# UNIVERSAL APPROXIMATION OF INPUT-OUTPUT MAPS AND DYNAMICAL SYSTEMS BY NEURAL NETWORK ARCHITECTURES

BY

JOSHUA HANSON

THESIS

Submitted in partial fulfillment of the requirements
for the degree of Master of Science in Electrical and Computer Engineering
in the Graduate College of the
University of Illinois at Urbana-Champaign, 2020

Urbana, Illinois

Adviser:

Associate Professor Maxim Raginsky

# ABSTRACT

It is well known that feedforward neural networks can approximate any continuous function supported on a finite-dimensional compact set to arbitrary accuracy. However, many engineering applications require modeling infinite-dimensional functions, such as sequence-to-sequence transformations or input-output characteristics of systems of differential equations. For discrete-time input-output maps having limited long-term memory, we prove universal approximation guarantees for temporal convolutional nets constructed using only a finite number of computation units which hold on an infinite-time horizon. We also provide quantitative estimates for the width and depth of the network sufficient to achieve any fixed error tolerance. Furthemore, we show that discrete-time input-output maps given by state-space realizations satisfying certain stability criteria admit such convolutional net approximations which are accurate on an infinite-time scale. For continuous-time input-output maps induced by dynamical systems that are stable in a similar sense, we prove that continuous-time recurrent neural nets are capable of reproducing the original trajectories to within arbitrarily small error tolerance over an infinite-time horizon. For a subset of these stable systems, we provide quantitative estimates on the number of neurons sufficient to guarantee the desired error bound.

# ACKNOWLEDGMENTS

I would first like to thank my thesis adviser Associate Professor Maxim Raginsky for his instruction and guidance, which have contributed greatly to the development of my fundamental research skills and mathematical maturity. I would also like to acknowledge funding for my research and education from the National Science Foundation through the Center for Advanced Electronics through Machine Learning (CAEML) Industry - University Cooperative Research Program (I/UCRC) under award CNS-16-24811.

I also want to recognize my friend and colleague Yifeng Chu, who I have worked and studied alongside in Maxim's research group and the majority of our first- and second-year courses. My intuition and understanding of many foundational theorems and concepts have come from one or another of our numerous marathon homework sessions.

Lastly, I wish to express my appreciation for my mom and dad, who have always respected my interests and supported my inquisitiveness and fascination with the sciences from a young age.

# TABLE OF CONTENTS

# CHAPTER 1

# INTRODUCTION

Many sequence-to-sequence learning applications such as audio generation, natural language processing, time series forecasting, and system simulation are well-suited for parametric models that can be trained from measurement data with minimal knowledge of the original system structure. Traditionally, recurrent neural networks and other recursive architectures typified by internal state feedback have been favorable choices for approximating sequence-to-sequence transformations because they are naturally formulated for processing arbitrary length inputs. In comparison, fixed feedforward architectures must be designed for a predetermined input sequence length.

Convolutional neural networks and related convolutional architectures are also well-suited for variable length input sequences, but enjoy additional computational advantages during training and runtime that follow from the lack of feedback elements [1]. Without state feedback, shifted copies of the input sequence can be processed in parallel rather than consecutively. Furthermore, in practice, convolutional models have achieved competitive performance in tasks often approached with recurrent models [2, 3, 4, 5, 1, 6].

There is significant overlap in the properties of systems that are modeled efficiently by both convolutional and recurrent architectures. Recurrent models possess a theoretically unlimited memory because the output at a given time depends on the initial condition and the complete sequence of input values occurring up until that time. Therefore, such models are strictly more expressive than autoregressive models that only consider a finite number of past input values to determine the output at each time. Despite this infinite memory property, recurrent architectures often fail tests of their capability to learn artificially long input-output sequence relationships [7]. Additionally, infinite memory represents marginal value in practice [8, 1] as well as in theoretical considerations [9].

Recurrent models that exhibit exponential stability – that is, where the

dependence on the output of the initial condition decays exponentially in time – can be approximated arbitrarily accurately by feedforward models that only consider past input values occurring within a sufficiently large, but nonetheless finite, time horizon [10]. Additionally, for learning finite length input-output sequence relationships, feedforward models that implement temporal convolutions perform as well as or better than recurrent nets [2, 11, 7]. Evidently, a salient characteristic of systems that are well-suited for both convolutional and recurrent models is a limited long-term dependence on past input values.

By construction, the criterion for exponential stability of a recurrent model as considered in previous work [10] depends on the state-space representation. In order to verify if a system demonstrates this stability property, it is necessary to choose a particular realization in advance. This framework is not ideal because any sequence-to-sequence transformation can be assigned both stable and unstable state-space realizations. For example, a stable model can be made into an unstable one by augmenting the dynamics with unstable, unobservable states that have no influence on the output. Therefore, it is desired to formulate a more general characterization of limited long-term memory that abstracts away the notion of stability and instead uses intrinsic properties of the input-output map. This formalism can then be reconnected to the stability conditions expressed in the language of dynamical systems.

Characterizing limited long-term input dependence amounts to requiring that the output of the system depend predominantly on values of the input occurring within a short time horizon and negligibly on values of the input beyond that time horizon. This property can be defined more precisely as continuity of the input-output map with respect to a special norm that weights recent input values more heavily than past input values, called fading memory [12]. Expanding upon this definition, a more minimal characterization called approximately finite memory [13] eliminates the role of the weighted norm in fading memory. This property instead requires that the output of a system due to an infinitely long input sequence and the output due to the same input sequence after being truncated to a finite time horizon will be arbitrarily close given a sufficiently large horizon. Such systems can be modeled canonically by temporal convolutional nets, which by construction only operate on input values occurring within a finite horizon of the output time.

To facilitate a more rigorous comparison between recurrent and convolu-

tional models, we seek to quantitatively evaluate the approximation capability of temporal convolutional nets for modeling causal, time-invariant input-output maps with approximately finite memory and of recurrent neural nets for modeling such maps that admit state-space realizations (taking the form of a dynamical system). In Chapter 2, we establish the existence of the desired model and quantitative bounds on the context length, network dimensions, and total number of computation units sufficient to achieve any given error tolerance.

To compare the expressivity of temporal convolutional nets to recurrent models that demonstrate exponential stability, we can consider bounds developed in [10] which outline when truncated recurrent models approximate their untruncated counterparts. However, there still exist many recurrent models that demonstrate approximately finite memory, thus are comparable to convolutional models, but are not exponentially stable. To expand our comparison to these models, we adapt our analysis to a more relaxed incremental stability condition – described in [14] – which only requires the influence of the initial condition to be asymptotically negligible. Chapter 3 discusses the technical formulation of this property and its relationship to other stability conditions. Therein, we show that incremental stability implies approximately finite memory. For any recurrent system exhibiting this property, we derive bounds on the context length and dimensions of a temporal convolutional net which approximates the aforementioned recurrent system and derive the result of [10] as a special case.

Chapter 4 shifts focus to analyzing the expressivity of continuous-time recurrent neural nets for modeling incrementally stable systems. Existing results have already established that recurrent nets are capable of simulating trajectories to within any error tolerance over a finite time interval [15, 16, 17, 18]. These previous works apply Grönwall's inequality to control the difference between the paths of the original system and the simulated paths, which incurs an exponential degradation of the approximation accuracy over time. Without enforcing additional assumptions, simulating a system over a longer time horizon with the same error limit requires an exponentially more accurate approximation. This is problematic because the number of computation units in the network and number of training samples required to achieve the desired error tolerance will depend on the simulation time scale, which in many applications is unknown *a priori*.

However, for systems satisfying modest stability conditions, Grönwall's inequality is overly conservative, and a more detailed argument can prove that approximation error does not in general accumulate without bound. The same incremental stability property described earlier can be utilized to establish strict guarantees that the output of a simulating model will remain sufficiently close to the output of the original system over infinite time scales, rather than permitting performance degradation after a fixed time horizon. For stable systems satisfying the additional assumption that the gradient of the Fourier transform of the transition function is integrable, we derive quantitative bounds for the size of the recurrent net sufficient to achieve a desired error tolerance.

To summarize, given any input-output map having approximately finite memory, we provide quantitative estimates for how efficiently a temporal convolutional net can approximate such a map. In the special case where this input-output map is given explicitly by a state-space realization, we express our results in terms of stability conditions on the corresponding dynamical system. To compare convolutional to recurrent architectures, we can express the latter as a state-space realization and determine how large a convolutional model should be to approximate that realization. Finally, we study simulating continuous-time state-space models by recurrent neural nets and provide quantitative characterizations of the expressivity of this architecture for modeling a subset of stable dynamical systems.

# CHAPTER 2

# UNIVERSAL APPROXIMATION WITH TEMPORAL CONVOLUTIONAL NETS

Let $\mathcal{S}$ denote the set of all real-valued sequences $\mathbf{u} = (u_t)_{t \in \mathbb{Z}_+}$, where $\mathbb{Z}_+ := \{0, 1, 2, \ldots\}$. An *input-output map* (or i/o map, for short) is a nonlinear operator $\mathsf{F} : \mathcal{S} \to \mathcal{S}$ that maps an input sequence $\mathbf{u} \in \mathcal{S}$ to an output sequence $\mathbf{y} = \mathsf{F}\mathbf{u} \in \mathcal{S}$. (We are considering real-valued input and output sequences for simplicity; all our results carry over to vector-valued sequences at the expense of additional notation.) We will denote the application and the composition of i/o maps by concatenation. We are concerned with i/o maps $\mathsf{F}$ that are:

- *causal* — for any $t \in \mathbb{Z}_+$, $\mathbf{u}_{0:t} = \mathbf{v}_{0:t}$ implies $(\mathsf{F}\mathbf{u})_t = (\mathsf{F}\mathbf{v})_t$, where $\mathbf{u}_{0:t} := (u_0, \ldots, u_t)$;

- *time-invariant* — for any $k \in \mathbb{Z}_+$,

$$(\mathsf{F}\mathsf{R}^k \mathbf{u})_t = \begin{cases} (\mathsf{F}\mathbf{u})_{t-k}, & \text{for } t \geq k \\ 0, & \text{for } 0 \leq t < k \end{cases},$$

where $\mathsf{R} : \mathcal{S} \to \mathcal{S}$ is the right shift operator $(\mathsf{R}\mathbf{u})_t := u_{t-1}\mathbf{1}_{\{t \geq 1\}}$.

## 2.1 Approximately finite memory

A key notion we will work with is that of *approximately finite memory* [13]:

**Definition 2.1.** *An i/o map $\mathsf{F}$ has* approximately finite memory *on a set of inputs $\mathcal{M} \subseteq \mathcal{S}$ if for any $\epsilon > 0$ there exists $m \in \mathbb{Z}_+$, such that*

$$\sup_{\mathbf{u} \in \mathcal{M}} \sup_{t \in \mathbb{Z}_+} \left| (\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t \right| \leq \epsilon, \tag{2.1}$$

*where $\mathsf{W}_{t,m} : \mathcal{S} \to \mathcal{S}$ is the* windowing operator *$(\mathsf{W}_{t,m}\mathbf{u})_\tau := u_\tau \mathbf{1}_{\{\max\{t-m,0\} \leq \tau \leq t\}}$. We will denote by $m_{\mathsf{F}}^*(\epsilon)$ the smallest $m \in \mathbb{Z}_+$, for which (2.1) holds.*

If $m_{\mathsf{F}}^*(0) < \infty$, then we say that $\mathsf{F}$ has *finite memory* on $\mathcal{M}$. If $\mathsf{F}$ is causal and time-invariant, this is equivalent to the existence of an integer $m \in \mathbb{Z}_+$ and a nonlinear functional $f : \mathbb{R}^{m+1} \to \mathbb{R}$, such that $f(0, \ldots, 0) = 0$ and, for any $\mathbf{u} \in \mathcal{M}$ and any $t \in \mathbb{Z}_+$,

$$(\mathsf{F}\mathbf{u})_t = f(u_{t-m}, u_{t-m+1}, \ldots, u_t), \tag{2.2}$$

with the convention that $u_s = 0$ if $s < 0$. In this work, we will focus on the important case when $f$ is a feedforward neural net with rectified linear unit (ReLU) activations $\mathrm{ReLU}(x) := \max\{x, 0\}$. That is, there exist $k$ affine maps $A_i : \mathbb{R}^{d_i} \to \mathbb{R}^{d_{i+1}}$ with $d_1 = m + 1$ and $d_{k+1} = 1$, such that $f$ is given by the composition

$$f = A_k \circ \mathrm{ReLU} \circ A_{k-1} \circ \mathrm{ReLU} \circ \ldots \circ \mathrm{ReLU} \circ A_1,$$

where, for any $r \geq 1$, $\mathrm{ReLU}(x_1, \ldots, x_r) := (\mathrm{ReLU}(x_1), \ldots, \mathrm{ReLU}(x_r))$. Here, $k$ is the depth (number of layers) and $\max\{d_2, \ldots, d_k\}$ is the width (largest number of units in any hidden layer).

**Definition 2.2.** *An i/o map $\mathsf{F}$ is a ReLU temporal convolutional net (or ReLU TCN, for short) with context length $m$ if* (2.2) *holds for some feedforward ReLU neural net $f : \mathbb{R}^{m+1} \to \mathbb{R}$.*

**Remark 2.1.** While such an $\mathsf{F}$ is evidently causal, it is generally not time-invariant unless $f(0, \ldots, 0) = 0$.

## 2.2 The universal approximation theorem for input-output maps

In this section, we state and prove one of our main results: Any causal and time-invariant i/o map that has approximately finite memory and satisfies an additional continuity condition can be approximated arbitrarily well by a ReLU temporal convolutional net. In what follows, we will consider i/o maps with uniformly bounded inputs, i.e., inputs in the set

$$\mathcal{M}(R) := \{\mathbf{u} \in \mathcal{S} : \|\mathbf{u}\|_\infty := \sup_{t \in \mathbb{Z}_+} |u_t| \leq R\} \qquad \text{for some } R > 0.$$

For any $t \in \mathbb{Z}_+$ and any $\mathbf{u} \in \mathcal{M}(R)$, the finite subsequence $\mathbf{u}_{0:t} = (u_0, \ldots, u_t)$ can be considered as an element of the cube $[-R, R]^{t+1} \subset \mathbb{R}^{t+1}$; conversely, any vector $\mathbf{x} \in [-R, R]^{t+1}$ can be embedded into the sequence space $\mathcal{M}(R)$ by setting $u_s = x_s \mathbf{1}_{\{0 \leq s \leq t\}}$. To any causal and time-invariant i/o map $\mathsf{F}$ we can associate the nonlinear functional $\tilde{\mathsf{F}}_t : \mathbb{R}^{t+1} \to \mathbb{R}$ defined in the obvious way: for any $\mathbf{x} = (x_0, x_1, \ldots, x_t) \in \mathbb{R}^{t+1}$,

$$\tilde{\mathsf{F}}_t(\mathbf{x}) := (\mathsf{F}\mathbf{u})_t,$$

where $\mathbf{u} \in \mathcal{S}$ is any input such that $u_s = x_s$ for $s \in \{0, 1, \ldots, t\}$ (the values of $u_s$ for $s > t$ can be arbitrary by causality). We impose the following assumptions on $\mathsf{F}$:

**Assumption 2.1.** *The i/o map $\mathsf{F}$ has approximately finite memory on $\mathcal{M}(R)$.*

**Assumption 2.2.** *For any $t \in \mathbb{Z}_+$, the functional $\tilde{\mathsf{F}}_t : \mathbb{R}^{t+1} \to \mathbb{R}$ is uniformly continuous on $[-R, R]^{t+1}$ with modulus of continuity*

$$\omega_{t,\mathsf{F}}(\delta) := \sup \left\{ |\tilde{\mathsf{F}}_t(\mathbf{x}) - \tilde{\mathsf{F}}_t(\mathbf{x}')| : \mathbf{x}, \mathbf{x}' \in [-R, R]^{t+1}, \|\mathbf{x} - \mathbf{x}'\|_\infty \leq \delta \right\},$$

*and inverse modulus of continuity*

$$\omega_{t,\mathsf{F}}^{-1}(\epsilon) := \sup \left\{ \delta > 0 : \omega_{t,\mathsf{F}}(\delta) \leq \epsilon \right\},$$

*where $\|\mathbf{x}\|_\infty := \max_{0 \leq i \leq t} |x_i|$ is the $\ell^\infty$ norm on $\mathbb{R}^{t+1}$.*

The following qualitative universal approximation result was obtained by Sandberg in [13]: if a causal and time-invariant i/o map $\mathsf{F}$ satisfies the above two assumptions, then, for any $\epsilon > 0$, there exists an affine map $A : \mathbb{R}^{m+1} \to \mathbb{R}^d$ and a lattice map $\ell : \mathbb{R}^d \to \mathbb{R}$, such that

$$\sup_{\mathbf{u} \in \mathcal{M}(R)} \sup_{t \in \mathbb{Z}_+} \left| (\mathsf{F}\mathbf{u})_t - \ell \circ A(\mathbf{u}_{t-m:t}) \right| < \epsilon, \tag{2.3}$$

where we say that a map $\ell : \mathbb{R}^d \to \mathbb{R}$ is a *lattice map* if $\ell(x_0, \ldots, x_{d-1})$ is generated from $x = (x_0, \ldots, x_{d-1})$ by a finite number of min and max operations that do not depend on $x$. Any lattice map can be implemented using ReLU units, so (2.3) is a ReLU TCN approximation guarantee. The main result of this chapter is a quantitative version of Sandberg's theorem:

7

**Theorem 2.1.** *Let* $\mathsf{F}$ *be a causal and time-invariant i/o map satisfying Assumptions 2.1 and 2.2. Then, for any $\epsilon > 0$ and any $\gamma \in (0,1)$, there exists a ReLU TCN $\widehat{\mathsf{F}}$ with*

- *context length $m = m_{\mathsf{F}}^*(\gamma\epsilon)$*

- *width $m+2$*

- *depth $\left(\frac{O(R)}{\omega_{m,\mathsf{F}}^{-1}((1-\gamma)\epsilon)}\right)^{m+2}$*

*such that*

$$\sup_{\mathbf{u}\in\mathcal{M}(R)} \|\mathsf{F}\mathbf{u} - \widehat{\mathsf{F}}\mathbf{u}\|_\infty < \epsilon. \tag{2.4}$$

**Remark 2.2.** The role of the additional parameter $\gamma \in (0,1)$ is to trade off the context length and the depth of the ReLU TCN.

**Remark 2.3.** While the approximating ReLU TCN $\widehat{\mathsf{F}}$ is clearly causal, it may not be time-invariant unless $\widehat{f}(0,\ldots,0) = 0$, where $\widehat{f}$ is the ReLU net constructed in the proof below.

*Proof.* Let $m = m_{\mathsf{F}}^*(\gamma\epsilon)$. Since $\tilde{\mathsf{F}}_m : \mathbb{R}^{m+1} \to \mathbb{R}$ is continuous with modulus of continuity $\omega_{m,\mathsf{F}}(\cdot)$, there exists a ReLU net $\widehat{f} : \mathbb{R}^{m+1} \to \mathbb{R}$ of width $m+2$ and depth $\left(\frac{O(R)}{\omega_{m,\mathsf{F}}^{-1}((1-\gamma)\epsilon)}\right)^{m+2}$, such that

$$\sup_{\mathbf{x}\in[-R,R]^{m+1}} |\tilde{\mathsf{F}}_m(\mathbf{x}) - \widehat{f}(\mathbf{x})| < (1-\gamma)\epsilon,$$

as proved by Hanin and Sellke [19]. Now consider the TCN $\widehat{\mathsf{F}}$ defined by $(\mathsf{F}\mathbf{u})_t := \widehat{f}(u_{t-m}, \ldots, u_t)$. Fix an input $\mathbf{u} \in \mathcal{M}(R)$ and consider two cases:

1) If $t \geq m$, then $\mathbf{u}_{t-m:t} = (\mathsf{L}^{t-m}\mathsf{W}_{t,m}\mathbf{u})_{0:m}$, where $\mathsf{L} : \mathcal{S} \to \mathcal{S}$ is the left shift operator $(\mathsf{L}\mathbf{u})_t := u_{t+1}$. Therefore,

$$(\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t \overset{(a)}{=} (\mathsf{F}\mathsf{R}^{t-m}\mathsf{L}^{t-m}\mathsf{W}_{t,m}\mathbf{u})_t \overset{(b)}{=} (\mathsf{F}\mathsf{L}^{t-m}\mathsf{W}_{t,m}\mathbf{u})_m \overset{(c)}{=} \tilde{\mathsf{F}}_m(\mathbf{u}_{t-m:t}),$$

where (a) uses the fact that $t \geq m$, (b) is by time invariance of $\mathsf{F}$, and (c) is by the definition of $\tilde{\mathsf{F}}_m$.

2) If $t < m$, then $\mathbf{u}_{t-m:t} = (\mathsf{R}^{m-t}\mathsf{W}_{t,m}\mathbf{u})_{0:m}$ (recall the convention that, for any $\mathbf{v}$, we set $v_s = 0$ whenever $s < 0$). Therefore

$$(\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t \overset{(a)}{=} (\mathsf{R}^{m-t}\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_m \overset{(b)}{=} (\mathsf{F}\mathsf{R}^{m-t}\mathsf{W}_{t,m}\mathbf{u})_m \overset{(c)}{=} \tilde{\mathsf{F}}_m(\mathbf{u}_{t-m:t}),$$

8

where (a) uses the fact that $m > t$, (b) is by time invariance, and (c) is by the definition of $\tilde{\mathsf{F}}_m$.

In either case, the triangle inequality gives

$$
\begin{aligned}
|(\mathsf{F}\mathbf{u})_t - (\widehat{\mathsf{F}}\mathbf{u})_t| &\leq |(\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t| + |(\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t - (\widehat{\mathsf{F}}\mathbf{u})_t| \\
&= |(\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t| + |\tilde{\mathsf{F}}_m(\mathbf{u}_{t-m:t}) - \widehat{f}(\mathbf{u}_{t-m:t})| \\
&< \gamma\epsilon + (1-\gamma)\epsilon = \epsilon.
\end{aligned}
$$

Since this holds for all $t$ and all $\mathbf{u}$ with $\|\mathbf{u}\|_\infty \leq R$, the result follows. $\qquad\square$

## 2.3   The fading memory property

In order to apply Theorem 2.1, we need control over the context length $m_{\mathsf{F}}^*(\cdot)$ and over the modulus of continuity $\omega_{t,\mathsf{F}}(\cdot)$. In general, these quantities are difficult to estimate. However, it was shown by Park and Sandberg in [20] that the property of approximately finite memory is closely related to the notion of *fading memory*, first introduced by Boyd and Chua in [12]. Intuitively, an i/o map $\mathsf{F}$ has fading memory if the outputs at any time $t$ due to any two inputs $\mathbf{u}$ and $\mathbf{v}$ that were close to one another in recent past will also be close.

Let $\mathcal{W}$ denote the subset of $\mathcal{S}$ consisting of all sequences $\mathbf{w}$, such that $w_t \in (0, 1]$ for all $t$ and $w_t \downarrow 0$ as $t \to \infty$. We will refer to the elements of $\mathcal{W}$ as *weighting sequences*. Then we have the following definition, due to [20]:

**Definition 2.3.** *We say that an i/o map $\mathsf{F}$ has* fading memory *on $\mathcal{M} \subseteq \mathcal{S}$ with respect to $\mathbf{w} \in \mathcal{W}$ if for any $\epsilon > 0$ there exists $\delta > 0$ such that, for all $\mathbf{u}, \mathbf{v} \in \mathcal{M}$ and all $t \in \mathbb{Z}_+$,*

$$
\max_{s \in \{0,\dots,t\}} w_{t-s}|u_s - v_s| < \delta \quad \implies \quad |(\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathbf{v})_t| < \epsilon. \tag{2.5}
$$

The weighting sequence $\mathbf{w}$ governs the rate at which the past values of the input are discounted in determining the current output. To capture the best trade-offs in (2.5), we will also use a $\mathbf{w}$-dependent modulus of continuity:

$$
\alpha_{\mathbf{w},\mathsf{F}}(\delta) := \sup\left\{|(\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathbf{v})_t| : t \in \mathbb{Z}_+, \mathbf{u}, \mathbf{v} \in \mathcal{M}, \max_{s \in \{0,\dots,t\}} w_{t-s}|u_s - v_s| \leq \delta\right\}.
$$

It was shown by Park and Sandberg in [20] that an i/o map satisfies Assump-

tions 2.1 and (2.2) if and only if it has fading memory with respect to some (and hence any) $\mathbf{w} \in \mathcal{W}$. The following result provides a quantitative version of this equivalence:

**Proposition 2.1.** *Let* $\mathsf{F}$ *be an i/o map.*

1. *If* $\mathsf{F}$ *satisfies Assumptions 2.1 and 2.2, then it has fading memory on* $\mathcal{M}$ *with respect to any weighting sequence* $\mathbf{w} \in \mathcal{W}$, *and*

$$\alpha_{\mathbf{w},\mathsf{F}}^{-1}(\epsilon) \geq w_{m_\mathsf{F}^*(\epsilon/3)} \omega_{m_\mathsf{F}^*(\epsilon/3),\mathsf{F}}^{-1}(\epsilon/3). \tag{2.6}$$

2. *If* $\mathsf{F}$ *has fading memory on* $\mathcal{M}(R)$ *with respect to some* $\mathbf{w} \in \mathcal{W}$, *then it satisfies Assumptions 2.1 and 2.2, and*

$$m_\mathsf{F}^*(\epsilon; R) \leq \inf\left\{m \in \mathbb{Z}_+ : w_m \leq \frac{\alpha_{\mathbf{w},\mathsf{F}}^{-1}(\epsilon)}{R}\right\} \quad \text{and} \quad \omega_{t,\mathsf{F}}(\delta) \leq \alpha_{\mathbf{w},\mathsf{F}}(\delta). \tag{2.7}$$

*Proof.* Suppose $\mathsf{F}$ satisfies Assumptions 2.1 and 2.2. Fix some $\epsilon > 0$ and let $m = m_\mathsf{F}^*(\epsilon/3)$ and $\delta = w_m \omega_{m,\mathsf{F}}^{-1}(\epsilon/3)$. Now fix some $t \in \mathbb{Z}_+$ and consider any two $\mathbf{u}, \mathbf{v} \in \mathcal{M}(R)$ such that

$$\max_{s \in \{0,\ldots,t\}} w_{t-s}|u_s - v_s| < \delta. \tag{2.8}$$

Using the same reasoning as in the proof of Theorem 2.1, we can write $(\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t = \tilde{\mathsf{F}}_m(\mathbf{u}_{t-m:t})$ and $(\mathsf{F}\mathsf{W}_{t,m}\mathbf{v})_t = \tilde{\mathsf{F}}_m(\mathbf{v}_{t-m:t})$, where, as before, we set $u_s = v_s = 0$ for $s < 0$. From the monotonicity of $\mathbf{w}$ and (2.8) it follows that

$$\|\mathbf{u}_{t-m:t} - \mathbf{v}_{t-m:t}\|_\infty \leq \frac{1}{w_m} \max_{s \in \{t-m,\ldots,t\}} w_{t-s}|u_s - v_s| < \omega_{m,\mathsf{F}}^{-1}(\epsilon/3),$$

which implies that

$$|(\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t - (\mathsf{F}\mathsf{W}_{t,m}\mathbf{v})_t| = |\tilde{\mathsf{F}}_m(\mathbf{u}_{t-m:t}) - \tilde{\mathsf{F}}_m(\mathbf{v}_{t-m:t})| < \epsilon/3.$$

Altogether, we see that (2.8) implies that

$$|(\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathbf{v})_t|$$
$$\leq |(\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t| + |(\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t - (\mathsf{F}\mathsf{W}_{t,m}\mathbf{v})_t| + |(\mathsf{F}\mathbf{v})_t - (\mathsf{F}\mathsf{W}_{t,m}\mathbf{v})_t|$$
$$< \epsilon/3 + \epsilon/3 + \epsilon/3 = \epsilon,$$

which leads to (2.6).

Now suppose that $\mathsf{F}$ has fading memory with respect to $\mathbf{w}$. Given $\epsilon > 0$, let $\delta = \alpha_{\mathbf{w},\mathsf{F}}^{-1}(\epsilon)$ and choose any $m \in \mathbb{Z}_+$, such that $w_m < \delta/R$. If $t < m$, then $\mathbf{u}_{0:t} = (\mathsf{W}_{t,m}\mathbf{u})_{0:t}$, and thus $(\mathsf{F}\mathbf{u})_t = (\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t$. On the other hand, if $t \geq m$, then, for any $\mathbf{u} \in \mathcal{M}(R)$,

$$\max_{s \in \{0,\ldots,t\}} |u_s - (\mathsf{W}_{t,m}\mathbf{u})_s| = \begin{cases} 0, & t - m \leq s \leq t \\ |u_s|, & s < t - m \end{cases}$$

and therefore, by the monotonicity of $\mathbf{w}$ and the choice of $m$,

$$\max_{s \in \{0,\ldots,t\}} w_{t-s} |u_s - (\mathsf{W}\mathbf{u}_{t,m})_s| = \max_{s < t-m} w_{t-s} |u_s| \leq w_m \|\mathbf{u}\|_\infty < \delta,$$

which implies that $|(\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t| < \epsilon$. Consequently, $m_{\mathsf{F}}^*(\epsilon) \leq m$. Moreover, since the elements of $\mathbf{w}$ take values in $(0, 1]$, it follows from definitions that, for any $\mathbf{u}, \mathbf{v} \in \mathcal{M}(R)$ and any $t$,

$$\|\mathbf{u}_{0:t} - \mathbf{v}_{0:t}\|_\infty < \delta$$
$$\implies \max_{s \in \{0,\ldots,t\}} w_{t-s} |u_s - v_s| < \delta$$
$$\implies |(\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathbf{v})_t| \leq \alpha_{\mathbf{w},\mathsf{F}}(\delta).$$

This establishes (2.7). □

11

# CHAPTER 3

# DYNAMICAL SYSTEMS AND INCREMENTAL STABILITY

So far, we have considered arbitrary i/o maps $\mathsf{F} : \mathcal{S} \to \mathcal{S}$. However, many such maps admit *state-space realizations* [21] — that is, there exist a state transition map $f : \mathbb{R}^n \times \mathbb{R} \to \mathbb{R}^n$, an output map $g : \mathbb{R}^n \to \mathbb{R}$, and an initial condition $\xi \in \mathbb{R}^n$, such that the output sequence $\mathbf{y} = \mathsf{F}\mathbf{u}$ is determined recursively by the discrete-time dynamical system

$$x_{t+1} = f(x_t, u_t) \tag{3.1a}$$

$$y_t = g(x_t) \tag{3.1b}$$

with $x_0 = \xi$. The i/o map $\mathsf{F}$ realized in this way is evidently causal, and it is time-invariant if $f(\xi, 0) = \xi$ and $g(\xi) = 0$. In this chapter, we will identify the conditions under which such recurrent models satisfy Assumptions 2.1 and 2.2. Along the way, we will derive the approximation results of Miller and Hardt in [10] as a special case.

## 3.1 Approximately finite memory and incremental stability

Consider the system in (3.1). Given any input $\mathbf{u} \in \mathcal{S}$, any $\xi \in \mathbb{R}^n$, and any $s, t \in \mathbb{Z}_+$ with $t \geq s$, we denote by $\varphi_{s,t}^{\mathbf{u}}(\xi)$ the state at time $t$ when $x_s = \xi$. Let $\mathcal{M}$ be a subset of $\mathcal{S}$. We say that $\mathbb{X} \subseteq \mathbb{R}^n$ is a *positively invariant set* of (3.1) for inputs in $\mathcal{M}$ if, for all $\xi \in \mathbb{X}$, all $\mathbf{u} \in \mathcal{M}$, and all $0 \leq s \leq t$, $\varphi_{s,t}^{\mathbf{u}}(\xi) \in \mathbb{X}$. We will be interested in systems with the following property [14]:

**Definition 3.1.** *The system* (3.1) *is* uniformly asymptotically incrementally stable *for inputs in $\mathcal{M}$ on a positively invariant set $\mathbb{X}$ if there exists a function*

12

$\beta : \mathbb{R}_+ \times \mathbb{R}_+ \to \mathbb{R}_+$ *of class* $\mathcal{KL}$,[1] *such that the inequality*

$$\|\varphi_{s,t}^{\mathbf{u}}(\xi) - \varphi_{s,t}^{\mathbf{u}}(\xi')\| \leq \beta(\|\xi - \xi'\|, t - s) \qquad (3.2)$$

*holds for all inputs* $\mathbf{u} \in \mathcal{M}$, *all initial conditions* $\xi, \xi' \in \mathbb{X}$, *and all* $0 \leq s \leq t$, *where* $\| \cdot \|$ *is the* $\ell^2$ *norm on* $\mathbb{R}^n$.

In other words, a system is incrementally stable if the influence of any initial condition in $\mathbb{X}$ on the state trajectory is asymptotically negligible. A key consequence is the following estimate:

**Proposition 3.1.** *Let* $\mathbf{u}, \tilde{\mathbf{u}}$ *be two input sequences in* $\mathcal{M}$. *Then, for any* $\xi \in \mathbb{X}$ *and any* $t \in \mathbb{Z}_+$,

$$\|\varphi_{0,t}^{\mathbf{u}}(\xi) - \varphi_{0,t}^{\tilde{\mathbf{u}}}(\xi)\| \leq \sum_{s=0}^{t-1} \beta \left( \|f(\tilde{x}_s, u_s) - f(\tilde{x}_s, \tilde{u}_s)\|, t - s - 1 \right), \qquad (3.3)$$

*where* $x_s$ *and* $\tilde{x}_s$ *denote the states at time* $s$ *due to inputs* $\mathbf{u}$ *and* $\tilde{\mathbf{u}}$, *respectively, with* $x_0 = \tilde{x}_0 = \xi$.

*Proof of Proposition 3.1.* The family of mappings $\varphi_{s,t}^{\mathbf{u}}(\cdot)$ has the following *semiflow property*: For any input $\mathbf{u}$ and any $0 \leq r \leq s \leq t$,

$$\varphi_{r,t}^{\mathbf{u}}(\xi) = \varphi_{s,t}^{\mathbf{u}}(\varphi_{r,s}^{\mathbf{u}}(\xi)). \qquad (3.4)$$

By telescoping and by the semiflow property (3.4), we have

$$\varphi_{0,t}^{\mathbf{u}}(\xi) - \varphi_{0,t}^{\tilde{\mathbf{u}}}(\xi) = \sum_{s=0}^{t-1} \left( \varphi_{s,t}^{\mathbf{u}}(\varphi_{0,s}^{\tilde{\mathbf{u}}}(\xi)) - \varphi_{s+1,t}^{\mathbf{u}}(\varphi_{0,s+1}^{\tilde{\mathbf{u}}}(\xi)) \right)$$

$$= \sum_{s=0}^{t-1} \left( \varphi_{s+1,t}^{\mathbf{u}}(\varphi_{s,s+1}^{\mathbf{u}}(\varphi_{0,s}^{\tilde{\mathbf{u}}}(\xi))) - \varphi_{s+1,t}^{\mathbf{u}}(\varphi_{0,s+1}^{\tilde{\mathbf{u}}}(\xi)) \right). \qquad (3.5)$$

Using the fact that $\varphi_{s,s+1}^{\mathbf{u}}(\varphi_{0,s}^{\tilde{\mathbf{u}}}(\xi)) = \varphi_{s,s+1}^{\mathbf{u}}(f(\varphi_{0,s}^{\tilde{\mathbf{u}}}(\xi), u_s))$ and the stability property (3.2),

$$\left\| \varphi_{s+1,t}^{\mathbf{u}}(\varphi_{s,s+1}^{\mathbf{u}}(\varphi_{0,s}^{\tilde{\mathbf{u}}}(\xi))) - \varphi_{s+1,t}^{\mathbf{u}}(\varphi_{0,s+1}^{\tilde{\mathbf{u}}}(\xi)) \right\| \leq \beta \left( \|f(\tilde{x}_s, u_s) - f(\tilde{x}_s, \tilde{u}_s)\|, t - s - 1 \right).$$

---

[1]A function $\beta : \mathbb{R}_+ \times \mathbb{R}_+ \to \mathbb{R}_+$ is of class $\mathcal{KL}$ if it is continuous and strictly increasing in its first argument, continuous and strictly decreasing in its second argument, $\beta(0, t) = 0$ for any $t$, and $\lim_{t \to \infty} \beta(r, t) = 0$ for any $r$ [21].

Substituting this into (3.5), we get (3.3). □

Consider a state-space model (3.1) with a positively invariant set $\mathbb{X}$, with the following assumptions:

**Assumption 3.1.** *The state transition map $f(x, u)$ is $L_f$-Lipschitz in $u$ for all $x \in \mathbb{X}$ and the output map $g(x)$ is $L_g$-Lipschitz in $x \in \mathbb{X}$.*

**Assumption 3.2.** *For any initial condition $\xi \in \mathbb{X}$ there exists a compact set $\mathbb{S}_\xi \subseteq \mathbb{X}$ such that $\varphi^{\mathbf{u}}_{0,t}(\xi) \in \mathbb{S}_\xi$ for all $\mathbf{u} \in \mathcal{M}(R)$ and all $t \in \mathbb{Z}_+$.*

**Assumption 3.3.** *The system (3.1) is uniformly asymptotically incrementally stable on $\mathbb{X}$ for inputs in $\mathcal{M}(R)$, and the function $\beta$ in (3.2) satisfies the summability condition*

$$\sum_{t \in \mathbb{Z}_+} \beta(C, t) < \infty \tag{3.6}$$

*for any $C \geq 0$. (For example, if $\beta(C, k) = C k^{-\alpha}$ for some $\alpha > 1$, then this condition is satisfied.)*

We are now in position to prove the main result of this section:

**Theorem 3.1.** *Suppose that Assumptions 3.1–3.3 are satisfied. Then the i/o map $\mathsf{F}$ of the system (3.1) satisfies Assumptions 2.1 and 2.2 with*

$$m^*_{\mathsf{F}}(\epsilon) \leq \min \left\{ m \in \mathbb{Z}_+ : \sum_{k \geq m} \beta(\mathrm{diam}(\mathbb{S}_\xi), k) < \epsilon / L_g \right\} \tag{3.7}$$

*and*

$$\omega_{t,\mathsf{F}}(\delta) \leq L_g \sum_{s=0}^{t-1} \beta(L_f \delta, s), \qquad \forall t \in \mathbb{Z}_+. \tag{3.8}$$

*Proof.* Fix some $t, m \in \mathbb{Z}_+$. For an arbitrary input $\mathbf{u} \in \mathcal{M}(R)$, let $\tilde{\mathbf{u}} = \mathsf{W}_{t,m}\mathbf{u}$, where we may assume without loss of generality that $t \geq m$. Then we have

$\tilde{u}_s = u_s \mathbf{1}_{\{t-m \leq s \leq t\}}$, and therefore

$$\sum_{s=0}^{t-1} \beta\left(\|f(\tilde{x}_s, u_s) - f(\tilde{x}_s, \tilde{u}_s)\|, t-s-1\right)$$

$$= \sum_{s=0}^{t-m-1} \beta\left(\|f(\tilde{x}_s, u_s) - f(\tilde{x}_s, 0)\|, t-s-1\right)$$

$$\leq \sum_{s=0}^{t-m-1} \beta(\operatorname{diam}(\mathbb{S}_\xi), t-s-1)$$

$$= \sum_{s=m}^{t-1} \beta(\operatorname{diam}(\mathbb{X}), s)$$

$$\leq \sum_{s=m}^{\infty} \beta(\operatorname{diam}(\mathbb{S}_\xi), s). \tag{3.9}$$

By the summability condition (3.6), the summation in (3.9) converges to 0 as $m \uparrow \infty$. Thus, if we choose $m$ so that the right-hand side of (3.9) is smaller than $\epsilon/L_g$, it follows from Proposition 3.1 that

$$|(\mathsf{F}\mathbf{u})_t - (\mathsf{F}\mathsf{W}_{t,m}\mathbf{u})_t| = |g(\varphi_{0,t}^{\mathbf{u}}(\xi)) - g(\varphi_{0,t}^{\tilde{\mathbf{u}}}(\xi))| \leq L_g \|\varphi_{0,t}^{\mathbf{u}}(\xi) - \varphi_{0,t}^{\tilde{\mathbf{u}}}(\xi)| < \epsilon.$$

This proves (3.7). Now fix any two $\mathbf{u}, \tilde{\mathbf{u}} \in \mathcal{M}(R)$ with $\|\mathbf{u}_{0:t} - \tilde{\mathbf{u}}_{0:t}\|_\infty < \delta$. Then $\max_{0 \leq s \leq t} \|f(x, u_s) - f(x, \tilde{u}_s)\| \leq L_f \delta$ for all $x \in \mathbb{X}$, so Proposition 3.1 gives

$$\begin{aligned}
|\tilde{\mathsf{F}}_t(\mathbf{u}_{0:t}) - \tilde{\mathsf{F}}_t(\tilde{\mathbf{u}}_{0:t})| &= |g(\varphi_{0,t}^{\mathbf{u}}(\xi)) - g(\varphi_{0,t}^{\tilde{\mathbf{u}}}(\xi))| \\
&\leq L_g \|\varphi_{0,t}^{\mathbf{u}}(\xi) - \varphi_{0,t}^{\tilde{\mathbf{u}}}(\xi)\| \\
&\leq L_g \sum_{s=0}^{t-1} \beta(L_f \delta, s),
\end{aligned}$$

which proves (3.8). $\qquad\square$

15

## 3.2 Exponential incremental stability and the Demidovich criterion

Miller and Hardt [10] consider the case of contracting systems: There exists some $\lambda \in (0,1)$ and a set $\mathbb{U} \subseteq \mathbb{R}^m$, such that

$$\|f(x,u) - f(x',u)\| \leq \lambda \|x - x'\| \tag{3.10}$$

for all $x, x' \in \mathbb{R}^n$ and all $u \in \mathbb{U}$. Such a system is *uniformly exponentially incrementally stable* on any positively invariant set $\mathbb{X}$, with $\beta(C,t) = C\lambda^t$. In this section, we obtain their result as a special case of a more general stability criterion, known in the literature on nonlinear system stability as the *Demidovich criterion* [22]. The following result is a simplified version of a more general result of [14]:

**Proposition 3.2** (the discrete-time Demidovich criterion)**.** *Consider the recurrent system* (3.1) *with a convex positively invariant set* $\mathbb{X}$*, where the state transition map* $f(x,u)$ *is differentiable in* $x$ *for any* $u \in \mathbb{U}$*. Suppose that there exists a symmetric positive definite matrix* $P$ *and a constant* $\mu \in (0,1)$*, such that*

$$\frac{\partial}{\partial x} f(x,u)^\top P \frac{\partial}{\partial x} f(x,u) - \mu P \preceq 0 \tag{3.11}$$

*for all* $x \in \mathbb{X}$ *and all* $u \in \mathbb{U}$*, where* $\frac{\partial}{\partial x} f(x,u)$ *is the Jacobian of* $f(\cdot, u)$ *with respect to* $x$*. Then the system* (3.1) *is uniformly exponentially incrementally stable with* $\beta(C,t) = \sqrt{\kappa(P)} C \mu^{t/2}$*, where* $\kappa(P)$ *is the condition number of* $P$*.*

*Proof.* Fix any $u \in \mathbb{U}$ and $\xi, \xi' \in \mathbb{X}$, and define the function $\Phi : [0,1] \to \mathbb{R}$ by

$$\Phi(s) := (f(\xi,u) - f(\xi',u))^\top P f(s\xi + (1-s)\xi', u).$$

Then

$$\Phi(1) - \Phi(0) = (f(\xi,u) - f(\xi',u))^\top P (f(\xi,u) - f(\xi',u)). \tag{3.12}$$

16

By the mean-value theorem, there exists some $\bar{s} \in [0, 1]$, such that

$$\Phi(1) - \Phi(0) = \frac{\mathrm{d}}{\mathrm{d}s}\Phi(s)\Big|_{s=\bar{s}} = (f(\xi, u) - f(\xi', u))^\top P \frac{\partial}{\partial x} f(\bar{\xi}, u)(\xi - \xi'), \tag{3.13}$$

where $\bar{\xi} = \bar{s}\xi + (1 - \bar{s})\xi' \in \mathbb{X}$, since $\mathbb{X}$ is convex. From (3.11), (3.12), and (3.13) it follows that

$$\begin{aligned}
(f(\xi, u) &- f(\xi', u))^\top P(f(\xi, u) - f(\xi', u)) \\
&\leq (\xi - \xi')^\top \frac{\partial}{\partial x} f(\bar{\xi}, u)^\top P \frac{\partial}{\partial x} f(\bar{\xi}, u)(\xi - \xi') \\
&\leq \mu(\xi - \xi')^\top P(\xi - \xi').
\end{aligned}$$

Define the function $V : \mathbb{X} \times \mathbb{X} \to \mathbb{R}_+$ by $V(\xi, \xi') := (\xi - \xi')^\top P(\xi - \xi')$. From the above estimate, it follows that $V$ is a *Lyapunov function* for the dynamics, i.e., for any $u \in \mathbb{U}$ and $\xi, \xi' \in \mathbb{X}$,

$$V(f(\xi, u), f(\xi', u)) \leq \mu V(\xi, \xi'). \tag{3.14}$$

Consequently, for any input $\mathbf{u}$ with $u_t \in \mathbb{U}$ for all $t$ and any $\xi, \xi' \in \mathbb{X}$,

$$\begin{aligned}
V(\varphi^{\mathbf{u}}_{0,t+1}(\xi), \varphi^{\mathbf{u}}_{0,t+1}(\xi')) &= V(f(\varphi^{\mathbf{u}}_{0,t}(\xi), u_t), f(\varphi^{\mathbf{u}}_{0,t}(\xi'), u_t)) \\
&\leq \mu V(\varphi^{\mathbf{u}}_{0,t}(\xi), \varphi^{\mathbf{u}}_{0,t}(\xi')).
\end{aligned}$$

Iterating, we obtain the inequality $V(\varphi^{\mathbf{u}}_{0,t}(\xi), \varphi^{\mathbf{u}}_{0,t}(\xi')) \leq \mu^t V(\xi, \xi')$. Finally, since $P \succ 0$,

$$\|\varphi^{\mathbf{u}}_{0,t}(\xi) - \varphi^{\mathbf{u}}_{0,t}(\xi)\|^2 \leq \frac{\lambda_{\max}(P)}{\lambda_{\min}(P)}\mu^t\|\xi - \xi'\|^2 = \kappa(P)\|\xi - \xi'\|^2\mu^t,$$

and the proof is complete. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \square$

**Theorem 3.2.** *Suppose the system* (3.1) *satisfies Assumption 3.1 and the Demidovich criterion with* $\mathbb{U} = [-R, R]$, *its positively invariant set* $\mathbb{X}$ *contains* $0$, *and* $f(0, 0) = 0$. *Then its i/o map* $\mathsf{F}$ *with zero initial condition* $x_0 = 0$

*satisfies Assumptions 2.1 and 2.2 with*

$$m_{\mathsf{F}}^*(\epsilon) \leq \frac{2\log(\frac{2\kappa(P)L_f L_g R}{(1-\sqrt{\mu})^2 \epsilon})}{\log\frac{1}{\mu}} \qquad and \qquad \omega_{t,\mathsf{F}}(\delta) \leq \frac{\sqrt{\kappa(P)}L_f L_g \delta}{1-\sqrt{\mu}}. \qquad (3.15)$$

*Proof.* Since $P$ is symmetric and positive definite, $\|x\|_P := \sqrt{x^\top P x}$ is a norm on $\mathbb{R}^n$ with $\lambda_{\min}(P)\|\cdot\|^2 \leq \|\cdot\|_P^2 \leq \lambda_{\max}(P)\|\cdot\|^2$. Then, for all $\xi \in \mathbb{X}$, $\mathbf{u} \in \mathcal{M}(R)$, and $t$,

$$\begin{aligned}
\|\varphi_{0,t+1}^{\mathbf{u}}(\xi)\|_P &= \|f(\varphi_{0,t}^{\mathbf{u}}(\xi), u_t)\|_P \\
&\leq \|f(\varphi_{0,t}^{\mathbf{u}}(\xi), u_t) - f(0, u_t)\|_P + \|f(0, u_t) - f(0,0)\|_P \\
&\leq \sqrt{\mu}\|\varphi_{0,t}^{\mathbf{u}}(\xi)\|_P + \sqrt{\lambda_{\max}(P)}L_f R,
\end{aligned}$$

where we have used the Lyapunov bound (3.14). Unrolling the recursion gives the estimate

$$\sup_{t\in\mathbb{Z}_+} \sup_{\mathbf{u}\in\mathcal{M}(R)} \|\varphi_{0,t}^{\mathbf{u}}(\xi)\|_P \leq \sqrt{\mu}\|\xi\|_P + \frac{\sqrt{\lambda_{\max}(P)}L_f R}{1-\sqrt{\mu}}.$$

Thus, Assumption 3.2 is satisfied, where $\mathbb{S}_\xi$ is the ball centered at 0 with $\ell^2$-radius $\sqrt{\kappa(P)}\left(\|\xi\| + \frac{L_f R}{1-\sqrt{\mu}}\right)$. Assumption 3.3 is also satisfied by Proposition 3.2. The estimates in (3.15) follow from Theorem 3.1. $\qquad\square$

The following result now follows as a direct consequence of Theorems 2.1 and 3.2:

**Corollary 3.1.** *If the system (3.1) satisfies the conditions of Theorem 3.2, then its i/o map $\mathsf{F}$ with zero initial condition can be $\epsilon$-approximated in the sense of Theorem 2.1 by a ReLU TCN $\widehat{\mathsf{F}}$ with width $\mathrm{polylog}(\frac{1}{\epsilon})$ and depth $\mathrm{quasipoly}(\frac{1}{\epsilon})$.[2]*

## 3.3 Contractivity vs. the Demidovich criterion

If the contractivity condition (3.10) holds and $f(x,u)$ is differentiable in $x$, then the Demidovich criterion is satisfied with $P = I_n$ and $\mu = \lambda^2$. In that case,

---

[2]We say that a given quantity $N$ has *quasipolynomial growth* in $1/\epsilon$, and we write $N \leq \mathrm{quasipoly}(1/\epsilon)$, if $N = O(\exp(\mathrm{polylog}(\frac{1}{\epsilon})))$.

we immediately obtain the exponential estimate $\beta(C,t) \leq C\lambda^t$. However, the Demidovich criterion covers a wider class of nonlinear systems. As an example, consider a discrete-time nonlinear system of *Lur'e type* (cf. [23, 24, 25] and references therein):

$$x_{t+1} = Ax_t + B\psi(u_t - y_t) \tag{3.16a}$$

$$y_t = Cx_t. \tag{3.16b}$$

Here, the state $x_t$ is $n$-dimensional while the input $u_t$ and the output $y_t$ are scalar, so $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times 1}$, and $C \in \mathbb{R}^{1 \times n}$. The map $\psi : \mathbb{R} \to \mathbb{R}$ is a fixed differentiable nonlinearity. The system in (3.16) has the form (3.1) with $f(x,u) = Ax + B\psi(u - Cx)$ and $g(x) = Cx$, and can be realized as the negative feedback interconnection of the discrete-time linear system

$$x_{t+1} = Ax_t + Bv_t \tag{3.17a}$$

$$y_t = Cx_t \tag{3.17b}$$

and the nonlinear element $\psi$ using the feedback law $v_t = \psi(u_t - y_t)$. We make the following assumptions (see, e.g., [21] for the requisite control-theoretic background):

**Assumption 3.4.** *The nonlinearity $\psi : \mathbb{R} \to \mathbb{R}$ satisfies $\psi(0) = 0$, and there exist real numbers $-\infty < a \leq b < \infty$ such that $a \leq \psi'(\cdot) \leq b$.*

**Assumption 3.5.** *$A$ is a* Schur *matrix, i.e., its spectral radius $\rho(A)$ is strictly smaller than 1; the pair $(A,B)$ is* controllable, *i.e., the $n \times n$ matrix $[B \,|\, AB \,|\, \ldots \,|\, A^{n-1}B]$ has rank $n$; and the pair $(A,C)$ is* observable, *i.e., the $n \times n$ matrix $[C^\top \,|\, A^\top C^\top \,|\, \ldots \,|\, (A^\top)^{n-1}C^\top]$ has rank $n$.*

**Assumption 3.6.** *Let $\mathbb{T} := \{z \in \mathbb{C} : |z| = 1\}$ denote the unit circle in the complex plane. The rational function $G(z) := C(zI_n - A)^{-1}B$ satisfies*

$$\|G\|_{\mathcal{H}_\infty(\mathbb{T})} := \sup_{z \in \mathbb{T}} |G(z)| < \gamma^{-1} \tag{3.18}$$

*for some $\gamma > 0$ such that $r^2 \leq \gamma^2$ for all $a \leq r \leq b$.*

**Remark 3.1.** Assumption 3.4 imposes a *slope condition* on $\psi$ and is standard in the analysis of Lur'e systems [26, 13, 24]. The function $G(z)$ is the *transfer*

*function* of the linear system (3.17). Assumption 3.5 states that the triple $(A, B, C)$ is a *minimal realization* of $G$. The quantity $\|G\|_{\mathcal{H}_\infty(\mathbb{T})}$ appearing in Eq. (3.18) in Assumption 3.6 is the $\mathcal{H}_\infty$-*norm* of $G$ on the unit circle in the complex plane. Assumptions 3.5 and 3.6 are also common and are in the spirit of the well-known *circle criterion* [26, 23].

With these preliminaries out of the way, we have the following:

**Proposition 3.3.** *Suppose that system* (3.16) *satisfies Assumptions 3.4–3.6. Then it satisfies the discrete-time Demidovich criterion with* $\mathbb{X} = \mathbb{R}^n$ *and* $\mathbb{U} = \mathbb{R}$, *and moreover* $\mu > \rho(A)^2$.

*Proof.* Since the matrix $A$ is Schur, the function

$$g(r) := \sup_{z \in \mathbb{T}} |G(rz)| = \|G(r\cdot)\|_{\mathcal{H}_\infty(\mathbb{T})}, \qquad r > \rho(A)$$

is continuous. In particular, there exists some $r_0 \in (\rho(A), 1)$, such that $g(r_0) < g(1) < \gamma^{-1}$. Consequently, the rational function

$$H(z) := \gamma G(r_0 z) = \frac{\gamma C}{r_0} \left( z I_n - \frac{A}{r_0} \right)^{-1} B$$

is well-defined for all $z \in \mathbb{C}$ with $|z| \geq r_0$, and we have the following:

- $\frac{A}{r_0}$ is a Schur matrix;

- the pair $(\frac{A}{r_0}, B)$ is controllable;

- the pair $(\frac{A}{r_0}, \frac{\gamma C}{r_0})$ is observable;

- $\|H\|_{\mathcal{H}_\infty(\mathbb{T})} < 1$.

Then, by the Discrete-Time Bounded-Real Lemma [27], there exist real matrices $L, W$ and a symmetric positive definite matrix $P \in \mathbb{R}^{n \times n}$, such that

$$A^\top P A + \gamma^2 C^\top C + r_0^2 L^\top L = r_0^2 P \tag{3.19a}$$

$$B^\top P B + W^\top W = I_n \tag{3.19b}$$

$$A^\top P B + r_0 L^\top W = r_0 I_n. \tag{3.19c}$$

From (3.19), for any $\theta \in \mathbb{R}$ we have

$$(A - \theta BC)^\top P(A - \theta BC) - r_0^2 P$$
$$= A^\top PA - \theta(C^\top B^\top PA + A^\top PBC) + \theta^2 C^\top B^\top PBC - r_0^2 P$$
$$= (\theta^2 - \gamma^2)C^\top C - (r_0 L - \theta WC)^\top (r_0 L - \theta WC).$$

Let $\mu := r_0^2$. Then, since $\gamma^2 \geq \theta^2$ for all $\theta \in [a, b]$, it follows that

$$(A - \theta BC)^\top P(A - \theta BC) - \mu P \preceq 0, \qquad a \leq \theta \leq b.$$

Since

$$\frac{\partial}{\partial x} f(x, u) = \frac{\partial}{\partial x} \left( Ax + B\psi(u - Cx) \right) = A - \psi'(u - Cx)BC$$

and $\psi'(u - Cx) \in [a, b]$ for all $x$ and $u$, the proposition is proved. $\qquad\square$

The crucial ingredient in the proof is the Discrete-Time Bounded-Real Lemma [27], which guarantees the existence of the matrix $P$ appearing in the Demidovich criterion. The main takeaway here is that the function $f(x, u) = Ax + B\psi(u - Cx)$ need not be contractive (i.e., it may be the case that $P \neq I_n$), but it will be contractive in the $\| \cdot \|_P$ norm.

## 3.4 Continuous-time dynamical systems

Much of what has already been developed for discrete-time state-space models can be adapted to continuous-time state space models. Consider dynamical systems with input of the following form:

$$\begin{aligned} \dot{x} &= f(x, u) & x(t) \in \mathbb{R}^n \quad u(t) \in \mathbb{R}^m \\ y &= h(x) & y(t) \in \mathbb{R}^p \end{aligned} \tag{3.20}$$

where the state transition map $f : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^n$ and the output map $h : \mathbb{R}^n \to \mathbb{R}^p$ are continuously differentiable. By augmenting the state space to include the equations $\dot{y} = \frac{\partial h}{\partial x} f(x, u)$ if necessary, we may assume without loss of generality that $h$ is given by a linear map $x \mapsto Hx$ for some $H \in \mathbb{R}^{p \times n}$.

The set of uniformly bounded inputs $u : \mathbb{R}_+ \to \mathbb{R}^m$ now becomes

$$\mathcal{U} := \{u : \mathbb{R}_+ \to \mathbb{R}^m : \sup_{t \geq 0} |u(t)| \leq R\},$$

where $|\cdot|$ denotes the Euclidean norm. For an input $u \in \mathcal{U}$ and times $0 \leq s \leq t$, we again denote the state $x(t)$ at time $t$ that results from initial condition $x(s) = \xi$ at time $s$ by $\varphi_{s,t}^u(\xi)$, referred to as the *flow* or *trajectory* generated by the system (3.20). Like before, we call a set $\mathcal{X} \subseteq \mathbb{R}^n$ *positively invariant* for inputs in $\mathcal{U}$ if, for all $\xi \in \mathcal{X}$, all $u \in \mathcal{U}$, and all $0 \leq s \leq t$, we have $\varphi_{s,t}^u(\xi) \in \mathcal{X}$. The incremental stability property phrased for discrete-time systems is now stated for continuous-time systems as follows:

**Definition 3.2.** *A dynamical system is* uniformly asymptotically incrementally stable *for inputs in $\mathcal{U}$ on a positively invariant set $\mathcal{X}$ if there exists a function $\beta : \mathbb{R}_+ \times \mathbb{R}_+ \to \mathbb{R}_+$ of class $\mathcal{KL}$[3] such that*

$$|\varphi_{s,t}^u(\xi) - \varphi_{s,t}^u(\xi')| \leq \beta(|\xi - \xi'|, t - s) \tag{3.21}$$

*holds for all $u \in \mathcal{U}$, all $\xi, \xi' \in \mathcal{X}$, and all $0 \leq s \leq t$.*

This property quantitatively captures the idea that perturbations to the initial condition have asymptotically negligible influence on the long-term behavior of the system trajectory. For systems satisfying this definition, imperfect system models may still be capable of generating outputs that uniformly approximate the outputs of the original system over infinite time intervals. We can formulate the necessary assumptions of desired approximation and simulation results as regularity conditions on the function $\beta$ which are now suited for the continuous-time setting. For systems not satisfying this stability condition, a sharp bound on the approximation error degrades exponentially with time [28, 29].

---

[3]A function $\beta : \mathbb{R}_+ \times \mathbb{R}_+ \to \mathbb{R}_+$ is of class $\mathcal{KL}$ if 1) for any $t$, the map $h \mapsto \beta(h, t)$ is continuous and strictly increasing and $\beta(0, t) = 0$ and 2) for any $h$, the map $t \mapsto \beta(h, t)$ is continuous and strictly decreasing and $\lim_{t \to \infty} \beta(h, t) = 0$.

# CHAPTER 4

# UNIVERSAL SIMULATION WITH RECURRENT NEURAL NETS

In many practical applications of dynamical systems modeling, the main criterion for an effective model is that it approximately reproduces both the correct input/output relationships and the internal state dynamics. There are many different ways of expressing this criterion; in this work, we use the following formulation [15]:

Consider two systems $\Sigma$ and $\tilde{\Sigma}$ described by the following dynamics:

$$\Sigma : \quad \begin{aligned} \dot{x} &= f(x, u) \\ y &= Hx \end{aligned}$$

$$\tilde{\Sigma} : \quad \begin{aligned} \dot{\tilde{x}} &= \tilde{f}(\tilde{x}, u) \\ \tilde{y} &= \tilde{H}\tilde{x} \end{aligned}$$

with inputs $u(t) \in \mathbb{R}^m$, outputs $y(t) \in \mathbb{R}^p$, and states $x(t) \in \mathbb{R}^n$ and $\tilde{x}(t) \in \mathbb{R}^{\tilde{n}}$. Suppose we are given a compact set $\mathcal{K} \subset \mathbb{R}^n$, a set $\mathcal{U}$ of admissible inputs, and a time interval $T \subseteq \mathbb{R}_+$. We say that $\tilde{\Sigma}$ *simulates* $\Sigma$ *on sets* $\mathcal{K}$ *and* $\mathcal{U}$ *up to accuracy* $\epsilon$ *for times* $t \in T$ if there exist two continuous maps $\alpha : \mathbb{R}^{\tilde{n}} \to \mathbb{R}^n$ and $\gamma : \mathbb{R}^n \to \mathbb{R}^{\tilde{n}}$ such that, when $\Sigma$ is initialized at $x(s) = \xi \in \mathcal{K}$, $\tilde{\Sigma}$ is initialized at $\tilde{x}(s) = \gamma(\xi)$, where $s := \inf T$, and any common input $u(\cdot) \in \mathcal{U}$ is supplied to both $\Sigma$ and $\tilde{\Sigma}$, we have

$$|x(t) - \alpha(\tilde{x}(t))| < \epsilon \qquad \text{and} \qquad |y(t) - \tilde{y}(t)| < \epsilon$$

for all $t \in T$. We consider the case when the simulating system $\tilde{\Sigma}$ is a (continuous-time) *recurrent neural net*, i.e., $\tilde{f}$ has the form

$$\tilde{f}(\tilde{x}, u) = -\frac{1}{\tau}\tilde{x} + \sigma_{\tilde{n}}(A\tilde{x} + Bu),$$

where $\tau > 0$ is a positive constant, $A \in \mathbb{R}^{\tilde{n} \times \tilde{n}}$ and $B \in \mathbb{R}^{\tilde{n} \times m}$ are time-

invariant matrices, and $\sigma_{\tilde{n}} : \mathbb{R}^{\tilde{n}} \to \mathbb{R}^{\tilde{n}}$ is the diagonal map defined by $\sigma_{\tilde{n}}(\tilde{x}) := [\sigma(\tilde{x}_1) \cdots \sigma(\tilde{x}_{\tilde{n}})]^{\mathsf{T}}$, where $\sigma : \mathbb{R} \to (0, 1)$ is a continuous, strictly increasing function with $\lim_{h \to -\infty} \sigma(h) = 0$ and $\lim_{h \to \infty} \sigma(h) = 1$. Such functions are referred to as *sigmoidal* in the literature on neural nets [30].

## 4.1 Simulating stable systems with recurrent neural nets

Consider system (3.20) with an open positively invariant set $\mathcal{X} \subseteq \mathbb{R}^n$. We impose the following assumptions:

**Assumption 4.1.** *There exists a compact subset $\mathcal{K} \subset \mathcal{X}$ such that, for any initial condition $\xi \in \mathcal{K}$, there exists a compact subset $\mathcal{X}_\xi \subset \mathcal{X}$, such that $\varphi_{s,t}^u(\xi) \in \mathcal{X}_\xi$ for all $u \in \mathcal{U}$ and all $t \geq s \geq 0$.*

**Assumption 4.2.** *System (3.20) is uniformly asymptotically incrementally stable on $\mathcal{X}$ for inputs in $\mathcal{U}$, and the function $\beta$ in equation (3.21) satisfies the following conditions:*

1. *For any $t \geq 0$, the map $h \mapsto \beta(h, t)$ is differentiable from the right at $h = 0$.*

2. $\displaystyle \int_0^\infty \frac{\partial}{\partial h} \beta(h, t) \Big|_{h=0^+} \mathrm{d}t =: b < \infty.$

Assumption 4.2 is evidently satisfied by exponentially stable systems with $\beta(h, t) = c h e^{-\kappa t}$ for some $c, \kappa > 0$, but it also holds for systems with much longer transients, e.g., when $\beta(h, t) = \frac{ch}{(t+1)^{1+\kappa}}$.

**Theorem 4.1.** *Consider system (3.20) and suppose that Assumptions 4.1 and 4.2 are satisfied. Then, for any $\epsilon > 0$, there exists a recurrent neural net of the form*

$$\dot{\tilde{x}} = -\frac{1}{\tau}\tilde{x} + \sigma_{\tilde{n}}(\tilde{A}\tilde{x} + \tilde{B}u)$$
$$\tilde{y} = \tilde{H}\tilde{x}$$

*for some $\tau > 0$, $\tilde{A} \in \mathbb{R}^{\tilde{n} \times \tilde{n}}$, $\tilde{B} \in \mathbb{R}^{\tilde{n} \times m}$, and $\tilde{H} \in \mathbb{R}^{p \times \tilde{n}}$ that simulates system (3.20) on sets $\mathcal{K}$ and $\mathcal{U}$ up to accuracy $\epsilon$ for all $t \in \mathbb{R}_+$. Moreover, the mappings $\alpha : \mathbb{R}^{\tilde{n}} \to \mathbb{R}^n$ and $\gamma : \mathbb{R}^n \to \mathbb{R}^{\tilde{n}}$ that implement the approximate simulation are linear.*

## 4.2 Technical lemmas

To prove the theorem, we will make use of the following lemmas.

**Lemma 4.1.** *Let $D\varphi_{s,t}^u(\xi) \cdot v$ denote the directional derivative of $\varphi_{s,t}^u(\xi)$ with respect to $\xi$ in the direction of $v$. Suppose that Assumptions 4.1 and 4.2 are satisfied. Then for any $\xi \in \mathcal{X}$, the induced norm*

$$\|D\varphi_{s,t}^u(\xi)\| := \sup_{|v|=1} |D\varphi_{s,t}^u(\xi) \cdot v|$$

*is integrable with respect to $t$ on $[s, \infty)$.*

*Proof.* From definitions,

$$
\begin{aligned}
\|D\varphi_{s,t}^u(\xi)\| &= \sup_{|v|=1} |D\varphi_{s,t}^u(\xi) \cdot v| \\
&= \sup_{|v|=1} \lim_{h\downarrow 0} \frac{1}{|hv|} |\varphi_{s,t}^u(\xi + hv) - \varphi_{s,t}^u(\xi)| \\
&\leq \sup_{|v|=1} \lim_{h\downarrow 0} \frac{1}{|hv|} \beta(|hv|, t-s) \\
&= \lim_{h\downarrow 0} \frac{1}{h} \beta(h, t-s) \\
&= \frac{\partial}{\partial h} \beta(h, t-s)\Big|_{h=0^+}
\end{aligned}
$$

and, by Assumption 4.2, $\frac{\partial}{\partial h}\beta(h, t-s)|_{h=0^+}$ is integrable with respect to $t$ on $[s, \infty)$. $\qquad\square$

**Lemma 4.2.** *Consider two dynamical systems $\dot{x} = f(x, u)$ and $\dot{\hat{x}} = \hat{f}(\hat{x}, u)$ with $x(t), \hat{x}(t) \in \mathbb{R}^n$, which generate flows $\varphi_{s,t}^u(\xi)$ and $\hat{\varphi}_{s,t}^u(\xi)$, respectively. Then the following inequality holds for all $t \geq s \geq 0$:*

$$|\varphi_{s,t}^u(\xi) - \hat{\varphi}_{s,t}^u(\xi)| \leq \int_s^t \|D\varphi_{r,t}^u(\hat{\varphi}_{s,r}^u(\xi))\| \cdot |f(\hat{\varphi}_{s,r}^u(\xi), u(r)) - \hat{f}(\hat{\varphi}_{s,r}^u(\xi), u(r))| \, \mathrm{d}r.$$
(4.1)

The proof can be found in [31], Chapter 3, Proposition 3.1.3.

**Lemma 4.3.** *Consider the $C^1$ map $f : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^m$ from system (3.20) satisfying Assumption 4.1. Then, for any $\epsilon > 0$, we can construct:*

- *compact sets $\mathcal{X}_0 \subset \mathcal{X}_1 \subset \mathcal{X}_2 \subset \mathcal{X}$;*

- *two $C^\infty$ bump functions $\rho_0, \rho_1 : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$ that satisfy $\rho_0|_{\mathcal{Z}_0} \equiv 1$, $\rho_0|_{(\mathcal{Z}_1)^c} \equiv 0$, $\rho_1|_{\mathcal{Z}_1} \equiv 1$, $\rho_1|_{(\mathcal{Z}_2)^c} \equiv 0$ for $\mathcal{Z}_i := \mathcal{X}_i \times B_R^m(0)$, $i \in \{1, 2, 3\}$, where $B_R^m(0) := \{v \in \mathbb{R}^m : |v| \le R\}$;*

- *a $C^1$ map $\hat{f}$ that vanishes outside $\mathcal{Z}_2$, such that, for $(x, u) \in \mathcal{Z}_1$,*

$$\hat{f}(x, u) = -\frac{1}{\tau}x + T\sigma_\ell(Ax + Bu + \mu) + \nu$$

*for some $\tau > 0$, $T \in \mathbb{R}^{n \times \ell}$, $A \in \mathbb{R}^{\ell \times n}$, $B \in \mathbb{R}^{\ell \times m}$, $\mu \in \mathbb{R}^\ell$, and $\nu \in \mathbb{R}^n$, and*

$$\sup_{(x,u)\in\mathbb{R}^n\times\mathbb{R}^m} |\rho(x, u)f(x, u) - \hat{f}(x, u)| \le \epsilon.$$

*Proof.* By Assumption 4.1, we have a compact subset $\mathcal{K} \subset \mathcal{X}$, and for each $\xi \in \mathcal{K}$ we have a compact subset $\mathcal{X}_\xi \subset \mathcal{X}$ with $\xi \in \mathcal{X}_\xi$. Since $\mathcal{K}$ is compact, the set $\cup_{\xi \in \mathcal{K}} \mathcal{X}_\xi$ is bounded. Therefore, the set

$$\mathcal{X}_\mathcal{K}^{(0)} := \mathrm{cl}\,\big( \bigcup_{\xi \in \mathcal{K}} \mathcal{X}_\xi \big),$$

is closed and bounded, hence compact. Now fix any $\epsilon > 0$ and let $\eta < \min(\epsilon, \lambda)$ where

$$\lambda := \mathrm{dist}(\mathcal{X}_\mathcal{K}^{(0)}, \partial\mathcal{X}) := \inf\{|x - y| : x \in \mathcal{X}_\mathcal{K}^{(0)},\ y \in \partial\mathcal{X}\}.$$

Let $B_r^k(0) \subset \mathbb{R}^k$ denote the closed Euclidean ball of radius $r$ centered at the origin and define the following sets:

$$\mathcal{X}_\mathcal{K}^{(1)} := \mathcal{X}_\mathcal{K}^{(0)} + B_{\frac{\eta}{2}}^n(0)$$
$$\mathcal{X}_\mathcal{K}^{(2)} := \mathcal{X}_\mathcal{K}^{(0)} + B_\eta^n(0)$$
$$\mathcal{Z}_\mathcal{K}^{(i)} := \mathcal{X}_\mathcal{K}^{(i)} \times B_R^m(0) \quad \text{for } i = 0, 1, 2$$

where $+$ denotes Minkowski addition and where $R$ is the uniform bound on inputs $u \in \mathcal{U}$. Then we have

$$\mathcal{X}_\mathcal{K}^{(0)} \subset \mathcal{X}_\mathcal{K}^{(1)} \subset \mathcal{X}_\mathcal{K}^{(2)} \subset \mathcal{X}$$
$$\mathcal{Z}_\mathcal{K}^{(0)} \subset \mathcal{Z}_\mathcal{K}^{(1)} \subset \mathcal{Z}_\mathcal{K}^{(2)} \subset \mathcal{X} \times B_R^m(0)$$
$$\mathcal{X}_\mathcal{K}^{(i)},\ \mathcal{Z}_\mathcal{K}^{(i)} \text{ are compact for } i = 0, 1, 2$$

Construct a $C^\infty$ bump function $\rho_0 : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$, so that $\rho_0|_{\mathcal{Z}_\mathcal{K}^{(0)}} \equiv 1$ and $\rho_0|_{\mathbb{R}^n \times \mathbb{R}^m \setminus \mathcal{Z}_\mathcal{K}^{(1)}} \equiv 0$. Then we have $\rho_0 f|_{\mathcal{Z}_\mathcal{K}^{(0)}} \equiv f|_{\mathcal{Z}_\mathcal{K}^{(0)}}$. By the universal approximation theorem [32], there exists a feedforward neural net

$$g(x, u) = T\sigma_\ell(Ax + Bu + \mu) + \nu,$$

where $T \in \mathbb{R}^{n \times \ell}$, $A \in \mathbb{R}^{\ell \times n}$, $B \in \mathbb{R}^{\ell \times m}$, $\mu \in \mathbb{R}^\ell$, $\nu \in \mathbb{R}^n$, such that

$$\sup_{(x,u) \in \mathcal{Z}_\mathcal{K}^{(2)}} |(\rho_0 f)(x, u) - g(x, u)| \leq \frac{\epsilon}{2}.$$

Choose $\tau > 0$ sufficiently large that $|x| \leq \frac{\tau\epsilon}{2}$ for all $x \in \mathcal{X}_\mathcal{K}^{(2)}$. Construct a second $C^\infty$ bump function $\rho_1 : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$, so that $\rho_1|_{\mathcal{Z}_\mathcal{K}^{(1)}} \equiv 1$ and $\rho_1|_{\mathbb{R}^n \times \mathbb{R}^m \setminus \mathcal{Z}_\mathcal{K}^{(2)}} \equiv 0$. Now we can construct a $C^1$ map $\hat{f} : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^n$ such that the following hold:

1. $\hat{f}|_{\mathcal{Z}_\mathcal{K}^{(1)}}(x, u) = -\frac{1}{\tau}x + g(x, u)$.

2. $\displaystyle\sup_{(x,u) \in \mathcal{Z}_\mathcal{K}^{(2)} \setminus \mathcal{Z}_\mathcal{K}^{(1)}} |\hat{f}(x, u)| < \epsilon$.

3. $\hat{f}|_{\mathbb{R}^n \times \mathbb{R}^m \setminus \mathcal{Z}_\mathcal{K}^{(2)}} \equiv 0$.

Evidently, $\hat{f}$ is given by multiplying the map $(x, u) \mapsto -\frac{1}{\tau}x + g(x, u)$ by the bump function $\rho_1$. Then we have

1. For all $(x, u) \in \mathcal{Z}_\mathcal{K}^{(0)}$, $(\rho_0 f)(x, u) = f(x, u)$ and $|f(x, u) - \hat{f}(x, u)| \leq \epsilon$.

2. For all $(x, u) \in \mathcal{Z}_\mathcal{K}^{(1)} \setminus \mathcal{Z}_\mathcal{K}^{(0)}$, $|(\rho_0 f)(x, u) - \hat{f}(x, u)| \leq \epsilon$.

3. For all $(x, u) \in \mathcal{Z}_\mathcal{K}^{(2)} \setminus \mathcal{Z}_\mathcal{K}^{(1)}$, $(\rho_0 f)(x, u) = 0$ and $|0 - \hat{f}(x, u)| \leq \epsilon$.

4. For all $(x, u) \in \mathbb{R}^n \times \mathbb{R}^m \setminus \mathcal{Z}_\mathcal{K}^{(2)}$, $(\rho_0 f)(x, u) = 0$ and $\hat{f}(x, u) = 0$.

Therefore $\|\rho_0 f - \hat{f}\|_\infty \leq \epsilon$ and $\hat{f}(x, u) = 0$ for $(x, u)$ outside the compact set $\mathcal{Z}_\mathcal{K}^{(2)}$, so $\mathcal{X}_0 = \mathcal{X}_\mathcal{K}^{(0)}$, $\mathcal{X}_1 = \mathcal{X}_\mathcal{K}^{(1)}$, $\mathcal{X}_2 = \mathcal{X}_\mathcal{K}^{(2)}$, $\rho_0$, $\rho_1$, and $\hat{f}$ are the objects we wished to construct. $\qquad\square$

**Lemma 4.4.** *The state-space dynamics*

$$\dot{x} = -\frac{1}{\tau}\hat{x} + T\sigma_\ell(A\hat{x} + Bu + \mu) + \nu \tag{4.2}$$

can be simulated with zero loss in accuracy by a system in the form of a recurrent net

$$\dot{\tilde{x}} = -\frac{1}{\tau}\tilde{x} + \sigma_{\tilde{n}}(\tilde{A}\tilde{x} + \tilde{B}u) \tag{4.3}$$

for some $\tilde{n} \in \mathbb{N}$, $\tilde{A} \in \mathbb{R}^{\tilde{n} \times \tilde{n}}$, $\tilde{B} \in \mathbb{R}^{\tilde{n} \times m}$. That is, there exist matrices $F \in \mathbb{R}^{n \times \tilde{n}}$ and $G \in \mathbb{R}^{\tilde{n} \times n}$, such that $\hat{x}(t) = F\tilde{x}(t)$ for all $t \geq 0$, with initial conditions $\hat{x}(0) = \xi$ and $\tilde{x}(0) = G\xi$.

*Proof.* Following [15], we will construct the recurrent net (4.3) and the matrices $F$ and $G$ in several steps.

**Step 1 - Eliminating $T$:** We may assume without loss of generality that the matrix $T$ takes the form $[T_1^{\mathsf{T}} \ 0]^{\mathsf{T}}$ with $T_1$ having full row rank. Then the system (4.2) can be written as

$$\dot{x}_1 = -\frac{1}{\tau}x_1 + T_1\sigma_\ell(A_1x_1 + A_2x_2 + Bu + \mu) + \nu_1, \qquad x_1(0) = \xi_1 \tag{4.4}$$
$$\dot{x}_2 = -\frac{1}{\tau}x_2 + \nu_2, \qquad\qquad\qquad\qquad\qquad x_2(0) = \xi_2$$

where $\xi = [\xi_1^{\mathsf{T}} \ \xi_2^{\mathsf{T}}]^{\mathsf{T}}$ and $\nu = [\nu_1^{\mathsf{T}} \ \nu_2^{\mathsf{T}}]^{\mathsf{T}}$. Since $T_1$ is surjective, there exist vectors $\tilde{\nu}_1, \tilde{\xi}_1 \in \mathbb{R}^\ell$ such that $T_1\tilde{\nu}_1 = \nu_1$ and $T_1\tilde{\xi}_1 = \xi_1$. Consider the following transformed system:

$$\dot{z}_1 = -\frac{1}{\tau}z_1 + \sigma_\ell(A_1T_1z_1 + A_2x_2 + Bu + \mu) + \tilde{\nu}_1, \qquad z_1(0) = \tilde{\xi}_1 \tag{4.5}$$
$$\dot{x}_2 = -\frac{1}{\tau}x_2 + \nu_2. \qquad\qquad\qquad\qquad\qquad x_2(0) = \xi_2$$

The trajectory $(x_1(t), x_2(t))$ of system (4.4) can be recovered from the trajectory $(z_1(t), x_2(t))$ of system (4.5) via the transformation $x_1(t) := T_1z_1(t)$. Let $\kappa := \sigma(0)$; then the equation for the dynamics of $x_2$ may be rewritten as

$$\dot{x}_2 = -\frac{1}{\tau}x_2 + \sigma_{n-r}(0x + 0u) + (\nu_2 - [\kappa \ \cdots \ \kappa]^{\mathsf{T}}),$$

where $r := \text{rank}(T_1)$. This permits us to combine the two equations in (4.5) into

$$\dot{\bar{x}} = -\frac{1}{\tau}\bar{x} + \sigma_{\bar{n}}(\bar{A}\bar{x} + \bar{B}u + \bar{\mu}) + \bar{\nu}$$

for suitable matrices $\bar{A} \in \mathbb{R}^{\bar{n} \times \bar{n}}, \bar{B} \in \mathbb{R}^{\bar{n} \times m}$, vectors $\bar{\mu}, \bar{\nu} \in \mathbb{R}^{\bar{n}}$, and the initial condition $\bar{x}(0) = \bar{\xi}$, where $\bar{n} := \ell + n - r$ and $\bar{\xi} := [\tilde{\xi}_1^{\mathsf{T}} \ \xi_2^{\mathsf{T}}]^{\mathsf{T}}$.

28

**Step 2 - Eliminating $\bar{\nu}$:** Define $\underline{x} := \bar{x} - \tau\bar{\nu}$ and $\theta := \tau\bar{A}\bar{\nu} + \bar{\mu}$. It follows that

$$\dot{\underline{x}} = -\frac{1}{\tau}\underline{x} + \sigma_{\bar{n}}(\bar{A}\underline{x} + \bar{B}u + \theta)$$

with $\underline{x}(0) = \underline{\xi} := \bar{\xi} - \tau\bar{\nu}$, and the trajectory $\bar{x}(t)$ from Step 1 is recovered via $\bar{x}(t) = \underline{x}(t) + \tau\bar{\nu}$.

**Step 3 - Eliminating $\theta$:** Since $\sigma : \mathbb{R} \to (0,1)$ is bounded, positive, and continuous, the fixed-point equation $z = \tau\sigma(z)$ has at least one nonzero solution $\zeta$, by Brouwer's fixed-point theorem. Consider the following dynamics for $\tilde{x}(t) \in \mathbb{R}^{\bar{n}+1}$:

$$\dot{\tilde{x}}_{1:\bar{n}} = -\frac{1}{\tau}\tilde{x}_{1:\bar{n}} + \sigma_{\bar{n}}(\bar{A}\tilde{x}_{1:\bar{n}} + \frac{1}{\zeta}\theta\tilde{x}_{\bar{n}+1} + \bar{B}u), \qquad \tilde{x}_{1:\bar{n}}(0) = \underline{\xi}$$

$$\dot{\tilde{x}}_{\bar{n}+1} = -\frac{1}{\tau}\tilde{x}_{\bar{n}+1} + \sigma(\tilde{x}_{\bar{n}+1}), \qquad\qquad\qquad \tilde{x}_{\bar{n}+1}(0) = \zeta$$

where evidently $\underline{x}(t) = \tilde{x}_{1:\bar{n}}(t)$ and $\tilde{x}_{\bar{n}+1}(t) \equiv \zeta$ for all $t$. With $\tilde{n} := \bar{n} + 1$, this system can be represented in the desired form $\dot{\tilde{x}} = -\frac{1}{\tau}\tilde{x} + \sigma_{\tilde{n}}(\tilde{A}\tilde{x} + \tilde{B}u)$ by choosing

$$\tilde{A} := \begin{bmatrix} \bar{A} & \frac{1}{\zeta}\theta \\ 0 & 1 \end{bmatrix}, \quad \tilde{B} := \begin{bmatrix} \bar{B} \\ 0 \end{bmatrix}.$$

Altogether, we have shown that the trajectory of the system (4.2) with $x(0) = \xi$ can be reproduced with zero loss in accuracy by a recurrent net (4.3) by expanding the dimension of the state space from $n$ to $n + \ell - r$ and adding one more neuron, for a total of $\tilde{n} = \ell + n - r + 1$ neurons. The matrices $F$ and $G$ can be constructed by retracing the above steps backwards from $\tilde{x}$ to $\hat{x}$ and then forwards from $\xi$ to $\tilde{\xi} := [\underline{\xi}^\mathsf{T} \ \zeta]^\mathsf{T}$. (The affine map $\bar{x}(t) = \underline{x}(t) + \tau\bar{\nu}$ can be implemented as a linear map $\bar{x}(t) = \tilde{x}_{1:\bar{n}}(t) + \frac{\tau}{\zeta}\bar{\nu}\tilde{x}_{\bar{n}+1}$, since $\tilde{x}_{\bar{n}+1}(t) \equiv \zeta$ for all $t$.) $\qquad\square$

## 4.3 Proof of the universal simulation theorem for stable dynamical systems

Fix some $\eta > 0$ to be chosen later. Consider system (3.20) with an open positively invariant set $\mathcal{X} \subseteq \mathbb{R}^n$ satisfying Assumptions 4.1 and 4.2. By

Lemma 4.3 (with $\epsilon \leftarrow \frac{\eta}{b}$), there exist a map $\hat{f} : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^n$ and compact sets $\mathcal{Z}_0 \subset \mathcal{Z}_1 \subset \mathcal{Z}_2 \subset \mathcal{X} \times \mathbb{R}^m$, such that $\hat{f}|_{\mathcal{Z}_1}(x, u) = -\frac{1}{\tau}x + T\sigma_\ell(Ax + Bu + \mu) + \nu$ and $\|\rho_0 f - \hat{f}\|_\infty \leq \frac{\eta}{b}$, where $\rho_0$ is a $C^\infty$ bump function satisfying $\rho_0|_{\mathcal{Z}_0} = 1$ and $\rho_0|_{(\mathcal{Z}_1)^c} = 0$.

Moreover, by Lemma 4.3, $\mathcal{Z}_i = \mathcal{X}_i \times B_R^m(0)$. Let $\hat{\varphi}_{s,t}^u(\xi)$ denote the flow generated by the system $\dot{\hat{x}} = \hat{f}(\hat{x}, u)$. Since $\hat{f}$ is a $C^1$ map that vanishes outside $\mathcal{Z}_2$, we clearly have $\hat{\varphi}_{s,t}^u(\xi) \in \mathcal{X}_2 \subset \mathcal{X}$ for all $\xi \in \mathcal{K}$, all $u \in \mathcal{U}$, and all $t \geq s \geq 0$, since if the trajectory reaches the boundary $\partial \mathcal{X}_2$, it must stop and remain there permanently because $\hat{f}|_{\partial \mathcal{X}_2 \times B_R^m(0)} \equiv 0$. Furthermore, since $\varphi_{s,t}^u(\xi) \in \mathcal{X}_\mathcal{K} = \mathrm{cl}(\cup_{\xi \in \mathcal{K}} \mathcal{X}_\xi)$ for all $t \geq s \geq 0$ by Assumption 4.1, the flow generated by the system $\dot{x} = (\rho_0 f)(x, u)$ is identically equal to $\varphi_{s,t}^u(\xi)$ because $\rho_0 f|_{\mathcal{X}_\mathcal{K} \times B_R^m(0)} \equiv f|_{\mathcal{X}_\mathcal{K} \times B_R^m(0)}$. Therefore by applying Lemmas 4.1 and 4.2, we have

$$
\begin{aligned}
&|\varphi_{s,t}^u(\xi) - \hat{\varphi}_{s,t}^u(\xi)| \\
&\leq \int_s^t \|D\varphi_{r,t}^u(\hat{\varphi}_{s,r}^u(\xi))\| \cdot |(\rho_0 f)(\hat{\varphi}_{s,r}^u(\xi), u(r)) - \hat{f}(\hat{\varphi}_{s,r}^u(\xi), u(r))| \, \mathrm{d}r \\
&\leq \int_s^t \frac{\partial}{\partial h}\beta(h, r - s)\Big|_{h=0^+} \sup_{(x,u) \in \mathcal{Z}_2} |(\rho_0 f)(x, u) - \hat{f}(x, u)| \, \mathrm{d}r \\
&\leq b \cdot \frac{\eta}{b} = \eta.
\end{aligned}
$$

By Lemma 4.4, the system

$$
\dot{\hat{x}} = -\frac{1}{\tau}\hat{x} + T\sigma_\ell(A\hat{x} + Bu + \mu) + \nu
$$

can be simulated with zero loss in accuracy by a system in the form of a recurrent net

$$
\dot{\tilde{x}} = -\frac{1}{\tau}\tilde{x} + \sigma_{\tilde{n}}(\tilde{A}\tilde{x} + \tilde{B}u).
$$

For the above recurrent net, let $\tilde{y}(t) := \tilde{H}\tilde{x}(t)$ with $\tilde{H} := HF$, where $H \in \mathbb{R}^{p \times n}$ is the linear output map of the original system (3.20) and $F \in \mathbb{R}^{n \times \tilde{n}}$ is the linear map given by Lemma 4.4. Then $H\hat{x}(t) = HF\tilde{x}(t) = \tilde{H}\tilde{x}(t)$ for all

$t \geq 0$, and consequently

$$
\begin{aligned}
|y(t) - \tilde{y}(t)| &= |Hx(t) - \tilde{H}\tilde{x}(t)| \\
&= |Hx(t) - H\hat{x}(t)| \\
&\leq \|H\| |x(t) - \hat{x}(t)| \\
&\leq \|H\| \eta.
\end{aligned}
$$

Choosing $\eta < \min(\epsilon, \frac{\epsilon}{\|H\|})$ gives $|x(t) - F\tilde{x}(t)| < \epsilon$ and $|y(t) - \tilde{y}(t)| < \epsilon$ for all $t \geq 0$, with $x(0) = \xi$ and $\tilde{x}(0) = G\xi$, which completes the proof.

## 4.4 Quantitative approximation bounds for Barron-class systems

Utilizing quantitative approximation bounds developed for feedforward nets, we can develop similar results for recurrent nets. For these bounds to hold, it is necessary for the vector field $f(x, u)$ of the original system (3.20) to satisfy certain regularity conditions [30]:

**Definition 4.1.** *We say that a continuous function $f : \mathbb{R}^d \to \mathbb{R}$ belongs to the* Barron class *if*

$$
C_f := \int_{\mathbb{R}^d} |\omega| |\tilde{f}(\omega)| \, d\omega < \infty,
$$

*where $\tilde{f} : \mathbb{R}^d \to \mathbb{R}$ is the Fourier transform of $f$.*

**Proposition 4.1.** *Let a continuous function $f : \mathbb{R}^d \to \mathbb{R}$ be given, with $C_f < \infty$. Then for every $r > 0$ and every $N \in \mathbb{N}$, there exists a feedforward neural net $g : \mathbb{R}^d \to \mathbb{R}$ of the form*

$$
g(z) = \sum_{k=1}^{N} c_k \sigma(a_k \cdot z + b_k) + c_0,
$$

*such that*

$$
\sup_{z \in B_r^d(0)} |f(z) - g(z)| \leq \frac{2rC_f}{\sqrt{N}}.
$$

The proof can be found in [30] or in [33]. Note that the constant $C_f$ depends implicitly on the input-space dimension $d$. If each coordinate of the state transition map $f : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^n$ from system (3.20) belongs to the Barron

class, then we can bound the number of computation units (neurons) in a recurrent net that simulates system (3.20).

**Proposition 4.2.** *The number of computation units sufficient to guarantee the result of Theorem 4.1 with accuracy $\epsilon$ is*

$$\tilde{n} \geq n + 1 + \frac{16(C_f b \|H\| \Delta)^2 n}{\epsilon^2},$$

*where $C_f$ and $b$ are defined earlier, and $\Delta := \sup_{x \in \mathcal{X}_\mathcal{K}} |x| + R + \frac{\epsilon}{2\|H\|}$.*

**Remark 4.1.** The constant $C_f$ may implicitly depend on the total dimension $n + m$.

*Proof.* The desired underlying feedforward net is constructed in the proof of Lemma 4.3, such that

$$\sup_{(x,u) \in \mathcal{Z}_1} |(\rho_0 f)(x, u) - g(x, u)| \leq \frac{\eta}{2b},$$

where $\mathcal{Z}_1$ is a compact subset of $\mathbb{R}^n \times \mathbb{R}^m$ contained in the ball of radius $\Delta$. On the other hand, Proposition 4.1 gives

$$\sup_{(x,u) \in \mathcal{Z}_1} |(\rho f)(x, u) - g(x, u)| \leq \frac{2C_f \Delta \sqrt{n}}{\sqrt{\ell}},$$

where $\ell$ is the number of neurons in $g$. To achieve the desired inequality, it suffices to take $\ell \geq \frac{16 C_f^2 b^2 \Delta^2 n}{\eta^2}$. From the proof of Theorem 4.1, we set $\eta < \frac{\epsilon}{\|H\|}$, and from Lemma 4.4 we know that $\tilde{n} \geq n + \ell + 1$ neurons suffice. $\square$

# REFERENCES

[1] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin, "Convolutional sequence to sequence learning," in *International Conference on Machine Learning*, 2017.

[2] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier, "Language modeling with gated convolutional networks," in *International Conference on Machine Learning*, 2017.

[3] N. Kalchbrenner, L. Espeholt, K. Simonyan, A. van den Oord, A. Graves, and K. Kavukcuoglu, "Neural machine translation in linear time," 2016. [Online]. Available: https://arxiv.org/abs/1610.10099

[4] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," 2016. [Online]. Available: https://arxiv.org/abs/1609.03499

[5] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, J. Klingner, A. Shah, M. Johnson, X. Liu, Ł. Kaiser, S. Gouws, Y. Kato, T. Kudo, H. Kazawa, K. Stevens, G. Kurian, N. Patil, W. Wang, C. Young, J. Smith, J. Riesa, A. Rudnick, O. Vinyals, G. Corrado, M. Hughes, and J. Dean, "Google's neural machine translation system: Bridging the gap between human and machine translation," 2016. [Online]. Available: https://arxiv.org/abs/1609.08144

[6] R. Johnson and T. Zhang, "Deep pyramid convolutional neural networks for text categorization," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Vancouver, Canada: Association for Computational Linguistics, July 2017. [Online]. Available: https://www.aclweb.org/anthology/P17-1052 pp. 562–570.

[7] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," 2018. [Online]. Available: https://arxiv.org/abs/1803.01271

[8] C. Chelba, M. Norouzi, and S. Bengio, "N-gram language modeling using recurrent neural network estimation," 2017. [Online]. Available: https://arxiv.org/abs/1703.10724

[9] V. Sharan, S. Kakade, P. Liang, and G. Valiant, "Prediction with a short memory," in *Symposium on Theory of Computing*, 2018.

[10] J. Miller and M. Hardt, "Stable recurrent models," in *International Conference on Learning Representations*, 2019.

[11] W. Yin, K. Kann, M. Yu, and H. Schütze, "Comparative study of CNN and RNN for natural language processing," 2017. [Online]. Available: https://arxiv.org/abs/1702.01923

[12] S. Boyd and L. O. Chua, "Fading memory and the problem of approximating nonlinear operators with Volterra series," *IEEE Transactions on Circuits and Systems*, vol. CAS-32, no. 11, pp. 1150–1161, 1985.

[13] I. W. Sandberg, "Structure theorems for nonlinear systems," *Multidimensional Systems and Signal Processing*, vol. 2, pp. 267–286, 1991.

[14] D. N. Tran, B. S. Rüffler, and C. M. Kellett, "Convergence properties for discrete-time nonlinear systems," *IEEE Transactions on Automatic Control*, 2017.

[15] E. D. Sontag, "Neural nets as systems models and controllers," in *Workshop on Adaptive and Learning Systems*, 1992.

[16] K. Funahashi and Y. Nakamura, "Approximation of dynamical systems by continuous time recurrent neural networks," *Neural Networks*, vol. 6, no. 6, pp. 801 – 806, 1993.

[17] T. W. S. Chow and Xiao-Dong Li, "Modeling of continuous time dynamical systems with input by recurrent neural networks," *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, vol. 47, no. 4, pp. 575–578, April 2000.

[18] Xiao-Dong Li, J. K. L. Ho, and T. W. S. Chow, "Approximation of dynamical time-variant systems by continuous-time recurrent neural networks," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 52, no. 10, pp. 656–660, Oct 2005.

[19] B. Hanin and M. Sellke, "Approximating continuous functions by ReLU nets of minimal width," 2018. [Online]. Available: http://arxiv.org/abs/1710.11278

[20] J. Park and I. W. Sandberg, "Criteria for the approximation of nonlinear systems," *IEEE Transactions on Circuits and Systems — I: Fundamental Theory and Applications*, vol. 39, no. 8, pp. 673–676, 1992.

[21] E. D. Sontag, *Mathematical Control Theory: Deterministic Finite Dimensional Systems.* Springer-Verlag, 1998.

[22] A. Pavlov, N. van de Wouw, and H. Nijmeijer, *Uniform Output Regulation of Nonlinear Systems: A Convergent Dynamics Approach.* Birkhäuser, 2006.

[23] I. W. Sandberg and L. Y. Xu, "Steady-state errors in discrete-time control systems," *Automatica*, vol. 29, no. 2, pp. 523–526, 1993.

[24] K. K. Kim and R. D. Braatz, "Observer-based output feedback control of discrete-time Lur'e systems with sector-bounded slope-restricted nonlinearities," *International Journal of Robust and Nonlinear Control*, vol. 24, pp. 2458–2472, 2014.

[25] E. Sarkans and H. Logemann, "Input-to-state stability of discrete-time Lur'e systems," *SIAM Journal on Control and Optimization*, vol. 54, no. 3, pp. 1739–1768, 2016.

[26] Y. Z. Tsypkin, "A criterion of absolute stability for sampled-data systems with monotone characteristics of the nonlinear element," *Doklady Akademii Nauk SSSR*, vol. 155, no. 5, pp. 1029–1032, 1964, in Russian.

[27] P. P. Vaidyanathan, "The discrete-time bounded-real lemma in digital filtering," *IEEE Transactions on Circuits and Systems*, vol. CAS-32, no. 9, pp. 918–924, September 1985.

[28] M. W. Hirsch and S. Smale, "Nonautonomous equations and differentiability of flows," in *Differential Equations, Dynamical Systems, and Linear Algebra.* Academic Press, 1974, ch. 15, pp. 296–303.

[29] E. D. Sontag, *Mathematical Control Theory: Deterministic Finite Dimensional Systems*, 2nd ed. Springer, 1998.

[30] A. R. Barron, "Universal approximation bounds for superpositions of a sigmoidal function," *IEEE Transactions on Information Theory*, vol. 39, no. 3, pp. 930–945, May 1993.

[31] R. van Handel, "Filtering, stability, and robustness," Ph.D. dissertation, California Institute of Technology, 2007.

[32] A. Pinkus, "Approximation theory of the MLP model in neural networks," *Acta Numerica*, pp. 143–195, 1999.

[33] J. E. Yukich, M. B. Stinchcombe, and H. White, "Sup-norm approximation bounds for networks through probabilistic methods," *IEEE Transactions on Information Theory*, vol. 41, no. 4, pp. 1021–1027, July 1995.