})$$

2 then:

$$|y(t; x^N, r, w, b) - z^N| < \frac{\delta}{2} \quad \forall t \in \left[0, \frac{1}{2}\right]. \quad (\text{C.5})$$

3 Indeed,  $|z^N - x^N| < \frac{\delta}{4}$  by (C.4) and, if  $|y(\cdot; x^N, w, b)| < |z^N| + \frac{\delta}{2}$  on  $[0, t]$ , for  $t \leq 1/2$ , then:

$$\int_0^t \sigma(w(s)y(s; x^N, r, w, b) + b(s)) ds \leq \frac{\|w\|_{L^\infty(0,t)} (|z^N| + \frac{\delta}{2}) + \|b\|_{L^\infty(0,t)}}{2} < \frac{\delta}{4},$$

4 considering (C.2) and that  $\tilde{\varepsilon} < \frac{\delta}{C(|z^N|+1)}$ . In a similar way, we can prove that:

$$|y(1/2; x^N, r, w, b) - z^N| \leq C \sum_{i=1}^N |z^i - x^i|. \quad (\text{C.6})$$

5 Indeed, for  $t \in [0, 1/2]$  by (C.2):

$$|\sigma(wy(1/2; x^N, r, w, b) + b)| \leq \|w\|_{L^\infty(0,1/2)} (|z^N| + \delta) + \|b\|_{L^\infty(0,1/2)} \leq C \sum_{i=1}^{N-1} |z^i - x^i|.$$

6 **Step 2.3: controlling  $y(1/2; x^N, r, w, b)_1$  in  $[1/2, 3/4]$ .** We seek to obtain that:

$$y(3/4; x^N, r, w, b)_1 = z_1^N. \quad (\text{C.7})$$

7 If  $y(1/2; x^N, w, b)_1 = z_1^N$ , it suffices to consider  $r_1 = 1$ ,  $r_2 = 1$ ,  $w = 0$  and  $b = 0$ , so we  
8 may restrict to the case  $y(1/2; x^N, w, b) \neq z_1^N$ . To obtain (C.7) we consider the controls  $r_1 =$   
9  $\text{sign}(z_1 - x_1)$ ,  $r_2 = 1$ ,  $w_1 = \mathbf{c} \sum_{i=1}^N |z^i - x^i| e$ ,  $b_1 = \mathbf{c} \sum_{i=1}^N |z^i - x^i| (-|z^N| + \delta)$ ,  $w_2 = 0$ ,  $b_2 = 0$   
10 in  $[1/2, 3/4]$ , for  $\mathbf{c}$  to be fixed later. These controls are constant in  $[1/2, 3/4]$ . First, we remark  
11 that  $\sigma(w_1 \cdot x + b) = 0$  for all  $x$  such that  $x \cdot e \leq |z^N| - \delta$ . In particular, from (C.1) and (C.3)  
12 we derive:

$$y(3/4; (x^1, \dots, x^{N-1}), r, w, b) = y(1/2; (x^1, \dots, x^{N-1}), r, w, b) = (z^1, \dots, z^{N-1}). \quad (\text{C.8})$$

Moreover,  $|y(t; x^N, r, w, b)_1 - z_1^N|$  is decreasing on  $[1/2, T_*]$ , for:

$$T_* := \inf \{ T_* \geq 1/2 : y(T_*; x^N, r, w, b)_1 = z_1^N \}.$$

In addition, thanks to (C.5) in  $[1/2, T_*]$  the following inequality is satisfied:

$$|\dot{y}(t; x^N, r, w, b)_1| = |w_1 \cdot y(t; x^N, r, w, b) - b_1| \geq \mathbf{c} \sum_{i=1}^N |z^i - x^i| \frac{\delta}{2}.$$

Combining this with (C.6) we obtain that  $T_* < 3/4$  if  $\mathbf{c} \geq C$ ; i.e. if  $\mathbf{c}$  is sufficiently large just with respect to  $\mathbf{z}$  (recall that  $\delta$  is a fixed parameter depending only on  $\mathbf{z}$ ). Consequently, since  $T_*$  is continuous with  $\mathbf{c}$  and  $\lim_{\mathbf{c} \rightarrow 0} T_* = \infty$  there is some  $\mathbf{c} \in (0, C]$  such that  $T_* = 3/4$ . In particular, there are some controls  $(r, w, b)$  such that (C.7), (C.8) hold, and such that:

$$\|(w, b)\|_{L^\infty(0, 3/4)} < C \sum_{i=1}^N |z^i - x^i|.$$

Finally, arguing as in the final part of Step 2.3 we obtain for  $C$  large enough that:

$$|y(3/4; x^N, r, w, b)_2 - z_2^N| \leq \min \left\{ C \sum_{i=1}^N |z^i - x^i|, \frac{3\delta}{4} \right\}.$$

**Step 2.4: controlling  $y(\cdot; x^N, r, w, b)$  in  $[3/4, 1]$ .** In a similar way, we can prolong the controls in  $[3/4, 1]$  so that  $y(1; \mathbf{x}, r, w, b) = \mathbf{z}$  and:

$$\|(w, b)\|_{L^\infty(0, 1)} < C \sum_{i=1}^N |z^i - x^i|,$$

1 are satisfied. This can be done as in Step 2.3 by fixing  $r_1 = 1$ ,  $r_2 = \text{sign}(z_2 - x_2)$ ,  $w_1 = 0$ ,  
 2  $b_1 = 0$ ,  $w_2 = \mathbf{c} \sum_{i=1}^N |z^i - x^i| e$ ,  $b_2 = \mathbf{c} \sum_{i=1}^N |z^i - x^i| (-|z^N| + \delta)$  for some  $\mathbf{c} > 0$ . Indeed, since  
 3  $w_1 = 0$  and  $b_1 = 0$ , the function  $t \mapsto y(t; x^N, r, w, b)_1$  is constant in  $[3/4, 1]$ .  $\square$

## 4 References

- 5 [AB19] J. B. Amara and E. Beldi. Simultaneous controllability of two vibrating strings  
 6 with variable coefficients. *Evol. Equ. Control The.*, 8(4):687–694, 2019.
- 7 [BP20] T. Breiten and L. Pfeiffer. On the turnpike property and the receding-horizon  
 8 method for linear-quadratic optimal control problems. *SIAM J. Control Optim.*,  
 9 58(2):1077–1102, 2020.
- 10 [Ces12] L. Cesari. *Optimization—theory and applications: problems with ordinary differ-*  
 11 *ential equations*, volume 17. Springer Science & Business Media, 2012.
- 12 [DGSW14] T. Damm, L. Grüne, M. Stieler, and K. Worthmann. An exponential turnpike  
 13 theorem for dissipative discrete time optimal control problems. *SIAM J. Control*  
 14 *Optim.*, 52(3):1935–1957, 2014.
- 15 [EGPZ20a] C. Esteve, B. Geshkovski, D. Pighin, and E. Zuazua. Large-time asymptotics in  
 16 deep learning. *arXiv preprint arXiv:2008.02491*, 2020.

- 1 [EGPZ20b] C. Esteve, B. Geshkovski, D. Pighin, and E. Zuazua. Turnpike in lipschitz-nonlinear  
2 optimal control. *arXiv preprint arXiv:2011.11091*, 2020.
- 3 [EKPZ20] C. Esteve, H. Kouhkouh, D. Pighin, and E. Zuazua. The turnpike property  
4 and the long-time behavior of the Hamilton-Jacobi equation. *arXiv preprint*  
5 *arXiv:2006.10430*, 2020.
- 6 [EYG21] C. Esteve-Yagüe and B. Geshkovski. Sparse approximation in learning via neural  
7 odes. *arXiv preprint arXiv:2102.13566*, 2021.
- 8 [FHS21] T. Faulwasser, A.-J. Hempel, and S. Streif. On the turnpike to design of deep  
9 neural nets: Explicit depth bounds. *arXiv preprint arXiv:2101.03000*, 2021.
- 10 [FS18] E. Fathi and B. M. Shoja. Deep neural networks for natural language processing.  
11 In *Handbook of statistics*, volume 38, pages 229–316. Elsevier, 2018.
- 12 [GBB11] X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In  
13 *Proceedings of the fourteenth international conference on artificial intelligence and*  
14 *statistics*, pages 315–323. JMLR Workshop and Conference Proceedings, 2011.
- 15 [GG18] L. Gruüne and R. Guglielmi. Turnpike properties and strict dissipativity for dis-  
16 crete time linear quadratic optimal control problems. *SIAM J. Control Optim.*,  
17 56(2):1282–1302, 2018.
- 18 [GH19] M. Gugat and F. M. Hante. On the turnpike phenomenon for optimal boundary  
19 control problems with hyperbolic systems. *SIAM J. Control Optim.*, 57(1):264–289,  
20 2019.
- 21 [GSZ21] M. Gugat, M. Schuster, and E. Zuazua. The Finite-Time Turnpike Phenomenon  
22 for Optimal Control Problems: Stabilization by Non-smooth Tracking Terms. In  
23 G. Sklyar and A. Zuyev, editors, *Stabilization of Distributed Parameter Systems:*  
24 *Design Methods and Applications*, pages 17–41, Cham, 2021. Springer International  
25 Publishing.
- 26 [GTZ16] M. Gugat, R. Trélat, and E. Zuazua. Optimal Neumann control for the 1D wave  
27 equation: Finite horizon, infinite horizon, boundary tracking terms and the turn-  
28 pike property. *Syst. Control Lett.*, 90:61–70, 2016.
- 29 [HR17] E. Haber and L. Ruthotto. Stable architectures for deep neural networks. *Inverse*  
30 *Probl.*, 34(1):014004, 2017.
- 31 [HSV95] R. F Hartl, S. P. Sethi, and R. G. Vickson. A survey of the maximum principles for  
32 optimal control problems with state constraints. *SIAM Rev.*, 37(2):181–218, 1995.

- 1 [HZRS15] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing  
2 human-level performance on imagenet classification. In *Proceedings of the IEEE  
3 international conference on computer vision*, pages 1026–1034, 2015.
- 4 [Kom11] J. Komić. *International Encyclopedia of Statistical Science. Harmonic Mean.*, pages  
5 622–624. Springer, Heidelberg, 2011.
- 6 [Lio88] J-L Lions. Contrôlabilité exacte, perturbations et stabilisation de systèmes dis-  
7 tribués. Tome 1. *RMA*, 8, 1988.
- 8 [LZ16] J. Lohéac and E. Zuazua. From averaged to simultaneous controllability. In *Annales  
9 de la Faculté des sciences de Toulouse: Mathématiques*, volume 25, pages 785–828,  
10 2016.
- 11 [Mor14] M. Morancey. Simultaneous local exact controllability of 1D bilinear Schrödinger  
12 equations. *Ann. I. H. Poincaré-AN*, 31(3):501–529, 2014.
- 13 [MRB20] I. Mazari and D. Ruiz-Balet. Quantitative stability for eigenvalues of Schrödinger  
14 operator, Quantitative bathtub principle & Application to the turnpike property  
15 for a bilinear optimal control problem. *arXiv preprint arXiv:2010.10798*, 2020.
- 16 [MS95] J. Mira and F. Sandoval. *From Natural to Artificial Neural Computation: Interna-  
17 tional Workshop on Artificial Neural Networks, Malaga-Torremolinos, Spain, June  
18 7-9, 1995: Proceedings*, volume 930. Springer Science & Business Media, 1995.
- 19 [NH10] V. Nair and G. E Hinton. Rectified linear units improve restricted boltzmann ma-  
20 chines. In *27th International Conference on International Conference on Machine  
21 Learning, ICML 10*, pages 807–814, 2010.
- 22 [PZ13] A. Porretta and E. Zuazua. Long time versus steady state optimal control. *SIAM  
23 J. Control Optim.*, 51(6):4242–4273, 2013.
- 24 [RBZ21] D. Ruiz-Balet and E. Zuazua. Neural ODE control for classification, approximation  
25 and transport. *arXiv preprint arXiv:2104.05278*, 2021.
- 26 [Rus86] D. L Russell. The Dirichlet–Neumann boundary control problem associated with  
27 Maxwell’s equations in a cylindrical region. *SIAM J. Control Optim.*, 24(2):199–  
28 229, 1986.
- 29 [SN20] N. Sakamoto and M. Nagahara. The turnpike property in the maximum hands-off  
30 control. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pages  
31 2350–2355. IEEE, 2020.
- 32 [Son92] E. D. Sontag. Neural nets as systems models and controllers. In *Proc. Seventh  
33 Yale Workshop on Adaptive and Learning Systems*, pages 73–79, 1992.

- 1 [SPZ19] N. Sakamoto, D. Pighin, and E. Zuazua. The turnpike property in nonlinear op-  
2 timal control—a geometric approach. In *2019 IEEE 58th Conference on Decision*  
3 *and Control (CDC)*, pages 2422–2427. IEEE, 2019.
- 4 [SS97] E. Sontag and H. Sussmann. Complete controllability of continuous-time recurrent  
5 neural networks. *Syst. Control Lett.*, 30(4):177–183, 1997.
- 6 [Tré05] E. Trélat. *Contrôle optimal: théorie & applications*. Vuibert Paris, 2005.
- 7 [Tré20] E. Trélat. Linear turnpike theorem. *arXiv preprint arXiv:2010.13605*, 2020.
- 8 [TW00] M. Tucsnak and G. Weiss. Simultaneous exact controllability and some applica-  
9 tions. *SIAM J. Control Optim.*, 38(5):1408–1427, 2000.
- 10 [TZ15] E. Trélat and E. Zuazua. The turnpike property in finite-dimensional nonlinear  
11 optimal control. *J. Differ. Equations*, 258(1):81–114, 2015.
- 12 [TZZ18] E. Trélat, C. Zhang, and E. Zuazua. Steady-state and periodic exponential turnpike  
13 property for optimal control problems in Hilbert spaces. *SIAM J. Control Optim.*,  
14 56(2):1222–1252, 2018.
- 15 [Wei17] E. Weinan. A proposal on machine learning via dynamical systems. *Commun.*  
16 *Math. Stat.*, 5(1):1–11, 2017.
- 17 [WZ21] M. Warma and S. Zamorano. Exponential turnpike property for fractional parabolic  
18 equations with non-zero exterior data. *ESAIM:COCV*, 27(1):1–35, 2021.
- 19 [WZL17] J. Wu, X. Zhu, and S. Li. Simultaneous controllability of damped wave equations.  
20 *Math. Method Appl. Sci.*, 40(1):319–324, 2017.
- 21 [Yeh06] J. Yeh. *Real analysis: theory of measure and integration second edition*. World  
22 Scientific Publishing Company, 2006.
- 23 [Zam18] S. Zamorano. Turnpike property for two-dimensional Navier–Stokes equations. *J.*  
24 *Math. Fluid Mech.*, 20(3):869–888, 2018.
- 25 [Zbi93] R. Zbikowski. Lie algebra of recurrent neural networks and identifiability. In *1993*  
26 *American Control Conference*, pages 2900–2901. IEEE, 1993.