# A Quantifier-Reversal Approximation Paradigm for Recurrent Neural Networks

Clemens Hutter[*†]        Valentin Abadie[†]        Helmut Bölcskei[†]

November 20, 2025

## Abstract

Classical neural network approximation results take the form: for every function $f$ and every error tolerance $\epsilon > 0$, one constructs a neural network whose architecture and weights depend on $\epsilon$. This paper introduces a fundamentally different approximation paradigm that reverses this quantifier order. For each target function $f$, we construct a *single* recurrent neural network (RNN) with fixed topology and fixed weights that approximates $f$ to within any prescribed tolerance $\epsilon > 0$ when run for sufficiently many time steps.

The key mechanism enabling this quantifier reversal is temporal computation combined with weight sharing: rather than increasing network depth, the approximation error is reduced solely by running the RNN longer. This yields exponentially decaying approximation error as a function of runtime while requiring storage of only a small, fixed set of weights. Such architectures are appealing for hardware implementations where memory is scarce and runtime is comparatively inexpensive.

To initiate the systematic development of this novel approximation paradigm, we focus on univariate polynomials. Our RNN constructions emulate the structural calculus underlying deep feedforward ReLU network approximation theory—parallelization, linear combinations, affine transformations, and, most importantly, a clocked mechanism that realizes function composition within a single recurrent architecture. The resulting RNNs have size independent of the error tolerance $\epsilon$ and hidden-state dimension linear in the degree of the polynomial.

## 1   Introduction

The starting point of this work is the classical universal approximation theorem for feedforward neural networks [1], [2], [3], which asserts that every continuous function on a compact domain can be approximated arbitrarily well by a shallow neural network with sigmoidal activation function. Subsequent quantitative results relating the smoothness of the target function and the prescribed approximation error to the size of the approximating network were later obtained in [4], [5].

In the past two decades the rectified linear unit (ReLU) has become the dominant activation function in theory and practice. Beginning with [6], quantitative approximation results for deep ReLU networks have been developed [7], [8], [9], culminating in [10], which shows that deep ReLU networks approximate a wide range of function classes in metric-entropy–optimal manner.

The quantitative approximation results in the literature typically take the following form: for a given function $f$ and a given approximation error $\epsilon > 0$, there exists a neural network $\mathcal{N}$ that approximates $f$ to within error $\epsilon$, formalized as

$$\forall f : \forall \epsilon : \exists \mathcal{N} \text{ such that } \mathcal{N} \text{ approximates } f \text{ to within error } \epsilon. \tag{1}$$

Thus the network architecture and weights depend on the chosen value of $\epsilon$. If a smaller error is later required, a new (typically larger) network must be instantiated. For instance, [10, Proposition III.5] shows that for every $\epsilon > 0$, there exists a deep ReLU network of $\epsilon$-independent width and depth $\mathcal{O}(\log(\epsilon^{-1}))$ that approximates polynomials $f$ to within error $\epsilon$.

---

[*]Swiss National Bank, Börsenstrasse 15, 8001 Zürich, Switzerland

[†]ETH Zürich, Chair for Mathematical Information Science, Sternwartstrasse 7, 8092 Zürich, Switzerland

In the present paper, we propose a new approximation paradigm that avoids this (unstructured) dependence of the network on the approximation error $\epsilon$. Specifically, we aim to reverse the order of quantifiers $\forall \epsilon, \exists \mathcal{N}$ to get a statement of the form

$$\forall f : \exists \mathcal{N}, \text{such that } \forall \epsilon : \mathcal{N} \text{ approximates } f \text{ to within error } \epsilon. \tag{2}$$

Concretely, this will be achieved by repeatedly self-composing a single neural network of fixed topology and fixed weights until the desired accuracy is reached. We show that the resulting approximation error can be made arbitrarily small and, in fact, decays exponentially in the number of compositions. The nature of this construction lends itself to a temporal formulation, specifically in terms of recurrent neural networks (RNNs).

**Definition 1** (RNN [11], [12])**.** *We denote by $\rho : \mathbb{R} \to \mathbb{R}, \rho(x) := max(0, x)$ the rectified linear unit (ReLU) function which acts component-wise, i.e., $\rho(x_1, \ldots, x_m) := (\rho(x_1), \ldots, \rho(x_m))$. An RNN with input dimension $d \in \mathbb{N}$, output dimension $d' \in \mathbb{N}$, and hidden state size $m \in \mathbb{N}$ is parametrized by matrices $A_h \in \mathbb{R}^{m \times m}$, $A_x \in \mathbb{R}^{m \times d}$, $A_o \in \mathbb{R}^{d' \times m}$, and vectors $b_h \in \mathbb{R}^m$, $b_o \in \mathbb{R}^{d'}$. These quantities, collectively called the weights of the RNN, specify the hidden state operator $\mathcal{K} : (\mathbb{R}^d)^{\mathbb{N}_0} \to (\mathbb{R}^m)^{\mathbb{N}_0}$ mapping an input sequence $(x[t])_{t \in \mathbb{N}_0}$ recursively to the sequence of hidden states $(h[t])_{t \in \mathbb{N}_0}$ according to*

$$h[-1] := 0 \in \mathbb{R}^m$$
$$(\mathcal{K}x)[t] = h[t] = \rho(A_h h[t-1] + A_x x[t] + b_h),$$

*and the output mapping $\mathcal{Q} : \mathbb{R}^m \to \mathbb{R}^{d'}$,*

$$\mathcal{Q}(h) := A_o h + b_o.$$

*The associated RNN is the operator $\mathcal{R} : (\mathbb{R}^d)^{\mathbb{N}_0} \to (\mathbb{R}^{d'})^{\mathbb{N}_0}$ given by*

$$\mathcal{R} := \mathcal{Q}\mathcal{K},$$

*with*

$$(\mathcal{R}x)[t] = (\mathcal{Q}\mathcal{K}x)[t] = \mathcal{Q}\left((\mathcal{K}x)[t]\right), \quad t \in \mathbb{N}_0.$$

*We use the notation $\mathcal{M}_{\text{in}}(\mathcal{R}) = d$, $\mathcal{M}_{\text{out}}(\mathcal{R}) = \mathcal{M}_{\text{out}}(\mathcal{Q}) = d'$, $\mathcal{M}_{\text{hid}}(\mathcal{R}) = \mathcal{M}_{\text{hid}}(\mathcal{K}) = m$ to describe the size of the RNN.*

To conform with (2), we wish, for a given function $f$, to find an RNN that achieves any (arbitrarily small) approximation error $\epsilon$ provided we let it run sufficiently long. To this end, if we desire an approximation of the function $f$ at the point $x \in \mathbb{R}$, we take the input sequence of the RNN as $\widetilde{x}[t] = x \mathbb{1}_{\{t=0\}}, t \in \mathbb{N}_0$. Here, $\mathbb{1}_{\{\cdot\}}$ denotes the truth function which takes on the value 1 if the statement inside $\{\cdot\}$ is true and equals 0 otherwise. To formalize this, we introduce the following operator.

**Definition 2.** *The mapping $\mathcal{D} : \mathbb{R}^d \to (\mathbb{R}^d)^{\mathbb{N}_0}$ is defined according to*

$$(\mathcal{D}x)[t] = x \mathbb{1}_{\{t=0\}}, \qquad t \in \mathbb{N}_0.$$

The corresponding output sequence of the RNN then produces increasingly accurate approximations of $f(x)$ as time $t$ evolves. We can operationalize the approximation paradigm (2) in terms of RNNs as

$$\forall f : \exists \mathcal{R} : \forall \epsilon : \exists t_0 : \sup_{t \geq t_0} \sup_x |(\mathcal{R}\mathcal{D}x)[t] - f(x)| < \epsilon. \tag{3}$$

Every approximation theorem fitting this paradigm must have the size, topology, and weights of the approximating RNN $\mathcal{R}$ be independent of the approximation error $\epsilon$, simply by virtue of the quantifier order in (3). Only the runtime required to achieve the desired approximation error $\epsilon$ will depend on $\epsilon$. This approximation paradigm exhibits interesting practical properties as storing the fixed (in general small) RNN on digital devices requires little memory and the approximation error can be controlled simply by adjusting the runtime of the RNN.

In the presentation of the proposed new approximation paradigm, we have deliberately kept (1)-(3) slightly vague in order to convey the central idea of reversing the quantifier order unencumbered

by burdensome technicalities. We believe that neural network approximation results following the paradigm (3) are possible for a variety of function classes. The present paper aims at initiating such a neural network approximation theory program by considering the approximation of univariate polynomials. Note that, while the network is independent of the desired approximation error $\varepsilon$, it does depend on the function $f$ to be approximated. In particular, as Theorem 3 below shows, the size of the resulting RNN depends on the degree of the polynomial. The techniques and specific neural network constructions we develop could prove useful more generally. For example, we think that our results can be extended to multivariate polynomials with relative ease.

We next state the central result of this paper, which says that the class of univariate polynomials can be approximated by RNNs according to the paradigm (3). Moreover, and perhaps surprisingly, the approximation error decreases exponentially in the runtime of the RNN.

**Theorem 3.** *Let $N \in \mathbb{N}$, $a_0, \ldots, a_N \in \mathbb{R}$, and $D \geq 1$. There exists an RNN $\mathcal{R}_a$ such that for $t \geq 16 \log(N)$,*

$$\sup_{x \in [-D,D]} \left| (\mathcal{R}_a \mathcal{D} x)[t] - \sum_{i=0}^{N} a_i x^i \right| \leq \|a\|_1 C_1 4^{-C_2 t},$$

*with $C_1 = 16 N D^{2N}$ and $C_2 = \frac{1}{4\lceil \log(N) \rceil}$.*

A more detailed version of this result is presented in Theorem 33. The remainder of the paper builds up the technical material needed in the proof of Theorem 33. We close this section with some notation conventions.

**Notation.** $\mathbb{N}$ and $\mathbb{N}_0$ denote the set of natural numbers excluding and, respectively, including zero. $(\mathbb{R}^d)^{\mathbb{N}_0}$ stands for the space of sequences of vectors in $\mathbb{R}^d$ indexed by time in $\mathbb{N}_0$, that is for $x[\cdot] \in (\mathbb{R}^d)^{\mathbb{N}_0}$, we have $x[t] \in \mathbb{R}^d$, for $t \in \mathbb{N}_0$. For $x \in \mathbb{R}^d$, we let $\|x\|_\infty = \max_{i=1,\ldots,d} |x_i|$ and write $x \preccurlyeq c$, for some $c > 0$, to indicate that $x_i \leq c$, for all $i \in \{1, \ldots, d\}$. $\mathbb{I}_d \in \mathbb{R}^{d \times d}$ is the $d$-dimensional identity matrix and $\mathbf{1}_d \in \mathbb{R}^d$ stands for the $d$-dimensional vector with all entries equal to one. $\log$ denotes the logarithm to base 2. Constants are throughout understood to be in $\mathbb{R}^+$. For $a, b \in \mathbb{N}_0$, set $\{a, \ldots, b\} := \mathbb{N}_0 \cap [a, b]$. In particular, $\{a, \ldots, b\} = \varnothing$ whenever $b < a$.

## 1.1 Relationship to deep neural network approximation theory

As an aside, it is instructive to compare RNNs with deep feed-forward neural networks. This subsection is not needed for the remainder of the paper, but provides additional context. We begin by recalling the definition of a deep feed-forward neural network.

**Definition 4** (Deep Neural Network). *[10, Definition II.1] Let $L, n_0, n_1, \ldots, n_L \in \mathbb{N}$ with $L \geq 2$. A mapping $\Phi : \mathbb{R}^{n_0} \to \mathbb{R}^{n_L}$ defined by*

$$\Phi = \begin{cases} W_2 \circ \rho \circ W_1, & L = 2, \\ W_L \circ \rho \circ W_{L-1} \circ \rho \circ \cdots \circ \rho \circ W_1, & L \geq 3, \end{cases}$$

*with affine maps $W_\ell : \mathbb{R}^{n_{\ell-1}} \to \mathbb{R}^{n_\ell}$, $\ell \in \{1, \ldots, L\}$, and activation function $\rho$, is called a feed-forward neural network. The $W_\ell$ are given by $W_\ell(x) = A_\ell x + b_\ell$, with $A_\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$ and $b_\ell \in \mathbb{R}^{n_\ell}$.*

Now let $f$ be a fixed polynomial and let $\epsilon > 0$ be the desired approximation error. By Theorem 3, there exists an RNN $\mathcal{R}_f$ such that

$$\sup_{x \in [-D,D]} \left| (\mathcal{R}_f \mathcal{D} x)[T] - f(x) \right| \leq \epsilon,$$

for $T := \lceil c \log(\epsilon^{-1}) \rceil$ with a constant $c > 0$ depending only on $f$.

For this fixed $T$, the map $x \mapsto (\mathcal{R}_f \mathcal{D} x)[T]$ can be expressed as a deep feed-forward network with $T + 2$ layers. We now make this equivalence explicit. Recall that the RNN input is given by $(\mathcal{D} x)[0] = x$ and $(\mathcal{D} x)[t] = 0$ for $t \geq 1$, and that the hidden state is initialized as $h[-1] = 0$. The RNN updates therefore take the form

$$h[0] = \rho(A_x x + b_h), \qquad h[t] = \rho(A_h h[t-1] + b_h) \ (t \geq 1), \qquad (\mathcal{R}_f \mathcal{D} x)[T] = A_o h[T] + b_o.$$

Define the affine maps

$$W_1(x) := A_x x + b_h, \qquad W_{\text{share}}(h) := A_h h + b_h, \qquad W_{T+2}(h) := A_o h + b_o,$$

and set $W_\ell := W_{\text{share}}$, for $\ell = 2, \ldots, T+1$. Iterating the recurrence yields

$$h[T] = \rho \circ W_{T+1} \circ \cdots \circ \rho \circ W_1(x),$$

and hence

$$(\mathcal{R}_f \mathcal{D} x)[T] = W_{T+2}(\rho \circ W_{T+1} \circ \rho \circ \cdots \circ \rho \circ W_1(x)).$$

Thus unfolding $\mathcal{R}_f$ over $T$ time steps yields the deep network

$$\Phi = W_{T+2} \circ \rho \circ W_{T+1} \circ \rho \circ \cdots \circ \rho \circ W_1,$$

with $T+2$ layers and shared weights in layers 2 through $T+1$.

We may now compare this with [10, Proposition III.5], which shows that for every $\epsilon > 0$ there exists a deep ReLU network $\Phi'_f$ of $\epsilon$-independent (constant) width and $\mathcal{O}(\log(\epsilon^{-1}))$ layers that approximates $f$ to within error $\epsilon$. Since unfolding $\mathcal{R}_f$ for $T = \mathcal{O}(\log(\epsilon^{-1}))$ time steps likewise produces a network of constant width (as shown later in the paper) and depth $\mathcal{O}(\log(\epsilon^{-1}))$, we obtain the same connectivity–error trade-off directly from Theorem 3. The key difference is that the construction in [10] requires redesigning the entire network $\Phi'_f$ for each new value of $\epsilon$, whereas our RNN-based method uses a single fixed architecture: to achieve a smaller error, one simply runs the network longer, repeatedly applying the shared layer $W_{\text{share}}$.

**Conceptual remark.** Classical deep neural network approximation theory, particularly the structural calculus developed in [10], relies on a small set of fundamental operations: composition of networks, parallelization, linear combinations, pre-and post-composition with affine maps, and depth extension via identity mappings.

A central theme of the present paper is that each of these operations admits an RNN analogue, while strictly adhering to the weight sharing paradigm described above. As demonstrated by the unfolding argument, RNNs emulate deep feed-forward networks by reusing the same weights over time.

Section 3 develops this novel RNN calculus explicitly: parallelization (Lemma 9), affine pre- and post-composition (Lemmata 10 and 11), clocked concatenation for function composition (Theorem 15), and the structured assembly of multiple concatenations (Lemma 16 and Corollary 17). These constructions show that the expressive power of deep feed-forward ReLU networks can equivalently be attained temporally within a single recurrent architecture through weight sharing, with clocked concatenation providing the principal mechanism for implementing compositional depth.

## 2 Approximating the squaring function and multiplication

The constructions developed in this section correspond to the basic operations used in classical deep feed-forward ReLU network approximation theory, in particular squaring and multiplication. In the classical setting, these operations are assembled through a structural calculus—compositions, parallelization, and affine transformations—to produce approximations of more complex functions such as polynomials. Our RNN-based approach reproduces this calculus, but within a recurrent architecture. We design RNNs that approximate squaring and multiplication and then combine them into polynomial approximations using the mechanisms developed in Section 3, most notably the clocked concatenation and weight sharing ideas inherent to the RNN framework.

Before delving into the details of our RNN constructions, we point out a simplification of the RNN definition.

**Remark 5.** *As we shall only be concerned with inputs of the form $(\mathcal{D}x)[t] = x \mathbb{1}_{\{t=0\}}$, for $x \in \mathbb{R}^d$, Definition 1 simplifies to*

$$
\begin{aligned}
(\mathcal{K}\mathcal{D}x)[0] &= \rho(A_x x + b_h) \\
(\mathcal{K}\mathcal{D}x)[t] &= \rho(A_h ((\mathcal{K}\mathcal{D}x)[t-1]) + b_h), && t \in \mathbb{N},
\end{aligned}
$$

*and*

$$(\mathcal{R}\mathcal{D}x)[t] = (\mathcal{Q}\mathcal{K}\mathcal{D}x)[t] = \mathcal{Q}\left(\left(\mathcal{K}\mathcal{D}x\right)[t]\right), \qquad\qquad t \in \mathbb{N}_0.$$

*As a notational convention, we extend the sequences $((K\mathcal{D}x)[t])_{t\in\mathbb{N}_0}$ and $((R\mathcal{D}x)[t])_{t\in\mathbb{N}_0}$ to negative time indices by setting*

$$(K\mathcal{D}x)[t] = 0 \quad and \quad (R\mathcal{D}x)[t] = 0, \qquad for\ all\ t < 0.$$

*This convention is used for expositional convenience only and does not affect the RNN dynamics on $t \geq 0$.*

We start by devising an RNN that approximates the squaring function. To this end we make use of the following feedforward neural network construction from [6], [10].

**Lemma 6.** *Let $F(x) := x - x^2$, $x \in [0, 1]$. Further, with $m \in \mathbb{N}$, let $I_m : [0, 1] \to [0, 1]$ be the piecewise linear interpolant of $F$ at the points $\frac{k}{2^m}$, with $k \in \{0, \ldots, 2^m\}$, that is,*

$$I_m\left(\frac{k}{2^m}\right) = F\left(\frac{k}{2^m}\right), \quad for\ k \in \{0, \ldots, 2^m\},$$

*and $I_m$ is affine on the intervals $\left[\frac{k}{2^m}, \frac{k+1}{2^m}\right]$, $k \in \{0, \ldots, 2^m - 1\}$. It holds that*

$$\sup_{x \in [0,1]} |F(x) - I_m(x)| = 2^{-2m-2}.$$

*Furthermore, define*

$$s_\ell(\cdot) := 2^{-1}\rho(\cdot) - \rho(\cdot - 2^{-2\ell-1}), \quad \ell \in \mathbb{N}_0,$$

*and recursively $H_\ell = s_\ell \circ H_{\ell-1}$, for $\ell \in \mathbb{N}$, with $H_0 = s_0$. Then, it holds that*

$$I_m(x) = \sum_{\ell=0}^{m-1} H_\ell(x), \qquad for\ x \in [0, 1].$$

*Proof.* See the proof of Proposition III.2 in [10]. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

We furthermore need the following property of $I_m$, which is not stated in [10].

**Corollary 7.** *For every $m \in \mathbb{N}$ and all $x \in [0, 1]$,*

$$0 \leq I_m(x) \leq x,$$

*where $I_m$ is defined in Lemma 6.*

*Proof.* First, note that $F(x) = x - x^2 \in [0, 1/4]$, for all $x \in [0, 1]$. Thus, $I_m(x) \in [0, 1/4]$, for all $x \in [0, 1]$, and in particular $I_m(x) \geq 0$. Furthermore, since $F(x) = x - x^2$ is concave, we have $I_m(x) \leq F(x)$, for all $x \in [0, 1]$. This implies

$$I_m(x) \leq F(x) = x - x^2 \leq x, \quad for\ all\ x \in [0, 1],$$

thereby completing the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

We are now ready to construct an RNN that approximates the squaring function.

**Theorem 8.** *For $D \geq 1$, there is an RNN $\mathcal{R}^{Sq} = \mathcal{Q}^{Sq}\mathcal{K}^{Sq}$, with $\mathcal{M}_{\mathrm{in}}(\mathcal{R}^{Sq}) = \mathcal{M}_{\mathrm{out}}(\mathcal{R}^{Sq}) = 1$, $\mathcal{M}_{\mathrm{hid}}(\mathcal{R}^{Sq}) = 7$, such that, for all $x \in [-D, D]$ and all $t \in \mathbb{N}_0$,*

$$|(\mathcal{R}^{Sq}\mathcal{D}x)[t] - x^2| \leq \frac{D^2}{4}\, 4^{-t},$$

*as well as*

$$\left\|(\mathcal{K}^{Sq}\mathcal{D}x)[t]\right\|_\infty \leq 1 \ \ and \ \ 0 \leq (\mathcal{R}^{Sq}\mathcal{D}x)[t] \leq D^2.$$

*Proof.* The weights of $\mathcal{R}^{\mathsf{Sq}}$ are as follows:

$$
A_h = \begin{pmatrix}
0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 1 & 2^{-1} & -1 & 0 & 0 & 0 \\
1 & 1 & 2^{-1} & -1 & 0 & -1 & 0 \\
1 & 1 & -2^{-1} & 1 & 1 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 2^{-2} & -2^{-1} \\
0 & 0 & 0 & 0 & 0 & 0 & 0
\end{pmatrix},
\qquad
b_h = \begin{pmatrix}
0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 2^{-1} \\ 1
\end{pmatrix},
$$

$$
A_x = \frac{1}{D}\begin{pmatrix}
1 \\ -1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0
\end{pmatrix},
\qquad
A_o = D^2 \begin{pmatrix} 0 & 0 & -2^{-1} & 1 & 1 & 0 & 0 \end{pmatrix},
\qquad b_o = 0.
$$

Next, we arbitrarily fix $x \in [-D, D]$, let $z := \frac{|x|}{D}$, and compute the sequence $h[\cdot] := (\mathcal{K}\mathcal{D}x)[\cdot]$ of hidden states corresponding to the input sequence $\widetilde{x} := (\mathcal{D}x)[\cdot] = x\mathbb{1}_{\{\cdot=0\}}$ according to Definition 1. We start by proving through induction that

$$
h[t] = \begin{pmatrix}
0 \\
0 \\
\rho(H_{t-2}(z)) \\
\rho(H_{t-2}(z) - 2^{-2t+1}) \\
\dfrac{z - \sum_{i=0}^{t-2} H_i(z)}{2^{-2t-1}} \\
1
\end{pmatrix},
\qquad \text{for all } t \geq 2,
\tag{4}
$$

with $H_t$, $t \in \mathbb{N}_0$, as defined in Lemma 6. Starting from $h[-1] = 0$, we compute

$$
\begin{aligned}
h[0] &= \rho(A_h h[-1] + A_x \widetilde{x}[0] + b_h) \\
&= \rho(A_x x + b_h) \\
&= \begin{pmatrix}
\rho(\frac{x}{D}) \\
\rho(-\frac{x}{D}) \\
0 \\
0 \\
0 \\
2^{-1} \\
1
\end{pmatrix}.
\end{aligned}
$$

Using $z = \rho\left(\frac{x}{D}\right) + \rho\left(-\frac{x}{D}\right) \geq 0$, we get

$$h[1] = \rho(A_h h[0] + A_x \widetilde{x}[1] + b_h)$$
$$= \rho(A_h h[0] + b_h)$$
$$= \begin{pmatrix} 0 \\ 0 \\ \hline z \\ \rho(z - 2^{-1}) \\ z \\ \hline 2^{-3} \\ 1 \end{pmatrix}.$$

Further,

$$h[2] = \rho(A_h h[1] + b_h)$$

$$= \rho \begin{pmatrix} 0 \\ 0 \\ \hline 2^{-1}\rho(z) - \rho(z - 2^{-1}) \\ 2^{-1}\rho(z) - \rho(z - 2^{-1}) - 2^{-3} \\ z - (2^{-1}z - \rho(z - 2^{-1})) \\ \hline 2^{-2} \cdot 2^{-3} - 2^{-1} + 2^{-1} \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \hline \rho(s_0(z)) \\ \rho(s_0(z) - 2^{-3}) \\ z - s_0(z) \\ \hline 2^{-2 \cdot 2 - 1} \\ 1 \end{pmatrix},$$

where $s_t(\cdot), t \in \mathbb{N}_0$, was defined in Lemma 6. Hence, recalling that $H_0 = s_0$ as per Lemma 6, we established the base case $t = 2$ of the induction.

We proceed to the induction step. To this end, we assume that (4) holds for some $t \geq 2$. Now, note that

$$h[t+1] = \rho(A_h h[t] + b_h)$$

$$= \rho \begin{pmatrix} 0 \\ 0 \\ \hline 2^{-1}\rho(H_{t-2}(z)) - \rho(H_{t-2}(z) - 2^{-2t+1}) \\ 2^{-1}\rho(H_{t-2}(z)) - \rho(H_{t-2}(z) - 2^{-2t+1}) - 2^{-2t-1} \\ z - \sum_{i=0}^{t-2} H_i(z) - \left(2^{-1}\rho(H_{t-2}(z)) - \rho(H_{t-2}(z) - 2^{-2t+1})\right) \\ 2^{-2} \cdot 2^{-2t-1} - 2^{-1} + 2^{-1} \\ 1 \end{pmatrix}$$

$$= \rho \begin{pmatrix} 0 \\ 0 \\ \hline s_{t-1}(H_{t-2}(z)) \\ s_{t-1}(H_{t-2}(z)) - 2^{-2t-1} \\ z - \sum_{i=0}^{t-2} H_i(z) - s_{t-1}(H_{t-2}(z)) \\ \hline 2^{-2(t+1)-1} \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \hline \rho(H_{t-1}(z)) \\ \rho(H_{t-1}(z) - 2^{-2(t+1)+1}) \\ z - \sum_{i=0}^{t-1} H_i(z) \\ \hline 2^{-2(t+1)-1} \\ 1 \end{pmatrix},$$

where we used $s_{t-1}(\cdot) = 2^{-1}\rho(\cdot) - \rho(\cdot - 2^{-2t+1})$ and $H_t = s_t \circ H_{t-1}$, both as per Lemma 6. This completes the proof of (4). Furthermore, we see directly that $\|h[0]\|_\infty \leq 1$ and $\|h[1]\|_\infty \leq 1$. By (4), we obtain $\|h[t]\|_\infty \leq 1$, for all $t \geq 2$, upon using Corollary 7.

Next, we compute the RNN output $(\mathcal{R}^{\mathsf{Sq}}x)[t]$, for $t \geq 2$, as follows

$$
\begin{aligned}
(\mathcal{R}^{\mathsf{Sq}}x)[t] &= A_o h[t] + b_o \\
&= D^2 \left( z - \sum_{i=0}^{t-2} H_i(z) - \left( 2^{-1}\rho(H_{t-2}(z)) - \rho(H_{t-2}(z) - 2^{-2t+1}) \right) \right) \\
&= D^2 \left( z - \sum_{i=0}^{t-2} H_i(z) - H_{t-1}(z) \right) \\
&= D^2 \left( z - \sum_{i=0}^{t-1} H_i(z)) \right) \\
&= D^2 \left( z - I_t(z) \right).
\end{aligned}
$$

As $z \in [0,1]$, we have, by Corollary 7, that $z - I_t(z) \in [0,1]$. Hence, $0 \leq (\mathcal{R}^{\mathsf{Sq}}x)[t] \leq D^2$, for all $t \geq 2$. Furthermore,

$$
\begin{aligned}
\left| x^2 - (\mathcal{R}^{\mathsf{Sq}}\mathcal{D}x)[t] \right| &= D^2 \left| z^2 - (z - I_t(z)) \right| \\
&\leq D^2 \sup_{y \in [0,1]} \left| y^2 - (y - I_t(y)) \right| \\
&\leq D^2 \, 2^{-2t-2} = \frac{D^2}{4} 4^{-t},
\end{aligned}
$$

where the second inequality follows from Lemma 6. The proof is finalized upon noting that $x \in [-D, D]$ was arbitrary. □

A result that is, at first glance, similar to Theorem 8 was established in [13, Theorem 4.4], albeit using an RNN architecture that does not lead to quantifier reversal.

We continue our RNN construction program by approximating the multiplication operation. This will be done through the polarization identity

$$
x_1 \cdot x_2 = \left( \frac{x_1 + x_2}{2} \right)^2 - \left( \frac{x_1 - x_2}{2} \right)^2. \tag{5}
$$

Specifically, we will first map the input $(x_1, x_2)$ to $((x_1 + x_2)/2, (x_1 - x_2)/2)$ through an affine transformation and then apply two instances of $\mathcal{R}^{\mathsf{Sq}}$ in parallel followed by a linear combination of their outputs. We proceed to develop auxiliary results needed in the construction of the RNN approximating the multiplication operation. The first lemma provides a formal way to run several RNNs in parallel inside one larger RNN.

**Lemma 9.** *For $N \in \mathbb{N}$, let $\mathcal{R}^1 = \mathcal{Q}^1\mathcal{K}^1, \ldots, \mathcal{R}^N = \mathcal{Q}^N\mathcal{K}^N$ be RNNs. There exists an RNN $\mathcal{R} = \mathcal{Q}\mathcal{K}$ such that for all $x_1 \in \mathbb{R}^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^1)}, \ldots, x_N \in \mathbb{R}^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^N)}$, and all $t \in \mathbb{N}_0$,*

$$
\left( \mathcal{K}\mathcal{D} \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix} \right) [t] = \begin{pmatrix} (\mathcal{K}^1\mathcal{D}x_1)[t] \\ \vdots \\ (\mathcal{K}^N\mathcal{D}x_N)[t] \end{pmatrix}, \tag{6}
$$

*and*

$$
\left( \mathcal{R}\mathcal{D} \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix} \right) [t] = \begin{pmatrix} (\mathcal{R}^1\mathcal{D}x_1)[t] \\ \vdots \\ (\mathcal{R}^N\mathcal{D}x_N)[t] \end{pmatrix}. \tag{7}
$$

*Proof.* The weights of $\mathcal{R} = \mathcal{QK}$ are given by

$$A_h = \begin{pmatrix} A_h^1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & A_h^N \end{pmatrix}, \qquad A_x = \begin{pmatrix} A_x^1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & A_x^N \end{pmatrix}, \qquad b_h = \begin{pmatrix} b_h^1 \\ \vdots \\ b_h^N \end{pmatrix},$$

$$A_o = \begin{pmatrix} A_o^1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & A_o^N \end{pmatrix}, \qquad b_o = \begin{pmatrix} b_o^1 \\ \vdots \\ b_o^N \end{pmatrix}.$$

Now, arbitrarily fix $x = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix}$. We first show (6), by induction. Indeed, for $t = 0$, we have by Remark 5,

$$(\mathcal{KD}x)[0] = \rho(A_x x + b_h) = \rho \begin{pmatrix} A_x^1 x_1 + b_h^1 \\ \vdots \\ A_x^N x_N + b_h^N \end{pmatrix} = \begin{pmatrix} (\mathcal{K}^1 \mathcal{D}x_1)[0] \\ \vdots \\ (\mathcal{K}^N \mathcal{D}x_N)[0] \end{pmatrix}.$$

Next, assume that (6) holds for some $t \in \mathbb{N}_0$ and compute

$$(\mathcal{KD}x)[t+1] = \rho\left( A_h \left( (\mathcal{KD}x)[t] \right) + b_h \right)$$

$$= \rho \left( A_h \begin{pmatrix} (\mathcal{K}^1 \mathcal{D}x_1)[t] \\ \vdots \\ (\mathcal{K}^N \mathcal{D}x_N)[t] \end{pmatrix} + b_h \right)$$

$$= \rho \begin{pmatrix} A_h^1 \left( (\mathcal{K}^1 \mathcal{D}x_1)[t] \right) + b_h^1 \\ \vdots \\ A_h^N \left( (\mathcal{K}^N \mathcal{D}x_N)[t] \right) + b_h^N \end{pmatrix} = \begin{pmatrix} (\mathcal{K}^1 \mathcal{D}x_1)[t+1] \\ \vdots \\ (\mathcal{K}^N \mathcal{D}x_N)[t+1] \end{pmatrix}.$$

This completes the proof of (6). To establish (7), we note that

$$(\mathcal{RD}x)[t] = A_o(\mathcal{KD}x)[t] + b_o = A_o \begin{pmatrix} (\mathcal{K}^1 \mathcal{D}x_1)[t] \\ \vdots \\ (\mathcal{K}^N \mathcal{D}x_N)[t] \end{pmatrix} + b_o$$

$$= \begin{pmatrix} A_o^1(\mathcal{K}^1 \mathcal{D}x_1)[t] + b_o^1 \\ \vdots \\ A_o^N(\mathcal{K}^N \mathcal{D}x_N)[t] + b_o^N \end{pmatrix} = \begin{pmatrix} (\mathcal{R}^1 \mathcal{D}x_1)[t] \\ \vdots \\ (\mathcal{R}^N \mathcal{D}x_N)[t] \end{pmatrix}, \qquad t \in \mathbb{N}_0. \qquad \square$$

The next result shows how a fixed linear map can be absorbed into the first layer of an RNN.

**Lemma 10.** *Let $\mathcal{R} = \mathcal{QK}$ be an RNN with $\mathcal{M}_{\text{in}}(\mathcal{R}) = d$ and let $A \in \mathbb{R}^{d \times d'}$, $d' \in \mathbb{N}$. Then, there exists an RNN $\mathcal{R}' = \mathcal{QK}'$ such that, for all $x \in \mathbb{R}^{d'}$, and all $t \in \mathbb{N}_0$,*

$$(\mathcal{R}' \mathcal{D}x)[t] = (\mathcal{RD}(Ax))[t] \qquad and \qquad (\mathcal{K}' \mathcal{D}x)[t] = (\mathcal{KD}(Ax))[t].$$

*Proof.* Let $A_h, A_x, A_o, b_h, b_o$ be the weights of the RNN $\mathcal{R}$. We define the corresponding RNN $\mathcal{R}' = \mathcal{QK}'$ with weights $A'_h = A_h, A'_x = A_x A, A'_o = A_o, b'_h = b_h, b'_o = b_o$, in particular the output mapping of $\mathcal{R}'$ is identical to that of $\mathcal{R}$. First, we establish, through induction, that

$$(\mathcal{K}' \mathcal{D}x)[t] = (\mathcal{KD}(Ax))[t], \qquad \text{for all } t \in \mathbb{N}_0. \tag{8}$$

9

The base case follows from

$$(\mathcal{K}'\mathcal{D}x)[0] = \rho\left(A_x'x + b_h'\right)$$
$$= \rho\left(A_x Ax + b_h\right)$$
$$= \rho\left(A_x(Ax) + b_h\right) = (\mathcal{K}\mathcal{D}(Ax))[0].$$

Next, assume that (8) holds for some $t \in \mathbb{N}_0$. We have by Remark 5,

$$(\mathcal{K}'\mathcal{D}x)[t+1] = \rho\left(A_h'\left((\mathcal{K}'\mathcal{D}x)[t]\right) + b_h'\right)$$
$$= \rho\left(A_h\left((\mathcal{K}'\mathcal{D}x)[t]\right) + b_h\right)$$
$$= \rho\left(A_h\left((\mathcal{K}\mathcal{D}(Ax))[t]\right) + b_h\right)$$
$$= (\mathcal{K}\mathcal{D}(Ax))[t+1],$$

which concludes the induction argument. Finally, $(\mathcal{R}'\mathcal{D}x)[t] = (\mathcal{R}\mathcal{D}(Ax))[t]$, for all $t \in \mathbb{N}_0$, follows from (8) since $A_o' = A_o$ and $b_o' = b_o$. $\qquad\square$

The last auxiliary result needed in the approximation of the multiplication operation shows that the application of an affine transformation to the output of an RNN can be absorbed into the RNN.

**Lemma 11.** *Let $\mathcal{R} = \mathcal{Q}\mathcal{K}$ be an RNN with $\mathcal{M}_{\mathrm{out}}(\mathcal{R}) = d$ and let $A \in \mathbb{R}^{d' \times d}, b \in \mathbb{R}^{d'}, d' \in \mathbb{N}$. Then, there exists an RNN $\mathcal{R}' = \mathcal{Q}'\mathcal{K}$ such that, for all $x \in \mathbb{R}^{\mathcal{M}_{\mathrm{in}}(\mathcal{R})}$, and all $t \in \mathbb{N}_0$,*

$$(\mathcal{R}'\mathcal{D}x)[t] = A\left((\mathcal{R}\mathcal{D}x)[t]\right) + b.$$

*Proof.* For $\mathcal{R} = \mathcal{Q}\mathcal{K}$, where $\mathcal{Q}(h) = A_o h + b_o$, define

$$\mathcal{R}' = \mathcal{Q}'\mathcal{K}, \quad \mathcal{Q}'(h) = A_o'h + b_o', \text{ with } A_o' = AA_o, b_o' = Ab_o + b.$$

That is, we take the hidden state operator of the modified RNN $\mathcal{R}'$ to be identical to that of $\mathcal{R}$. Noting that

$$(\mathcal{R}'\mathcal{D}x)[t] = (\mathcal{Q}'\mathcal{K}\mathcal{D}x)[t] = A_o'\left((\mathcal{K}\mathcal{D}x)[t]\right) + b_o'$$
$$= AA_o\left((\mathcal{K}\mathcal{D}x)[t]\right) + Ab_o + b$$
$$= A\left(A_o(\mathcal{K}\mathcal{D}x)[t] + b_o\right) + b$$
$$= A(\mathcal{R}\mathcal{D}x)[t] + b, \qquad \text{for all } t \in \mathbb{N}_0,$$

the proof is completed. $\qquad\square$

We can now proceed to the derivation of the approximation result for the multiplication operation.

**Theorem 12.** *For $D \geq 1$, there is an RNN $\mathcal{R}^\times = \mathcal{Q}^\times\mathcal{K}^\times$ such that, for all $x = (x_1, x_2) \in [-D, D]^2$, and all $t \in \mathbb{N}_0$,*

$$\left|(\mathcal{R}^\times\mathcal{D}x)[t] - (x_1 \cdot x_2)\right| \leq \frac{D^2}{2}\, 4^{-t},$$

*and*

$$\|\mathcal{K}^\times\mathcal{D}x)[t]\|_\infty \leq 1, \quad |(\mathcal{R}^\times\mathcal{D}x)[t]| \leq D^2.$$

*Furthermore, $\mathcal{M}_{\mathrm{hid}}(\mathcal{R}^\times) = 14$.*

*Proof.* We start by employing Lemma 9 with $N = 2$ and $\mathcal{R}^1 = \mathcal{R}^{\mathrm{Sq}}, \mathcal{R}^2 = \mathcal{R}^{\mathrm{Sq}}$ according to Theorem 8. This yields an RNN $\mathcal{R} = \mathcal{Q}\mathcal{K}$ satisfying, for all $z = (z_1, z_2) \in [-D, D]^2$,

$$\left\|(\mathcal{R}\mathcal{D}z)[t] - \begin{pmatrix} z_1^2 \\ z_2^2 \end{pmatrix}\right\|_\infty \leq \frac{D^2}{4}4^{-t}, \quad \|(\mathcal{K}\mathcal{D}z)[t]\|_\infty \leq 1, \quad \text{and } 0 \preccurlyeq (\mathcal{R}\mathcal{D}z)[t] \preccurlyeq D^2, \quad \text{for all } t \in \mathbb{N}_0. \quad (9)$$

Next, consider the matrix

$$A = \frac{1}{2}\begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}, \tag{10}$$

and observe that, for all $x \in [-D, D]^2$, $Ax \in [-D, D]^2$. We now set $z = Ax$ in (9) and apply Lemma 10 with $A$ in (10) and $\mathcal{R}$ as per (9) to obtain the RNN $\widetilde{\mathcal{R}} = \mathcal{Q}\widetilde{\mathcal{K}}$ satisfying, for all $x = (x_1, x_2) \in [-D, D]^2$,

$$\left\| (\tilde{\mathcal{R}}\mathcal{D}x)[t] - \begin{pmatrix} (\frac{x_1 + x_2}{2})^2 \\ (\frac{x_1 - x_2}{2})^2 \end{pmatrix} \right\|_\infty \leq \frac{D^2}{4} 4^{-t}, \tag{11}$$

$$\left\| (\widetilde{\mathcal{K}}\mathcal{D}x)[t] \right\|_\infty \leq 1, \quad \text{and}$$

$$0 \preccurlyeq (\widetilde{\mathcal{R}}\mathcal{D}x)[t] \preccurlyeq D^2, \tag{12}$$

for all $t \in \mathbb{N}_0$. Finally, we apply Lemma 11 to $\widetilde{\mathcal{R}}$, with $d' = 1, d = 2, A = \begin{pmatrix} 1 & -1 \end{pmatrix}$, and $b = 0$ to obtain the RNN $\mathcal{R}^\times = \mathcal{Q}^\times \mathcal{K}^\times$, with $\mathcal{K}^\times = \widetilde{\mathcal{K}}$, and note that

$$(\mathcal{R}^\times \mathcal{D}x)[t] = ((\widetilde{\mathcal{R}}\mathcal{D}x)[t])_1 - ((\widetilde{\mathcal{R}}\mathcal{D}x)[t])_2, \qquad \text{for all } t \in \mathbb{N}_0. \tag{13}$$

It thus follows from (12) that $-D^2 \preccurlyeq (\mathcal{R}^\times \mathcal{D}x)[t] \preccurlyeq D^2$, for all $t \in \mathbb{N}_0$. Finally, we have

$$|(\mathcal{R}^\times \mathcal{D}x)[t] - (x_1 \cdot x_2)| \overset{(5)}{=} \left| (\mathcal{R}^\times \mathcal{D}x)[t] - \left( \left( \frac{x_1 + x_2}{2} \right)^2 - \left( \frac{x_1 - x_2}{2} \right)^2 \right) \right|$$

$$\overset{(13)}{=} \left| ((\widetilde{\mathcal{R}}\mathcal{D}x)[t])_1 - \left( \frac{x_1 + x_2}{2} \right)^2 + \left( \frac{x_1 - x_2}{2} \right)^2 - \left( (\widetilde{\mathcal{R}}\mathcal{D}x)[t] \right)_2 \right|$$

$$\leq \left| ((\widetilde{\mathcal{R}}\mathcal{D}x)[t])_1 - \left( \frac{x_1 + x_2}{2} \right)^2 \right| + \left| \left( (\widetilde{\mathcal{R}}\mathcal{D}x)[t] \right)_2 - \left( \frac{x_1 - x_2}{2} \right)^2 \right|$$

$$\overset{(11)}{\leq} \frac{D^2}{4} 4^{-t} + \frac{D^2}{4} 4^{-t} = \frac{D^2}{2} 4^{-t}, \qquad \text{for all } x \in [-D, D], \text{ and all } t \in \mathbb{N}_0.$$

The proof is concluded upon noting that $\mathcal{M}_{\text{hid}}(\mathcal{R}^\times) = 14$ by Lemma 9 applied to $\mathcal{R}^\times$ which contains two parallel instances of $\mathcal{R}^{\text{Sq}}$ with $\mathcal{M}_{\text{hid}}(\mathcal{R}^{\text{Sq}}) = 7$. $\qquad \square$

## 3 Function composition through RNN concatenation

The present section develops the recurrent analogue of the structural operations that play a central role in classical deep feed-forward ReLU network approximation theory. A key objective is to realize function composition within a single recurrent architecture, something that cannot be achieved simply by stacking networks as in the feed-forward setting. Theorem 15 establishes a clocked concatenation mechanism in which two subnetworks operate in parallel and a time-schedule determines when the intermediate output of one serves as the input to the other. Through this mechanism, RNNs emulate feed-forward network composition while maintaining a fixed recursive architecture and reusing the same weights over time.

Our ultimate goal is the approximation of general polynomials in $x$, which requires RNNs that approximate higher powers. This will be realized by composing functions, e.g., the map $x \to x^4$ can be expressed by composing the squaring function with itself. More generally, consider the functions $f, g$ with associated RNNs $\mathcal{R}^f$ and $\mathcal{R}^g$ approximating $f$ and $g$, respectively, in the sense of

$$|(\mathcal{R}^f \mathcal{D}x)[t] - f(x)| \leq c_1 4^{-c_2 t}, \qquad \text{and}$$

$$|(\mathcal{R}^g \mathcal{D}x)[t] - g(x)| \leq c_1 4^{-c_2 t}, \qquad \text{for all } t \in \mathbb{N}_0, \qquad \text{for some } c_1, c_2 > 0.$$

We wish to construct a new RNN $\mathcal{R}'$ that approximates $g \circ f$ such that

$$|(\mathcal{R}'\mathcal{D}x)[t] - g(f(x))| \leq c_3 4^{-c_4 t}, \qquad \text{for all } t \in \mathbb{N}_0, \qquad \text{for some } c_3, c_4 > 0.$$

The core idea is to construct the RNN $\mathcal{R}'$ such that it internally runs $\mathcal{R}^f$ and $\mathcal{R}^g$ in parallel. Specifically, the hidden state vector of $\mathcal{R}'$ consists of two parts, as illustrated in Figure 1. The part labeled $h_f$ follows exactly the same recursion as the hidden state of $\mathcal{R}_f$ and therefore produces the same output
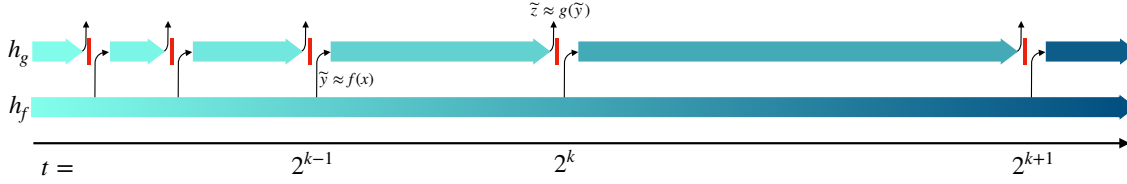
Figure 1: Evolution of the hidden state sequence of the RNN approximating $g \circ f$.

sequence as $\mathcal{R}_f$, yielding progressively more accurate approximations of $f(x)$. The other part, labeled $h_g$, evolves as prescribed by $\mathcal{R}_g$, but with the following modification. Whenever the time index is a power of 2, we reinitialize $h_g$ based on the present output of the subnetwork $\mathcal{R}^f$. Specifically, as visualized in Figure 1, at time index $t = 2^{k-1}$ the hidden state $h_g$ is set to 0. Then, the current value of $h_f$ is used to produce an approximation $\widetilde{y}$ of $f(x)$ which, in turn, is employed to reinitialize $h_g$. Subsequently, at time index $t = 2^k$, an output approximation $\widetilde{z}$ of $g(f(x))$ is produced, where $\mathcal{R}^f$ had been running for $2^{k-1}$ time steps to compute $\widetilde{y} \approx f(x)$ and $\mathcal{R}^g$ ran for $2^{k-1}$ time steps with input $\widetilde{y}$ to approximate $g(\widetilde{y})$. Importantly, $\mathcal{R}^f$ continues to run in parallel such that at time step $t = 2^k$ a refined approximation of $f(x)$ is available to initialize $h_g$ for the next readout at time step $t = 2^{k+1}$.

We next show how to realize the clocking mechanism.

**Lemma 13.** *Define the sequence*

$$\widehat{\delta}[t] = \begin{cases} 1, & \text{if } t = 2^k, \text{ for } k \in \mathbb{N} \text{ with } k \geq 2 \\ 0, & \text{else.} \end{cases}$$

*Let*

$$\widehat{A} := \begin{pmatrix} -4 & 2 & 0 & 0 & 0 \\ -4 & 2 & 0 & 2 & -1/2 \\ 0 & 0 & 1/2 & 0 & -1 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \widehat{b} := \begin{pmatrix} -1 \\ 1/2 \\ 1 \\ -2 \\ 1 \end{pmatrix}, \tag{14}$$

*and consider the recursively defined sequence $h[t] \in \mathbb{R}^5$, with $h[-1] = 0$ and $h[t] = \rho(\widehat{A}h[t-1] + \widehat{b})$, for $t \in \mathbb{N}_0$. It holds that*

$$(h[t])_1 = \widehat{\delta}[t+2]$$

*and $\|h[t]\|_\infty \leq 2$, both for all $t \in \mathbb{N}_0$.*

*Proof.* The proof is effected through induction over time. We start by showing that

$$h[2^k - 2] = \begin{pmatrix} 1 \\ 2 \\ 2^{2-2^k} \\ 0 \\ 1 \end{pmatrix}, \qquad \text{for all } k \in \mathbb{N}, \text{ with } k \geq 2. \tag{15}$$

The base case, $k = 2$, is established through direct computation as

$$h[0] = \rho \begin{pmatrix} -1 \\ 1/2 \\ 1 \\ -2 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 1/2 \\ 1 \\ 0 \\ 1 \end{pmatrix}, \qquad h[1] = \rho \begin{pmatrix} 0 \\ 1 \\ 1/2 \\ -1/2 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 1/2 \\ 0 \\ 1 \end{pmatrix}, \quad \text{and} \quad h[2] = \rho \begin{pmatrix} 1 \\ 2 \\ 2^{-2} \\ -1/2 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 2^{-2} \\ 0 \\ 1 \end{pmatrix}.$$

For the induction step, fix $k \geq 2$ and assume that (15) holds for this $k$. We first compute

$$h[2^k - 1] = \rho \begin{pmatrix} -1 \\ 0 \\ 2^{1-2^k} \\ 2^{2-2^k} \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 2^{1-2^k} \\ 2^{2-2^k} \\ 1 \end{pmatrix}, \qquad h[2^k] = \rho \begin{pmatrix} -1 \\ 2^{3-2^k} \\ 2^{-2^k} \\ 2^{1-2^k} - 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 2^{3-2^k} \\ 2^{-2^k} \\ 0 \\ 1 \end{pmatrix}. \tag{16}$$

Next, we show that

$$h[2^k + t] = \begin{pmatrix} 0 \\ 2^{t+3-2^k} \\ 2^{-2^k-t} \\ 0 \\ 1 \end{pmatrix}, \qquad \text{for all } t \in \{0, 1, \ldots, 2^k - 3\}, \tag{17}$$

by induction over $t$. The base case, $t = 0$, was already established in (16). For the induction step, assume that (17) holds for some $t \in \{0, 1, \ldots, 2^k - 4\}$ and compute

$$h[2^k + t + 1] = \rho \begin{pmatrix} 2^{t+4-2^k} - 1 \\ 2^{t+4-2^k} \\ 2^{-2^k-t-1} \\ 2^{t+3-2^k} + 2^{-2^k-t} - 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 2^{t+4-2^k} \\ 2^{-2^k-(t+1)} \\ 0 \\ 1 \end{pmatrix},$$

where we used $2^{t+4-2^k} \leq 1$ thanks to $t \leq 2^k - 4$. This completes the induction over $t$ and thus establishes (17). Particularizing (17) to $t = 2^k - 3$, we have

$$h[2^{k+1} - 3] = \begin{pmatrix} 0 \\ 1 \\ 2^{3-2^{k+1}} \\ 0 \\ 1 \end{pmatrix}.$$

Next, we compute

$$h[2^{k+1} - 2] = \rho \begin{pmatrix} 1 \\ 2 \\ 2^{2-2^{k+1}} \\ 1 + 2^{3-2^{k+1}} - 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 2^{2-2^{k+1}} \\ 0 \\ 1 \end{pmatrix},$$

where we used $2^{3-2^{k+1}} \leq 1$, for $k \geq 2$. This completes the induction over $k$ and thus establishes (15) and, in particular, that

$$(h[t])_1 = 1, \quad \text{if } t = 2^k - 2 \text{ for } k \in \mathbb{N}, \; k \geq 2.$$

Inspection of (16) and (17) reveals that, for all other $t$, $(h[t])_1 = 0$. Finally, $\|h[t]\|_\infty \leq 2$, for all $t \in \mathbb{N}_0$, follows from (15) through (17). $\qquad\square$

We next define a map, which, based on Lemma 13, will then be shown to produce an RNN that realizes the desired behavior.

**Definition 14.** *Let $D_1, D_h > 0$, and let $\mathcal{R}^f$ and $\mathcal{R}^g$ be RNNs such that $\mathcal{M}_{\mathrm{in}}(\mathcal{R}^g) = \mathcal{M}_{\mathrm{out}}(\mathcal{R}^f) =: d'_f$. We identify the weights of $\mathcal{R}^f$ as $A^f_x, A^f_h, b^f_h, A^f_o, b^f_o$ and the weights of $\mathcal{R}^g$ as $A^g_x, A^g_h, b^g_h, A^g_o, b^g_o$, further let $m_g := \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^g)$, $m_f := \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^f)$, and $d_f := \mathcal{M}_{\mathrm{in}}(\mathcal{R}^f)$. We define the mapping $\mathcal{R} = \mathring{\Psi}_{D_1, D_h}(\mathcal{R}^g, \mathcal{R}^f)$, with weights given by*

$$
A_h = \left( \begin{array}{c|c|c|c|c}
A^f_h & 0 & 0 & 0 & 0 \\ \hline
A^f_o & 0 & 0 & 0 & D_1 \mathbf{1}_{d'_f} \widehat{A}_o \\ \hline
-A^f_o & 0 & 0 & 0 & D_1 \mathbf{1}_{d'_f} \widehat{A}_o \\ \hline
0 & A^g_x & -A^g_x & A^g_h & -D_h \mathbf{1}_{m_g} \widehat{A}_o \\ \hline
0 & 0 & 0 & 0 & \widehat{A}
\end{array} \right), \quad
b_h = \left( \begin{array}{c}
b^f_h \\ \hline
b^f_o - D_1 \mathbf{1}_{d'_f} \\ \hline
-b^f_o - D_1 \mathbf{1}_{d'_f} \\ \hline
b^g_h \\ \hline
\widehat{b}
\end{array} \right), \quad
A_x = \left( \begin{array}{c}
A^f_x \\ \hline
0 \\ \hline
0 \\ \hline
0 \\ \hline
0
\end{array} \right),
\tag{18}
$$

$$
A_o = \left( 0 \ \middle| \ 0 \ \ 0 \ \middle| \ A^g_o \ \middle| \ 0 \right), \quad \text{and} \quad b_o = b^g_o, \quad \text{where} \ \widehat{A}_o := \left( 1 \ \ 0 \ \ 0 \ \ 0 \ \ 0 \right),
$$

*and $\widehat{A}$ and $\widehat{b}$ are as in (14).*

It holds that

$$
\mathcal{M}_{\mathrm{hid}}(\mathcal{R}) = \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^f) + 2\mathcal{M}_{\mathrm{out}}(\mathcal{R}^f) + \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^g) + 5.
$$

*Furthermore, define the mappings*

$$
\mathring{\mathcal{M}}(\mathcal{R}^g, \mathcal{R}^f) := \left( \mathbb{I}_{m_f} \ \middle| \ 0 \ \ 0 \ \middle| \ 0 \ \middle| \ 0 \right), \quad \text{and} \quad \mathring{\mathcal{W}}(\mathcal{R}^g, \mathcal{R}^f) := \left( 0 \ \middle| \ 0 \ \ 0 \ \middle| \ \mathbb{I}_{m_g} \ \middle| \ 0 \right).
$$

We next establish the properties of the RNN produced by the map $\mathring{\Psi}_{D_1, D_h}(\mathcal{R}^g, \mathcal{R}^f)$.

**Theorem 15.** *Let $D, D_2 > 0$, and let $D_1, D_h > 0$, and $\mathcal{R}^f = \mathcal{Q}^f \mathcal{K}^f$, $\mathcal{R}^g = \mathcal{Q}^g \mathcal{K}^g$ be as in Definition 14. Assume that*

*(A1)* $\left\| (\mathcal{R}^f \mathcal{D} x)[t] \right\|_\infty \le D_1$, *for $x \in [-D, D]^{d_f}$, $t \in \mathbb{N}_0$.*

*(A2)* $\| (\mathcal{K}^g \mathcal{D} y)[t] \|_\infty \le D_h$, *for $y \in [-D_1, D_1]^{d'_f}$, $t \in \mathbb{N}_0$.*

*Then, the RNN $\mathring{\Psi}_{D_1, D_h}(\mathcal{R}^g, \mathcal{R}^f) =: \mathcal{R} = \mathcal{Q}\mathcal{K}$ specified in Definition 14 satisfies, for all $x \in [-D, D]^{d_f}$,*

$$
(\mathcal{R}\mathcal{D}x)[2^k - 2] = \left( \mathcal{R}^g \mathcal{D}\left( (\mathcal{R}^f \mathcal{D}x)[2^{k-1} - 2] \right) \right)[2^{k-1} - 2], \quad \text{for all } k \in \mathbb{N}, \text{ with } k \ge 3.
\tag{19}
$$

*Furthermore, with $M = \mathring{\mathcal{M}}(\mathcal{R}^g, \mathcal{R}^f)$ and $W = \mathring{\mathcal{W}}(\mathcal{R}^g, \mathcal{R}^f)$, we have*

$$
\begin{align}
M(\mathcal{K}\mathcal{D}x)[t] &= (\mathcal{K}^f \mathcal{D}x)[t], & \text{for all } t \in \mathbb{N}_0, \tag{20} \\
W(\mathcal{K}\mathcal{D}x)[2^k - 2] &= \left( \mathcal{K}^g \mathcal{D}\left( (\mathcal{R}^f \mathcal{D}x)[2^{k-1} - 2] \right) \right)[2^{k-1} - 2], & \text{for all } k \in \mathbb{N}, \text{ with } k \ge 3. \tag{21}
\end{align}
$$

*Additionally assuming*

*(A3)* $\left\| (\mathcal{K}^f \mathcal{D}x)[t] \right\|_\infty \le D_h$, *for $x \in [-D, D]^{d_f}$, $t \in \mathbb{N}_0$,*

*(A4)* $\| (\mathcal{R}^g \mathcal{D} y)[t] \|_\infty \le D_2$, *for $y \in [-D_1, D_1]^{d'_f}$, $t \in \mathbb{N}_0$,*

*we have, for all $x \in [-D, D]^{d_f}$,*

*(B1)* $\| (\mathcal{K}\mathcal{D}x)[t] \|_\infty \le \max\{2, D_1, D_h\}$, *for all $t \in \mathbb{N}_0$,*

*(B2)* $\| (\mathcal{R}\mathcal{D}x)[t] \|_\infty \le D_2$, *for all $t \in \mathbb{N}_0$.*

*Proof.* Arbitrarily fix $x \in [-D, D]^{d_f}$ and consider the hidden state sequence in response to the input sequence $(\mathcal{D}x)[\cdot]$, i.e., $h[t] = (\mathcal{K}\mathcal{D}x)[t]$. We divide this hidden state sequence corresponding to the blocks in (18) according to

$$
h[t] = \left( \begin{array}{c}
h_1[t] \\ \hline
h_2[t] \\ \hline
h_3[t] \\ \hline
h_4[t]
\end{array} \right) = \rho \left( A_h h[t-1] + A_x((\mathcal{D}x)[t]) + b_h \right),
$$

and analyze each block separately. First, note that $h_4[\cdot]$ follows

$$h_4[t] = \rho(\widehat{A}h_4[t-1] + \widehat{b}), \; t \in \mathbb{N}_0, \; \text{with } h_4[-1] = 0. \tag{22}$$

This recursion implements the clocking mechanism defined in Lemma 13, generating the control pulses that govern the alternation between the $\mathcal{R}^f$ and $\mathcal{R}^g$ computations. By Lemma 13, we have, for all $t \in \mathbb{N}_0$,

$$\widehat{A}_o h_4[t] = \widehat{\delta}[t+2]. \tag{23}$$

Next, we note that $h_1[t]$ follows the recursion

$$h_1[t] = \rho\left(A_h^f h_1[t-1] + A_x^f(\mathcal{D}x)[t] + b_h^f\right), \; t \in \mathbb{N}_0, \quad \text{with } h_1[-1] = 0.$$

Thus, recalling Definition 1, we have

$$h_1[t] = (\mathcal{K}^f \mathcal{D}x)[t]. \tag{24}$$

Next, we consider the sequence $h_2[t]$, which is given by

$$
\begin{aligned}
h_2[t] &= \rho\left(\begin{array}{c} A_o^f h_1[t-1] + D_1 \mathbf{1}_{d'_f}\widehat{A}_o h_4[t-1] + b_o^f - D_1 \mathbf{1}_{d'_f} \\ -A_o^f h_1[t-1] + D_1 \mathbf{1}_{d'_f}\widehat{A}_o h_4[t-1] - b_o^f - D_1 \mathbf{1}_{d'_f} \end{array}\right) \\
&= \rho\left(\begin{array}{c} (A_o^f(\mathcal{K}^f\mathcal{D}x)[t-1] + b_o^f) - D_1 \mathbf{1}_{d'_f}(1 - \widehat{A}_o h_4[t-1]) \\ -(A_o^f(\mathcal{K}^f\mathcal{D}x)[t-1] + b_o^f) - D_1 \mathbf{1}_{d'_f}(1 - \widehat{A}_o h_4[t-1]) \end{array}\right) \\
&\overset{(23)}{=} \rho\left(\begin{array}{c} (\mathcal{Q}^f\mathcal{K}^f\mathcal{D}x)[t-1] - D_1 \mathbf{1}_{d'_f}(1 - \widehat{\delta}[t+1]) \\ -(\mathcal{Q}^f\mathcal{K}^f\mathcal{D}x)[t-1] - D_1 \mathbf{1}_{d'_f}(1 - \widehat{\delta}[t+1]) \end{array}\right) \\
&\overset{(i)}{=} \rho\left(\begin{array}{c} (\mathcal{R}^f\mathcal{D}x)[t-1] \\ -(\mathcal{R}^f\mathcal{D}x)[t-1] \end{array}\right)\widehat{\delta}[t+1],
\end{aligned}
\tag{25}
$$

where in (i) we used $\left\|(\mathcal{Q}^f\mathcal{K}^f\mathcal{D}x)[t-1]\right\|_\infty = \left\|(\mathcal{R}^f\mathcal{D}x)[t-1]\right\|_\infty \leq D_1$, for all $t \in \mathbb{N}_0$, by assumption (A1). From (25) it now follows that $h_3[t], t \in \mathbb{N}_0$, is given by the following recursion, with $h_3[-1] = 0$,

$$
\begin{aligned}
h_3[t] &= \rho\Big( A_x^g \rho((\mathcal{R}^f\mathcal{D}x)[t-2])\widehat{\delta}[t] - A_x^g \rho(-(\mathcal{R}^f\mathcal{D}x)[t-2])\widehat{\delta}[t] \\
&\qquad + A_h^g h_3[t-1] + b_h^g - D_h \mathbf{1}_{m_g}\widehat{A}_o h_4[t-1] \Big) \\
&\overset{(ii)}{=} \rho\left( A_h^g h_3[t-1] + b_h^g - D_h \mathbf{1}_{m_g}\widehat{A}_o h_4[t-1] + \widehat{\delta}[t] \cdot A_x^g(\mathcal{R}^f\mathcal{D}x)[t-2] \right) \\
&\overset{(23)}{=} \rho\left( A_h^g h_3[t-1] + b_h^g - D_h \mathbf{1}_{m_g}\widehat{\delta}[t+1] + \widehat{\delta}[t] \cdot A_x^g(\mathcal{R}^f\mathcal{D}x)[t-2] \right),
\end{aligned}
$$

where in (ii) we used the identity $x = \rho(x) - \rho(-x)$. As $\widehat{\delta}[\ell] = 0$, for $\ell \in \{0, \dots, 3\}$, it follows that

$$h_3[t] = (\mathcal{K}^g\mathcal{D}0)[t], \quad \text{for } t \in \{0, 1, 2\}.$$

Next, we compute

$$
\begin{aligned}
h_3[3] &= \rho\left( A_h^g h_3[2] + b_h^g - D_h \mathbf{1}_{m_g}\widehat{\delta}[4] + \widehat{\delta}[3] \cdot A_x^g(\mathcal{R}^f\mathcal{D}x)[1] \right) \\
&= \rho\left( A_h^g h_3[2] + b_h^g - D_h \mathbf{1}_{m_g} \right) = 0,
\end{aligned}
\tag{26}
$$

where we used $\|(\mathcal{K}^g\mathcal{D}0)[3]\|_\infty = \|\rho(A_h^g h_3[2] + b_h^g)\|_\infty \leq D_h$, thanks to Assumption (A2). We now prove by nested induction that

$$h_3[2^k - 1] = 0, \qquad\qquad\qquad\qquad\quad \text{for } k \in \mathbb{N}, \; k \geq 2, \tag{27}$$
$$\text{and} \quad h_3[2^k + \ell] = \left(\mathcal{K}^g\mathcal{D}\left((\mathcal{R}^f\mathcal{D}x)[2^k - 2]\right)\right)[\ell], \qquad \text{for } \ell \in \{0, \dots, 2^k - 2\}. \tag{28}$$

The base case for the induction over $k$, i.e., (27) for $k = 2$, was already established in (26). Next, we assume that (27) holds for some $k \geq 2$ and compute

$$h_3[2^k] = \rho\left(A_h^g h_3[2^k - 1] + b_h^g - D_h \mathbf{1}_{m_g}\widehat{\delta}[2^k + 1] + \widehat{\delta}[2^k] \cdot A_x^g (\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)$$

$$= \rho\left(b_h^g + A_x^g (\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)$$

$$= \left(\mathcal{K}^g \mathcal{D}\left((\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)\right)[0],$$

where in the last step we used Remark 5. This establishes the base case $\ell = 0$ for the induction over $\ell$ in (28). Next, assume that (28) holds for some $\ell \in \{0, \ldots, 2^k - 3\}$ and compute

$$h_3[2^k + \ell + 1] = \rho\Bigg(A_h^g h_3[2^k + \ell] + b_h^g - D_h \mathbf{1}_{m_g}\widehat{\delta}[2^k + \ell + 2]$$

$$+ \widehat{\delta}[2^k + \ell + 1] \cdot A_x^g (\mathcal{R}^f \mathcal{D}x)[2^k + \ell - 1]\Bigg)$$

$$= \rho\left(A_h^g h_3[2^k + \ell] + b_h^g\right)$$

$$\overset{(28)}{=} \rho\left(A_h^g \left(\mathcal{K}^g \mathcal{D}\left((\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)\right)[\ell] + b_h^g\right)$$

$$\overset{\text{Rem. } 5}{=} \left(\mathcal{K}^g \mathcal{D}\left((\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)\right)[\ell + 1].$$

This completes the induction over $\ell$ and thus establishes (28). We now proceed by noting that

$$h_3[2^{k+1} - 1] = \rho\Bigg(A_h^g h_3[2^k + 2^k - 2] + b_h^g - D_h \mathbf{1}_{m_g}\widehat{\delta}[2^{k+1}]$$

$$+ \widehat{\delta}[2^{k+1} - 1] \cdot A_x^g (\mathcal{R}^f \mathcal{D}x)[2^{k+1} - 3]\Bigg)$$

$$\overset{(iii)}{=} \rho\left(A_h^g \left(\mathcal{K}^g \mathcal{D}\left((\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)\right)[2^k - 2] + b_h^g - D_h \mathbf{1}_{m_g}\right)$$

$$\overset{(iv)}{=} 0,$$

where (iii) follows from (28) with $\ell = 2^k - 2$ and (iv) holds because

$$A_h^g \left(\mathcal{K}^g \mathcal{D}\left((\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)\right)[2^k - 2] + b_h^g \preccurlyeq \|\rho(A_h^g \left(\mathcal{K}^g \mathcal{D}\left((\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)\right)[2^k - 2] + b_h^g)\|_\infty$$

$$\overset{(A1)}{\leq} \max_{\|x'\|_\infty \leq D_1} \|\rho(A_h^g \left(\mathcal{K}^g \mathcal{D}x'\right)[2^k - 2] + b_h^g)\|_\infty$$

$$\overset{\text{Rem. } 5}{=} \max_{\|x'\|_\infty \leq D_1} \|(\mathcal{K}^g \mathcal{D}x')[2^k - 1]\|_\infty \overset{(A2)}{\leq} D_h.$$

This establishes that (27) holds for $k + 1$ as well and thus completes the nested induction. The overall network output is given by $(\mathcal{R}\mathcal{D}x)[t] = A_o^g h_3[t] + b_o^g$, which, using (28), yields

$$(\mathcal{R}\mathcal{D}x)[2^k + \ell] = \left(\mathcal{R}^g \mathcal{D}\left((\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)\right)[\ell], \quad \text{for all } k \in \mathbb{N}, \text{ with } k \geq 2, \text{ and all } \ell \in \{0, \ldots, 2^k - 2\}. \tag{29}$$

In particular, setting $\ell = 2^k - 2$, we obtain

$$(\mathcal{R}\mathcal{D}x)[2^{k+1} - 2] = \left(\mathcal{R}^g \mathcal{D}\left((\mathcal{R}^f \mathcal{D}x)[2^k - 2]\right)\right)[2^k - 2],$$

which establishes (19). Moreover, (20) is an immediate consequence of Definition 14 and (24). Finally, (21) follows from Definition 14 and (28) with $\ell = 2^k - 2$.

To establish (B1), we note that $\|h_1[t]\|_\infty \leq D_h$ by (24) and Assumption (A3), $\|h_2[t]\|_\infty \leq D_1$ by (25) and Assumption (A1), $\|h_3[t]\|_\infty \leq D_h$ by (27) and (28) together with Assumptions (A1) and (A2), and finally $\|h_4[t]\|_\infty \leq 2$ by (22) and Lemma 13. Taking the maximum over these bounds, we arrive at (B1). The proof is concluded upon noting that (B2) follows from (29) together with Assumption (A4). $\qquad\square$

We now illustrate, through a simple example, the concatenation mechanism developed in Theorem 15 and prepare for its generalization in Lemma 16. Consider the function $x \mapsto x^{16}$, which can be expressed as Sq $\circ$ Sq $\circ$ Sq $\circ$ Sq, where Sq $: x \mapsto x^2$. Accordingly, $x^{16}$ can be approximated by concatenating four copies of $\mathcal{R}^{\mathrm{Sq}}$, the RNN introduced in Theorem 8. Using the map from Definition 14, we construct the network

$$\mathring{\Psi}(\mathring{\Psi}(\mathcal{R}^{\mathrm{Sq}}, \mathcal{R}^{\mathrm{Sq}}), \, \mathring{\Psi}(\mathcal{R}^{\mathrm{Sq}}, \mathcal{R}^{\mathrm{Sq}})),$$

depicted in Figure 2, whose output sequence approximates $x^{16}$. Furthermore, suitably exploiting (20) and (21) makes it possible to also recover approximations of the intermediate functions $x^2$, $x^4$, and $x^8$. The following result, Lemma 16, formalizes this idea by extending the two-network concatenation, which implements function composition, of Theorem 15 to the concatenation of multiple RNNs.
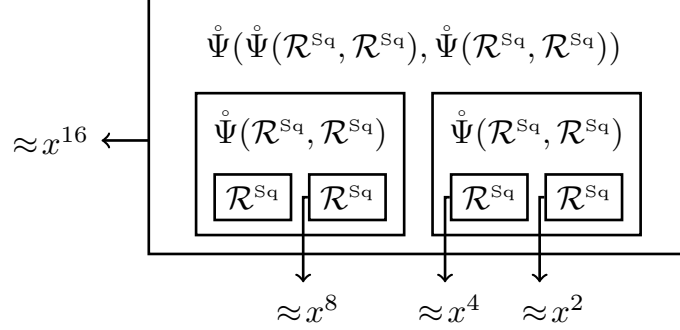


Figure 2: RNN approximating $x \to x^{16}$.

**Lemma 16.** *Let $D > 0$, $L \in \mathbb{N}$, and let $\mathcal{R}^1, \ldots, \mathcal{R}^{2^L}$ be RNNs with $\mathcal{M}_{\mathrm{in}}(\mathcal{R}^{\ell+1}) = \mathcal{M}_{\mathrm{out}}(\mathcal{R}^\ell)$, $\ell \in \{1, \ldots, 2^L - 1\}$. Assume that there are constants $D_1, \ldots, D_{2^L}$ and $D_h$ with $\max\{2, \max_{\ell \in \{1, \ldots, 2^L\}} D_\ell\} \leq D_h$ such that, for all $\ell \in \{1, \ldots, 2^L\}$,*

$$\sup_{\|x\|_\infty \leq D_{\ell-1}, t \in \mathbb{N}_0} \left\| (\mathcal{R}^\ell \mathcal{D}x)[t] \right\|_\infty \leq D_\ell, \tag{30}$$

$$\sup_{\|x\|_\infty \leq D_{\ell-1}, t \in \mathbb{N}_0} \left\| (\mathcal{K}^\ell \mathcal{D}x)[t] \right\|_\infty \leq D_h, \tag{31}$$

*where we set $D_0 = D$. Define the mappings*

$$g_\ell^k : \mathbb{R}^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^\ell)} \to \mathbb{R}^{\mathcal{M}_{\mathrm{out}}(\mathcal{R}^\ell)} \ \ as \ \ x \to \left( \mathcal{R}^\ell \mathcal{D}x \right)[2^k - 2], \qquad \ell \in \{1, \ldots, 2^L\}, \, k \geq 2. \tag{32}$$

*Then, there exists an RNN $\widehat{\mathcal{R}}^L = \widehat{\mathcal{Q}}^L \widehat{\mathcal{K}}^L$ such that*

$$\sup_{\|x\|_\infty \leq D, t \in \mathbb{N}_0} \left\| \left( \widehat{\mathcal{R}}^L \mathcal{D}x \right)[t] \right\|_\infty \leq D_{2^L}, \qquad \sup_{\|x\|_\infty \leq D, t \in \mathbb{N}_0} \left\| \left( \widehat{\mathcal{K}}^L \mathcal{D}x \right)[t] \right\|_\infty \leq D_h, \tag{33}$$

*and*

$$\left( \widehat{\mathcal{R}}^L \mathcal{D}x \right)[2^k - 2] = \left( g_{2^L}^{k-L} \circ \cdots \circ g_1^{k-L} \right)(x), \quad k \geq L + 2, \, x \in [-D, D]^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^1)}. \tag{34}$$

*Furthermore, for every $\ell \in \{1, \ldots, 2^L\}$, there are $A^\ell, b^\ell$, and $\widetilde{k}_1, \ldots, \widetilde{k}_\ell \in \{0, \ldots, \lceil \log(\ell) \rceil\}$ such that*

$$A^\ell \left( \widehat{\mathcal{K}}^L \mathcal{D}x \right)[2^k - 2] + b^\ell = \left( g_\ell^{k-\widetilde{k}_\ell} \circ \cdots \circ g_1^{k-\widetilde{k}_1} \right)(x), \qquad k \geq L + 2, \, x \in [-D, D]^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^1)}. \tag{35}$$

*Finally, we have*

$$\mathcal{M}_{\mathrm{hid}}(\widehat{\mathcal{R}}^L) = \sum_{\ell=1}^{2^L} \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^\ell) + 2 \sum_{\ell=1}^{2^L-1} \mathcal{M}_{\mathrm{out}}(\mathcal{R}^\ell) + 5(2^L - 1). \tag{36}$$

*Proof.* We prove the statement by induction over $L$. For the base case, $L = 1$, given $\mathcal{R}^1 = \mathcal{Q}^1 \mathcal{K}^1$ and $\mathcal{R}^2 = \mathcal{Q}^2 \mathcal{K}^2$, we take $\widehat{\mathcal{R}}^1 = \widehat{\mathcal{Q}}^1 \widehat{\mathcal{K}}^1 = \mathring{\Psi}_{D_h, D_h}(\mathcal{R}^2, \mathcal{R}^1)$ according to Definition 14 and $M^1 = \mathring{\mathcal{M}}(\mathcal{R}^2, \mathcal{R}^1)$, $W^1 = \mathring{\mathcal{W}}(\mathcal{R}^2, \mathcal{R}^1)$. By assumption, we have

$$\sup_{\|x\|_\infty \leq D, t \in \mathbb{N}_0} \left\| (\mathcal{R}^1 \mathcal{D} x)[t] \right\|_\infty \leq D_1 \leq D_h \text{ and } \sup_{\|x'\|_\infty \leq D_1, t \in \mathbb{N}_0} \left\| (\mathcal{K}^2 \mathcal{D} x')[t] \right\|_\infty \leq D_h, \tag{37}$$

as well as

$$\sup_{\|x\|_\infty \leq D, t \in \mathbb{N}_0} \left\| (\mathcal{K}^1 \mathcal{D} x)[t] \right\|_\infty \leq D_h \text{ and } \sup_{\|x'\|_\infty \leq D_1, t \in \mathbb{N}_0} \left\| (\mathcal{R}^2 \mathcal{D} x')[t] \right\|_\infty \leq D_2. \tag{38}$$

We now invoke Theorem 15 with $\mathcal{R}^f = \mathcal{R}^1$ and $\mathcal{R}^g = \mathcal{R}^2$ and note that its conditions (A1) and (A2) are satisfied thanks to (37), and (A3) and (A4) are met owing to (38). We thus have (33) for $L = 1$ as a consequence of (B1) and (B2) in Theorem 15. Furthermore, it follows from (19) in Theorem 15 that

$$(\widehat{\mathcal{R}}^1 \mathcal{D} x)[2^k - 2] = \left( \mathcal{R}^2 \mathcal{D} \left( (\mathcal{R}^1 \mathcal{D} x)[2^{k-1} - 2] \right) \right)[2^{k-1} - 2], \text{ for all } k \geq 3, \text{ and all } x \in [-D, D]^{\mathcal{M}_{\text{in}}(\mathcal{R}^1)},$$

which, upon invoking (32), yields (34). Next, we establish (35). To this end, let

$$A^1 := A_o^1 M^1, \ b^1 = b_o^1, \text{ and } A^2 := A_o^2 W^1, \ b^2 = b_o^2,$$

where $A_o^1, b_o^1$ and $A_o^2, b_o^2$ are the weights of the affine output mappings for the RNNs $\mathcal{R}^1$ and $\mathcal{R}^2$, respectively. We have, for all $x$ with $\|x\|_\infty \leq D$ and all $k \geq 3$,

$$\begin{aligned}
A^1 \left( \widehat{\mathcal{K}}^1 \mathcal{D} x \right) [2^k - 2] + b^1 &= A_o^1 M^1 \left( \widehat{\mathcal{K}}^1 \mathcal{D} x \right) [2^k - 2] + b_o^1 \\
&\stackrel{(20)}{=} A_o^1 \left( \mathcal{K}^1 \mathcal{D} x \right) [2^k - 2] + b_o^1 \\
&= \left( \mathcal{R}^1 \mathcal{D} x \right) [2^k - 2] \\
&\stackrel{(32)}{=} g_1^{k-0}(x).
\end{aligned}$$

Hence, this establishes (35) for $\ell = 1$, with $\widetilde{k}_1 = 0$. Furthermore, we have, for all $x$ with $\|x\|_\infty \leq D$ and all $k \geq 3$,

$$\begin{aligned}
A^2 \left( \widehat{\mathcal{K}}^1 \mathcal{D} x \right) [2^k - 2] + b^2 &= A_o^2 W^1 \left( \widehat{\mathcal{K}}^1 \mathcal{D} x \right) [2^k - 2] + b_o^2 \\
&\stackrel{(21)}{=} A_o^2 \left( \mathcal{K}^2 \mathcal{D} \left( \mathcal{R}^1 \mathcal{D} x \right) [2^{k-1} - 2] \right) [2^{k-1} - 2] + b_o^2 \\
&= \left( \mathcal{R}^2 \mathcal{D} \left( \mathcal{R}^1 \mathcal{D} x \right) [2^{k-1} - 2] \right) [2^{k-1} - 2] \\
&\stackrel{(32)}{=} \left( g_2^{k-1} \circ g_1^{k-1} \right)(x).
\end{aligned}$$

This establishes (35) for $\ell = 2$, with $\widetilde{k}_1 = \widetilde{k}_2 = 1$. Finally, by Definition 14, we obtain

$$\mathcal{M}_{\text{hid}}(\widehat{\mathcal{R}}^1) = \mathcal{M}_{\text{hid}}(\mathcal{R}^1) + 2\mathcal{M}_{\text{out}}(\mathcal{R}^1) + \mathcal{M}_{\text{hid}}(\mathcal{R}^2) + 5,$$

which yields (36) for $L = 1$. This finishes the proof of the base case $L = 1$ for the induction over $L$.

We proceed to the induction step and assume that Lemma 16 holds for some $L \geq 1$. Specifically, let $\mathcal{R}^1, \ldots, \mathcal{R}^{2^{L+1}}$ be RNNs satisfying (30) and (31) with constants $D_1, \ldots, D_{2^{L+1}}, D_h$. We first invoke Lemma 16 for $\mathcal{R}^1, \ldots, \mathcal{R}^{2^L}$ and denote the resulting overall RNN as $\widehat{\mathcal{R}}^a = \widehat{\mathcal{Q}}^a \widehat{\mathcal{K}}^a$. By (33) we have

$$\left\| \left( \widehat{\mathcal{R}}^a \mathcal{D} x \right) [t] \right\|_\infty \leq D_{2^L} \quad \text{and} \quad \left\| \left( \widehat{\mathcal{K}}^a \mathcal{D} x \right) [t] \right\|_\infty \leq D_h, \quad \text{for all } t \in \mathbb{N}_0, \text{ and all } x \in [-D, D]^{\mathcal{M}_{\text{in}}(\mathcal{R}^1)}. \tag{39}$$

Similarly, we invoke Lemma 16 for the $2^L$ RNNs $\mathcal{R}^{2^L+1}, \ldots, \mathcal{R}^{2^{L+1}}$ and denote the resulting overall RNN by $\widehat{\mathcal{R}}^b = \widehat{\mathcal{Q}}^b \widehat{\mathcal{K}}^b$, which satisfies

$$\left\| \left( \widehat{\mathcal{R}}^b \mathcal{D} x \right) [t] \right\|_\infty \leq D_{2^{L+1}} \text{ and } \left\| \left( \widehat{\mathcal{K}}^b \mathcal{D} x \right) [t] \right\|_\infty \leq D_h, \text{ for all } t \in \mathbb{N}_0, \text{ and all } x \in [-D_{2^L}, D_{2^L}]^{\mathcal{M}_{\text{in}}(\mathcal{R}^{2^L+1})}. \tag{40}$$

Now, set $\widehat{\mathcal{R}}^{L+1} = \widehat{\mathcal{Q}}^{L+1} \widehat{\mathcal{K}}^{L+1} = \mathring{\Psi}_{D_h, D_h}(\widehat{\mathcal{R}}^b, \widehat{\mathcal{R}}^a)$ and invoke Theorem 15 with the following correspondence of quantities:

| Here | Theorem 15 |
|:---:|:---:|
| $\widehat{\mathcal{R}}^a$ | $\mathcal{R}^f$ |
| $\widehat{\mathcal{R}}^b$ | $\mathcal{R}^g$ |
| $D$ | $D$ |
| $D_h$ | $D_h$ |
| $D_{2^L}$ | $D_1$ |
| $D_{2^{L+1}}$ | $D_2$ |

By (39), Conditions (A1) and (A3) of Theorem 15 are satisfied. Further, by (40), Conditions (A4) and (A2) of Theorem 15 are met. Hence, (33) holds for $\widehat{\mathcal{R}}^{L+1}$ as a consequence of (B1) and (B2) in Theorem 15. By (19) in Theorem 15 we have

$$(\widehat{\mathcal{R}}^{L+1}\mathcal{D}x)[2^k - 2] = \left(\widehat{\mathcal{R}}^b\mathcal{D}\left((\widehat{\mathcal{R}}^a\mathcal{D}x)[2^{k-1}-2]\right)\right)[2^{k-1}-2], \quad \text{for all } k \geq 3, \text{ and all } x \in [-D,D]^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^1)},$$

which, upon using (34) for $\widehat{\mathcal{R}}^b$ and $\widehat{\mathcal{R}}^a$, yields

$$(\widehat{\mathcal{R}}^{L+1}\mathcal{D}x)[2^k - 2] = \left(g_{2^{L+1}}^{k-(L+1)} \circ \cdots \circ g_{2^L+1}^{k-(L+1)}\right)\left(\left(g_{2^L}^{k-(L+1)} \circ \cdots \circ g_1^{k-(L+1)}\right)(x)\right)$$
$$= \left(g_{2^{L+1}}^{k-(L+1)} \circ \cdots \circ g_1^{k-(L+1)}\right)(x), \text{ for all } k \geq L+3, \text{ and all } x \in [-D,D]^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^1)}.$$

This proves (34) for $\widehat{\mathcal{R}}^{L+1}$.

Next, we establish (35) for $\widehat{\mathcal{R}}^{L+1}$ and first treat the case $\ell \in \{1,\ldots,2^L\}$. To this end arbitrarily fix $\ell \in \{1,\ldots,2^L\}$ and note that, using (35) for $\widehat{\mathcal{R}}^a$, which is possible by the induction assumption, we can conclude that there are $A_a^\ell, b_a^\ell$, and $\widetilde{k}_1^a, \ldots, \widetilde{k}_\ell^a \in \{0,\ldots,\lceil\log(\ell)\rceil\}$ such that

$$A_a^\ell\left(\widehat{\mathcal{K}}^a\mathcal{D}x\right)[2^k-2]+b_a^\ell = \left(g_\ell^{k-\widetilde{k}_\ell^a} \circ \cdots \circ g_1^{k-\widetilde{k}_1^a}\right)(x), \qquad \text{for all } k \geq L+2, \text{ and all } x \in [-D,D]^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^1)}. \tag{41}$$

Now, we let

$$A^\ell := A_a^\ell M, \quad b^\ell := b_a^\ell, \quad \text{and} \quad \widetilde{k}_1 := \widetilde{k}_1^a, \ldots, \widetilde{k}_\ell := \widetilde{k}_\ell^a,$$

with $M = \mathring{\mathcal{M}}(\widehat{\mathcal{R}}^b, \widehat{\mathcal{R}}^a)$, and compute

$$A^\ell\left(\widehat{\mathcal{K}}^{L+1}\mathcal{D}x\right)[2^k-2]+b^\ell = A_a^\ell M\left(\widehat{\mathcal{K}}^{L+1}\mathcal{D}x\right)[2^k-2]+b_a^\ell$$
$$\overset{(20)}{=} A_a^\ell\left(\widehat{\mathcal{K}}^a\mathcal{D}x\right)[2^k-2]+b_a^\ell$$
$$\overset{(41)}{=} \left(g_\ell^{k-\widetilde{k}_\ell^a} \circ \cdots \circ g_1^{k-\widetilde{k}_1^a}\right)(x)$$
$$= \left(g_\ell^{k-\widetilde{k}_\ell} \circ \cdots \circ g_1^{k-\widetilde{k}_1}\right)(x),$$

for all $k \geq L+2$ and all $x \in [-D,D]^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^1)}$. In particular, the same identity holds for all $k \geq L+3$, so (35) is established for $\ell \in \{1,\ldots,2^L\}$. It remains to treat the case $\ell \in \{2^L+1,\ldots,2^{L+1}\}$. To this end, we arbitrarily fix $\ell \in \{2^L+1,\ldots,2^{L+1}\}$ and set $\ell' = \ell - 2^L \in \{1,\ldots,2^L\}$. Now, upon application of (35) to $\widehat{\mathcal{R}}^b$, we can conclude the existence of $A_b^{\ell'}, b_b^{\ell'}$, and

$$\widetilde{k}_1^b, \ldots, \widetilde{k}_{\ell'}^b \in \{0,\ldots,\lceil\log(\ell')\rceil\} \tag{42}$$

such that

$$A_b^{\ell'}\left(\widehat{\mathcal{K}}^b\mathcal{D}x\right)[2^k-2]+b_b^{\ell'} = \left(g_{2^L+\ell'}^{k-\widetilde{k}_{\ell'}^b} \circ \cdots \circ g_{2^L+1}^{k-\widetilde{k}_1^b}\right)(x), \tag{43}$$

for all $k \geq L+2$, and all $x \in [-D_{2^L}, D_{2^L}]^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^{2^L+1})}$. Next, with $W = \mathring{\mathcal{W}}(\widehat{\mathcal{R}}^b, \widehat{\mathcal{R}}^a)$, let

$$A^\ell := A_b^{\ell'}W, \quad b^\ell := b_b^{\ell'}, \quad \text{and}$$
$$\widetilde{k}_1 := L+1, \ldots, \widetilde{k}_{2^L} := L+1, \widetilde{k}_{2^L+1} := \widetilde{k}_1^b+1, \ldots, \widetilde{k}_\ell := \widetilde{k}_{\ell'}^b+1, \tag{44}$$

19

and note that by (42), $\widetilde{k}_1, \ldots, \widetilde{k}_\ell \in \{0, \ldots, L+1\}$ because $\ell' = \ell - 2^L \leq 2^L$. As $\ell \in \{2^L + 1, \ldots, 2^{L+1}\}$, we have $\lceil \log(\ell) \rceil = L + 1$ and can thus equivalently write $\widetilde{k}_1, \ldots, \widetilde{k}_\ell \in \{0, \ldots, \lceil \log(\ell) \rceil\}$. Next, we compute

$$
\begin{aligned}
A^\ell \left( \widehat{\mathcal{K}}^{L+1} \mathcal{D}x \right) [2^k - 2] + b^\ell &= A_b^{\ell'} W \left( \widehat{\mathcal{K}}^{L+1} \mathcal{D}x \right) [2^k - 2] + b_b^{\ell'} \\
&\overset{(21)}{=} A_b^{\ell'} \left( \widehat{\mathcal{K}}^b \mathcal{D} \left( \widehat{\mathcal{R}}^a \mathcal{D}x \right) [2^{k-1} - 2] \right) [2^{k-1} - 2] + b_b^{\ell'} \\
&\overset{(43)}{=} \left( g_\ell^{k-1-\widetilde{k}_{\ell'}^b} \circ \cdots \circ g_{2^L+1}^{k-1-\widetilde{k}_1^b} \right) \left( \left( \widehat{\mathcal{R}}^a \mathcal{D}x \right) [2^{k-1} - 2] \right) \\
&\overset{(34) \text{ for } \widehat{\mathcal{R}}^a}{=} \left( g_\ell^{k-1-\widetilde{k}_{\ell'}^b} \circ \cdots \circ g_{2^L+1}^{k-1-\widetilde{k}_1^b} \right) \left( \left( g_{2^L}^{k-1-L} \circ \cdots \circ g_1^{k-1-L} \right) (x) \right), \\
&\overset{(44)}{=} \left( g_\ell^{k-\widetilde{k}_\ell} \circ \cdots \circ g_1^{k-\widetilde{k}_1} \right) (x),
\end{aligned}
$$

for all $k \in \mathbb{N}$ such that $k - 1 \geq L + 2$, i.e., $k \geq L + 3$, and all $x \in [-D, D]^{\mathcal{M}_{\text{in}}(\mathcal{R}^1)}$. This establishes (35) for the case $\ell \in \{2^L + 1, \ldots, 2^{L+1}\}$ as well.

The proof is concluded by establishing (36) as follows

$$
\begin{aligned}
\mathcal{M}_{\text{hid}}(\widehat{\mathcal{R}}^{L+1}) &= \mathcal{M}_{\text{hid}}(\widehat{\mathcal{R}}^a) + 2\mathcal{M}_{\text{out}}(\widehat{\mathcal{R}}^a) + \mathcal{M}_{\text{hid}}(\widehat{\mathcal{R}}^b) + 5 \\
&= \left( \sum_{\ell=1}^{2^L} \mathcal{M}_{\text{hid}}(\mathcal{R}^\ell) + 2 \sum_{\ell=1}^{2^L - 1} \mathcal{M}_{\text{out}}(\mathcal{R}^\ell) + 5(2^L - 1) \right) \\
&\quad + 2\mathcal{M}_{\text{out}}(\mathcal{R}^{2^L}) + 5 \\
&\quad + \left( \sum_{\ell=1}^{2^L} \mathcal{M}_{\text{hid}}(\mathcal{R}^{2^L+\ell}) + 2 \sum_{\ell=1}^{2^L - 1} \mathcal{M}_{\text{out}}(\mathcal{R}^{2^L+\ell}) + 5(2^L - 1) \right) \\
&= \sum_{\ell=1}^{2^{L+1}} \mathcal{M}_{\text{hid}}(\mathcal{R}^\ell) + 2 \sum_{\ell=1}^{2^{L+1} - 1} \mathcal{M}_{\text{out}}(\mathcal{R}^\ell) + 5 \left( (2^L - 1) + (2^L - 1) + 1 \right).
\end{aligned}
$$

This finalizes the induction step going from $\widehat{\mathcal{R}}^L$ to $\widehat{\mathcal{R}}^{L+1}$ and thereby completes the overall proof. $\quad\square$

We finally extend Lemma 16 to the concatenation of an arbitrary number—as opposed to a power of two—RNNs. This will be effected by suitably inserting dummy networks and is formalized as follows.

**Corollary 17.** *Fix $D \geq 1$, $L \in \mathbb{N}$, and let $\mathcal{R}^1, \ldots, \mathcal{R}^L$ be RNNs with $\mathcal{M}_{\text{in}}(\mathcal{R}^{\ell+1}) = \mathcal{M}_{\text{out}}(\mathcal{R}^\ell)$. Furthermore, assume that there are constants $D_1, \ldots, D_L$ and $D_h \geq \max\{2, \max_{\ell \in \{1,\ldots,L\}} D_\ell\}$ such that, for all $\ell \in \{1, \ldots, L\}$,*

$$
\sup_{\|x\|_\infty \leq D_{\ell-1}, t \in \mathbb{N}_0} \left\| (\mathcal{R}^\ell \mathcal{D}x)[t] \right\|_\infty \leq D_\ell,
$$

$$
\sup_{\|x\|_\infty \leq D_{\ell-1}, t \in \mathbb{N}_0} \left\| (\mathcal{K}^\ell \mathcal{D}x)[t] \right\|_\infty \leq D_h,
$$

*where we set $D_0 = D$. Then, there exists a hidden state operator (Definition 1) $\widehat{\mathcal{K}}$ with*

$$
\mathcal{M}_{\text{hid}}(\widehat{\mathcal{K}}) \leq \sum_{\ell=1}^{L} \mathcal{M}_{\text{hid}}(\mathcal{R}^\ell) + 2 \sum_{\ell=1}^{L} \mathcal{M}_{\text{out}}(\mathcal{R}^\ell) + 13L, \tag{45}
$$

*such that, for every $\ell \in \{1, \ldots, L\}$, there are $A^\ell, b^\ell$, and $\widetilde{k}_1, \ldots, \widetilde{k}_\ell \in \{1, \ldots, \lceil \log(\ell) \rceil\}$, so that*

$$
A^\ell \left( \widehat{\mathcal{K}} \mathcal{D}x \right) [2^k - 2] + b^\ell = \left( g_\ell^{k-\widetilde{k}_\ell} \circ \cdots \circ g_1^{k-\widetilde{k}_1} \right) (x), \text{ for all } k \geq \lceil \log(L) \rceil + 2, \text{ and all } x \in [-D, D]^{\mathcal{M}_{\text{in}}(\mathcal{R}^1)}, \tag{46}
$$

*with $g_\ell^k$ as defined in (32).*

*Proof.* With the goal of applying Lemma 16, we complete the specified collection $\mathcal{R}^1, \ldots, \mathcal{R}^L$ of RNNs by the dummy RNN $\mathcal{R}_d^{\circ} = \mathcal{Q}_d^{\circ}\mathcal{K}_d^{\circ}$ for input dimension $d \in \mathbb{N}$ with weights

$$A_h = \begin{pmatrix} 0 \end{pmatrix} \qquad A_x = \begin{pmatrix} 0\mathbf{1}_d^T \end{pmatrix} \qquad b_h = \begin{pmatrix} 0 \end{pmatrix} \qquad A_o = \begin{pmatrix} 0 \end{pmatrix} \qquad b_o = \begin{pmatrix} 0 \end{pmatrix}.$$

Since $(\mathcal{R}_d^{\circ}\mathcal{D}x)[t] = 0$ and $(\mathcal{K}_d^{\circ}\mathcal{D}x)[t] = 0$, for all $x \in \mathbb{R}^d$, the conditions (30) and (31) in Lemma 16 are trivially satisfied. Next, we let $L' := \lceil \log(L) \rceil$, invoke Lemma 16 for

$$\mathcal{R}^1, \ldots, \mathcal{R}^L, \underbrace{\mathcal{R}_{\mathcal{M}_{\mathrm{out}}(\mathcal{R}^L)}^{\circ}, \mathcal{R}_1^{\circ}, \ldots, \mathcal{R}_1^{\circ},}_{2^{L'}-L \text{ dummy networks}}$$

and denote the resulting network by $\widehat{\mathcal{R}} = \widehat{\mathcal{Q}}\widehat{\mathcal{K}}$. Now $\widehat{\mathcal{K}}$ is the desired hidden state operator since, by (35), for every $\ell \in \{1, \ldots, 2^{L'}\}$, there are $A^\ell, b^\ell$, and $\widetilde{k}_1, \ldots, \widetilde{k}_\ell \in \{0, \ldots, \lceil \log(\ell) \rceil\}$ such that

$$A^\ell \left( \widehat{\mathcal{K}}\mathcal{D}x \right) [2^k - 2] + b^\ell = \left( g_\ell^{k-\widetilde{k}_\ell} \circ \cdots \circ g_1^{k-\widetilde{k}_1} \right)(x), \quad \text{for all } k \geq L'+2, \text{ and all } x \in [-D, D]^{\mathcal{M}_{\mathrm{in}}(\mathcal{R}^1)}.$$

Restricting to $\ell \in \{1, \ldots, L\}$ yields (46).

It remains to establish that the operator $\widehat{\mathcal{K}}$ we have identified satisfies (45). When $L = 2^{L'}$ the statement follows straight from Lemma 16. Else, we have

$$\mathcal{M}_{\mathrm{hid}}(\widehat{\mathcal{K}}) \overset{(36)}{=} \sum_{\ell=1}^{L} \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^\ell) + (2^{L'} - L)\mathcal{M}_{\mathrm{hid}}(\mathcal{R}^{\circ}) + 2 \left( \sum_{\ell=1}^{L} \mathcal{M}_{\mathrm{out}}(\mathcal{R}^\ell) + (2^{L'} - L - 1)\mathcal{M}_{\mathrm{out}}(\mathcal{R}^{\circ}) \right)$$
$$+ 5(2^{L'} - 1)$$
$$\overset{(i)}{\leq} \sum_{\ell=1}^{L} \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^\ell) + L + 2 \left( \sum_{\ell=1}^{L} \mathcal{M}_{\mathrm{out}}(\mathcal{R}^\ell) + (L - 1) \right) + 5(2L - 1)$$
$$\leq \sum_{\ell=1}^{L} \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^\ell) + 2 \sum_{\ell=1}^{L} \mathcal{M}_{\mathrm{out}}(\mathcal{R}^\ell) + 13L,$$

where in (i) we used $2^{L'} \leq 2L$. $\qquad\square$

# 4  Approximation of monomials

This section develops the hierarchical constructions used to approximate higher-order monomials. In deep feed-forward ReLU networks, such hierarchies are formed by successive applications of squaring and multiplication. In our RNN framework, temporal depth plays the same role: higher powers are obtained through repeated use of the clocked concatenation mechanism, combined with parallelization and affine transformations as needed. This temporal realization of compositional structure is key to constructing RNNs that efficiently approximate the vector $(x, x^2, \ldots, x^N)$.

Our goal is to build a single RNN that, when run for sufficiently many time steps, simultaneously produces accurate approximations of $x^2, x^3, \ldots, x^N$, with $N = 2^L$, $L \in \mathbb{N}$. To this end, we combine the previously established RNNs for squaring (Theorem 8) and multiplication (Theorem 12) through the concatenation procedure formalized in Corollary 17. The resulting hierarchical construction, illustrated in Figure 3, organizes the powers in a pyramid-like structure, enabling efficient realization of all monomials up to degree $N = 2^L$ using $L = \log(N)$ successive concatenations. The restriction to $N = 2^L$ poses no difficulty, since unused monomials will later be discarded by assigning them zero coefficients. The $\ell$-th row in Figure 3 corresponds to the application of the function $f^\ell$ as defined next.

**Definition 18.** *We define the mappings $f^1 : \mathbb{R} \to \mathbb{R}^2$ and, for $\ell \in \mathbb{N}$, with $\ell \geq 2$, $f^\ell : \mathbb{R}^{2^{\ell-2}+1} \to \mathbb{R}^{2^{\ell-1}+1}$, as follows*

$$f^1 : x \quad \to \quad \left( x^2 \quad x \right)^T$$
$$f^\ell(x)_i = \begin{cases} x_{2^{\ell-2}}x_{2^{\ell-2}+1}, & \text{if } i = 1 \\ x_k^2, & \text{if } i = 2k, \text{ with } k \in \{1, \ldots, 2^{\ell-2}\} \\ x_k x_{k+1}, & \text{if } i = 2k + 1, \text{ with } k \in \{1, \ldots, 2^{\ell-2} - 1\} \\ x_{2^{\ell-2}+1}, & \text{if } i = 2^{\ell-1} + 1. \end{cases} \qquad (47)$$
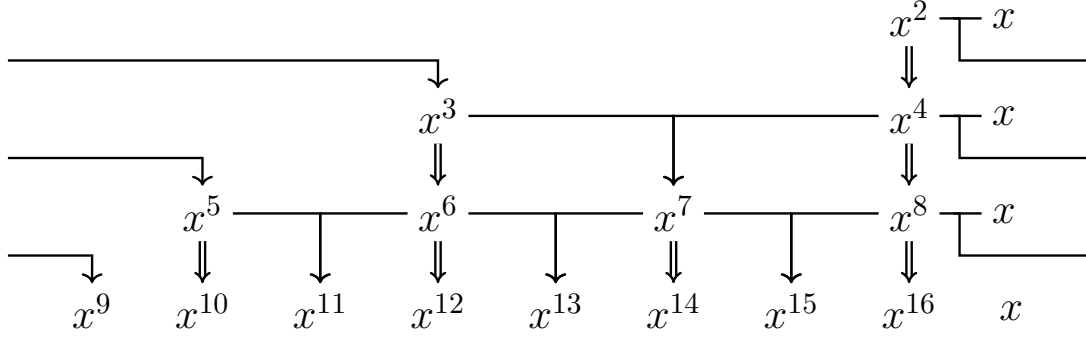
Figure 3: Expressing monomials of degree up to $N = 16$ in terms of iterated squaring ($\Rightarrow$) and multiplication ($\rightarrow$).

The functions $f_\ell$ just defined are designed to generate, through recursive application, all monomials of increasing degree. The next result makes this relationship explicit by characterizing the exact form of the compositions according to $(f^\ell \circ \cdots \circ f^1)(x) = \left( x^{2^{\ell-1}+1} \quad \ldots \quad x^{2^\ell} \quad x \right)^T$. In particular, it shows that each composition step doubles the range of available powers of $x$ while preserving the input variable itself.

**Lemma 19.** *Let $\ell \in \mathbb{N}$ and define $F^\ell := f^\ell \circ \cdots \circ f^1$. For all $x \in \mathbb{R}$, it holds that*

$$
\begin{aligned}
F^\ell(x)_i &= x^{2^{\ell-1}+i}, && \text{for } i \in \{1, \ldots, 2^{\ell-1}\}, \quad \text{and} \\
F^\ell(x)_{2^{\ell-1}+1} &= x.
\end{aligned}
\tag{48}
$$

*Proof.* Arbitrarily fix $x \in \mathbb{R}$. The proof proceeds by induction. To establish the base case, note that, for $\ell = 1$, we have

$$
\begin{aligned}
f^1(x)_1 &= x^2 = x^{2^0+1} \\
f^1(x)_{2^0+1} &= f^1(x)_2 = x.
\end{aligned}
$$

For the induction step, assume that (48) holds for some $\ell \in \mathbb{N}$ and set $z := (f^\ell \circ \cdots \circ f^1)(x) \in \mathbb{R}^{2^{\ell-1}+1}$. By the induction hypothesis

$$
z_i = x^{2^{\ell-1}+i}, \qquad \text{for } i \in \{1, \ldots, 2^{\ell-1}\}, \quad \text{and} \quad z_{2^{\ell-1}+1} = x.
$$

Using (47), we now compute

$$
\left( f^{\ell+1}(z) \right)_1 = z_{2^{\ell-1}} \cdot z_{2^{\ell-1}+1} = x^{2^{\ell-1}+2^{\ell-1}} \cdot x = x^{2^\ell+1}.
$$

For indices of the form $i = 2k$, with $k \in \{1, \ldots, 2^{\ell-1}\}$, we obtain

$$
\left( f^{\ell+1}(z) \right)_i = (z_k)^2 = \left( x^{2^{\ell-1}+k} \right)^2 = x^{2^\ell+i}.
$$

Similarly, for $i = 2k+1$, with $k \in \{1, \ldots, 2^{\ell-1} - 1\}$,

$$
\left( f^{\ell+1}(z) \right)_i = z_k \cdot z_{k+1} = x^{2^{\ell-1}+k} x^{2^{\ell-1}+k+1} = x^{2^\ell+i}.
$$

Finally,

$$
\left( f^{\ell+1}(z) \right)_{2^\ell+1} = z_{2^{\ell-1}+1} = x.
$$

Hence, (48) holds for $\ell + 1$ as well, completing the induction. $\qquad\square$

The following bound, quantifying how the range of the composed maps grows with $\ell$, is an immediate consequence of Lemma 19.

**Corollary 20.** *Let $\ell \in \mathbb{N}, D \geq 1$. For all $x \in [-D, D]$, it holds that*

$$\left\|F^\ell(x)\right\|_\infty \leq D^{2^\ell}.$$

In addition, we will make use of the following properties of the maps $f^\ell$ introduced in Definition 18.

**Lemma 21.** *Let $D \geq 1$. For all $x, y \in [-D, D]$, it holds that*

$$\left\|f^1(x)\right\|_\infty \leq D^2 \qquad and \qquad \left\|f^1(x) - f^1(y)\right\|_\infty \leq 2D|x - y|.$$

*For all $\ell \geq 2$ and all $x, y \in [-D^{2^{\ell-1}}, D^{2^{\ell-1}}]^{2^{\ell-2}+1}$,*

$$\left\|f^\ell(x)\right\|_\infty \leq D^{2^\ell} \qquad and \qquad \left\|f^\ell(x) - f^\ell(y)\right\|_\infty \leq 2D^{2^{\ell-1}}\|x - y\|_\infty.$$

*Proof.* We first consider the case $\ell = 1$. For $x \in [-D, D]$, we have $|x| \leq D \leq D^2$ and $x^2 \leq D^2$, hence $\left\|f^1(x)\right\|_\infty \leq D^2$. Moreover, for all $x, y \in [-D, D]$, $|x^2 - y^2| \leq |x||x - y| + |y||x - y| \leq 2D|x - y|$, which implies $\left\|f^1(x) - f^1(y)\right\|_\infty \leq 2D|x - y|$. We now turn to the case $\ell \geq 2$. From (47), each coordinate of $f^\ell(x)$ is either $x_{2^{\ell-2}}x_{2^{\ell-2}+1}$ (the first one), a product $x_i x_j$, with $i, j \in \{1, \ldots, 2^{\ell-2}\}$ (possibly $i = j$), or $x_{2^{\ell-2}+1}$ corresponding to the last component. For the latter the claim is immediate. Now, arbitrarily fix $x, y \in [-D^{2^{\ell-1}}, D^{2^{\ell-1}}]^{2^{\ell-2}+1}$ and $i, j \in \{1, \ldots, 2^{\ell-2} + 1\}$. We have

$$|x_i x_j| = |x_i||x_j| \leq D^{2^{\ell-1}} D^{2^{\ell-1}} = D^{2^\ell}$$

and

$$|x_i x_j - y_i y_j| \leq |x_i x_j - x_i y_j| + |x_i y_j - y_i y_j| = |x_i||x_j - y_j| + |y_j||x_i - y_i| \leq 2D^{2^{\ell-1}}\|x - y\|_\infty.$$

As $x, y$ and $i, j$ were arbitrary, the statement follows. $\qquad\square$

Next, we construct RNNs that approximate the mappings $f^\ell$. Besides the RNNs for squaring and multiplication, we also require one that realizes the identity map.

**Lemma 22.** *There exists an RNN $\mathcal{R}^{Id} = \mathcal{Q}^{Id}\mathcal{K}^{Id}$ such that, for all $x \in \mathbb{R}$, and all $t \in \mathbb{N}_0$,*

$$(\mathcal{R}^{Id}\mathcal{D}x)[t] = x, \quad \left\|(\mathcal{K}^{Id}\mathcal{D}x)[t]\right\|_\infty = |x|.$$

*Moreover, $\mathcal{M}_{\mathrm{hid}}(\mathcal{R}^{Id}) = 2$.*

*Proof.* Choose the weights as

$$A_h = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \qquad A_x = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \qquad b_h = 0, \qquad A_o = \begin{pmatrix} 1 & -1 \end{pmatrix}, \qquad b_o = 0.$$

The claim follows by direct verification based on Definition 1, using $x = \rho(x) - \rho(-x)$. $\qquad\square$

We start by building an RNN that approximates $f^1$.

**Lemma 23.** *Let $D \geq 1$. There exists an RNN $\mathcal{R}_D^1$ with $\mathcal{M}_{\mathrm{out}}(\mathcal{R}_D^1) = 2$ such that, for all $x \in [-D, D]$, and all $t \in \mathbb{N}_0$,*

$$\left|((\mathcal{R}_D^1\mathcal{D}x)[t])_1 - x^2\right| \leq \frac{D^2}{2}4^{-t}, \quad ((\mathcal{R}_D^1\mathcal{D}x)[t])_2 = x,$$

*and*

$$\left\|(\mathcal{K}_D^1\mathcal{D}x)[t]\right\|_\infty \leq D, \quad \left\|(\mathcal{R}_D^1\mathcal{D}x)[t]\right\|_\infty \leq D^2.$$

*Furthermore, $\mathcal{M}_{\mathrm{hid}}(\mathcal{R}_D^1) = 9$.*

*Proof.* Combine the RNN approximating $x^2$ from Theorem 8 with the RNN realizing the identity operator from Lemma 22. By Lemma 9, these two networks can be run in parallel, yielding an RNN that exhibits the desired behavior for inputs $\begin{pmatrix} x \\ x \end{pmatrix}$. Finally, noting that $\begin{pmatrix} x \\ x \end{pmatrix} = Ax$ with $A := \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, we apply Lemma 10 to incorporate this affine transformation into the input layer. This yields $\mathcal{R}_D^1$ with the desired behavior. $\qquad\square$

Next, we construct RNNs that approximate $f^\ell$, for arbitrary $\ell \geq 2$.

**Lemma 24.** *Let $D \geq 1$. For every $\ell \geq 2$, there exists an RNN $\mathcal{R}_D^\ell = \mathcal{Q}_D^\ell \mathcal{K}_D^\ell$ such that, for all $x \in \left[ -D^{2^{\ell-1}}, D^{2^{\ell-1}} \right]^{2^{\ell-2}+1}$, and all $t \in \mathbb{N}_0$,*

$$\left\| (\mathcal{R}_D^\ell \mathcal{D}x)[t] - f^\ell(x) \right\|_\infty \leq \frac{D^{2^\ell}}{2} 4^{-t}, \quad \left\| (\mathcal{R}_D^\ell \mathcal{D}x)[t] \right\|_\infty \leq D^{2^\ell}, \quad \left\| (\mathcal{K}_D^\ell \mathcal{D}x)[t] \right\|_\infty \leq D.$$

*Moreover, $\mathcal{M}_{\mathrm{hid}}(\mathcal{R}_D^\ell) \leq 10 \cdot 2^\ell$.*

*Proof.* Arbitrarily fix $\ell \in \mathbb{N}$, $\ell \geq 2$, and define the selector matrix $A \in \mathbb{R}^{(3 \cdot 2^{\ell-2}+1) \times (2^{\ell-2}+1)}$ that rearranges and duplicates the coordinates of $x \in \mathbb{R}^{2^{\ell-2}+1}$ according to the sequence of inputs used in (47). Explicitly,

$$Ax = \Big( \underbrace{x_{2^{\ell-2}}, x_{2^{\ell-2}+1}}_{\text{special product}}, \underbrace{x_1, \ldots, x_{2^{\ell-2}}}_{\text{squares}}, \underbrace{x_1, x_2, x_2, x_3, \ldots, x_{2^{\ell-2}-1}, x_{2^{\ell-2}}}_{\text{adjacent products}}, \underbrace{x_{2^{\ell-2}+1}}_{\text{identity}} \Big)^T. \tag{49}$$

Equivalently, the rows of $A$ are arranged as follows:

1. the first two rows select $x_{2^{\ell-2}}$ and $x_{2^{\ell-2}+1}$;

2. the next $2^{\ell-2}$ rows select $x_1, \ldots, x_{2^{\ell-2}}$;

3. for each $k \in \{1, \ldots, 2^{\ell-2} - 1\}$, two successive rows select $x_k$ and $x_{k+1}$;

4. the final row selects $x_{2^{\ell-2}+1}$.

Each row of $A$ contains exactly one nonzero entry equal to 1. For $\ell = 2$, for example, we have $x = (x_1, x_2)$ and

$$A = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

Next, let $\mathcal{R}^\times$ and $\mathcal{R}^{\mathrm{Sq}}$ denote the multiplication and squaring RNNs from Theorem 12 and Theorem 8, respectively, each instantiated with input bound $D^{2^{\ell-1}}$. Since $(D^{2^{\ell-1}})^2 = D^{2^\ell}$, it follows that, for all $x, x_1, x_2 \in [-D^{2^{\ell-1}}, D^{2^{\ell-1}}]$, and all $t \in \mathbb{N}_0$,

$$\left| \left( \mathcal{R}^\times \mathcal{D}\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right) [t] - (x_1 \cdot x_2) \right| \leq \frac{D^{2^\ell}}{2} 4^{-t}, \quad \left\| \left( \mathcal{K}^\times \mathcal{D}\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right) [t] \right\|_\infty \leq 1 \leq D, \quad \left| \left( \mathcal{R}^\times \mathcal{D}\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \right) [t] \right| \leq D^{2^\ell},$$

$$\left| \left( \mathcal{R}^{\mathrm{Sq}} \mathcal{D}x \right) [t] - x^2 \right| \leq \frac{D^{2^\ell}}{2} 4^{-t}, \quad \left\| \left( \mathcal{K}^{\mathrm{Sq}} \mathcal{D}x \right) [t] \right\|_\infty \leq 1 \leq D, \quad \left| \left( \mathcal{R}^{\mathrm{Sq}} \mathcal{D}x \right) [t] \right| \leq D^{2^\ell}.$$

To assemble the network approximating $f_\ell$, we use Lemma 9 to construct the RNN $\widetilde{\mathcal{R}}$ that runs the RNNs

$$\mathcal{R}^\times, \mathcal{R}^{\mathrm{Sq}}, \mathcal{R}^\times, \mathcal{R}^{\mathrm{Sq}}, \ldots, \mathcal{R}^\times, \mathcal{R}^{\mathrm{Sq}}, \mathcal{R}^{\mathrm{Id}}$$

in parallel, with $2^{\ell-2}$ copies of both $\mathcal{R}^\times$ and $\mathcal{R}^{\mathrm{Sq}}$. The desired RNN $\mathcal{R}_D^\ell$ is then obtained by applying Lemma 10 to $\widetilde{\mathcal{R}}$, with $A$ as defined in (49). Its hidden state size satisfies

$$\mathcal{M}_{\mathrm{hid}}(\mathcal{R}_D^\ell) = 2^{\ell-2}\left(\mathcal{M}_{\mathrm{hid}}(\mathcal{R}^{\mathrm{Sq}}) + \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^\times)\right) + \mathcal{M}_{\mathrm{hid}}(\mathcal{R}^{\mathrm{Id}}) = 2^{\ell-2}(7+14) + 2 = 2^\ell\frac{21}{4} + 2 \le 10 \cdot 2^\ell.$$

$\square$

The next step is to combine the RNNs approximating $f^\ell$ into a single RNN that approximates $F^L$. Its hidden state operator is defined as follows.

**Definition 25.** *Let $L \in \mathbb{N}$ and $D \ge 1$. Invoke Corollary 17 with the RNNs $\mathcal{R}_D^1, \dots, \mathcal{R}_D^L$ constructed in Lemmata 23 and 24, and denote the resulting hidden state operator by $\mathcal{K}_{D,L}^\pi$.*

The following lemma provides an explicit bound on the hidden-state dimension of $\mathcal{K}_{D,L}^\pi$.

**Lemma 26.** *Let $L \in \mathbb{N}$ and $D \ge 1$. The hidden state dimension of $\mathcal{K}_{D,L}^\pi$ according to Definition 25 satisfies $\mathcal{M}_{\mathrm{hid}}(\mathcal{K}_{D,L}^\pi) \le 40 \cdot 2^L$.*

*Proof.*

$$\begin{aligned}
\mathcal{M}_{\mathrm{hid}}(\mathcal{K}_{D,L}^\pi) &\overset{(45)}{\le} \sum_{\ell=1}^L \mathcal{M}_{\mathrm{hid}}(\mathcal{R}_D^\ell) + 2\sum_{\ell=1}^L \mathcal{M}_{\mathrm{out}}(\mathcal{R}_D^\ell) + 13L \\
&\le \sum_{\ell=1}^L 10 \cdot 2^\ell + 2\sum_{\ell=1}^L (2^{\ell-1} + 1) + 13L \\
&= 20\sum_{\ell=0}^{L-1} 2^\ell + 2\sum_{\ell=0}^{L-1} 2^\ell + 15L \\
&= 20(2^L - 1) + 2(2^L - 1) + 15L \\
&\le 22 \cdot 2^L + 15L \\
&\le 40 \cdot 2^L.
\end{aligned}$$

$\square$

Using the definition of $g_\ell^k$ in (32), we can now quantify how well the finite-time output of each network $\mathcal{R}_D^\ell$ approximates the corresponding target function $f^\ell$. The following corollary provides uniform bounds on the approximation error and the output magnitude.

**Corollary 27.** *Let $L \in \mathbb{N}$ and $D \ge 1$. For each $\ell \in \{1, \dots, L\}$ and every $k \ge 2$, let*

$$g_\ell^k(x) := (\mathcal{R}_D^\ell \mathcal{D}x)[2^k - 2]. \tag{50}$$

*Then, for all $x \in [-D, D]$ if $\ell = 1$ and all $x \in [-D^{2^{\ell-1}}, D^{2^{\ell-1}}]^{2^{\ell-2}+1}$ if $\ell \ge 2$, it holds that*

$$\|g_\ell^k(x) - f^\ell(x)\|_\infty \le 8D^{2^\ell}\, 4^{-2^k}, \qquad \|g_\ell^k(x)\|_\infty \le D^{2^\ell}.$$

*Proof.* Fix $L \in \mathbb{N}$, $D \ge 1$, and $k \ge 2$.

*Case $\ell = 1$.* By Lemma 23, for all $x \in [-D, D]$ and all $t \in \mathbb{N}_0$,

$$\left|((\mathcal{R}_D^1 \mathcal{D}x)[t])_1 - x^2\right| \le \frac{D^2}{2} 4^{-t}, \qquad ((\mathcal{R}_D^1 \mathcal{D}x)[t])_2 = x,$$

and hence

$$\left\|(\mathcal{R}_D^1 \mathcal{D}x)[t] - f^1(x)\right\|_\infty \le \frac{D^2}{2} 4^{-t}, \qquad \left\|(\mathcal{R}_D^1 \mathcal{D}x)[t]\right\|_\infty \le D^2.$$

Evaluating at $t = 2^k - 2$ and using $g_1^k(x) = (\mathcal{R}_D^1 \mathcal{D}x)[2^k - 2]$, yields

$$\|g_1^k(x) - f^1(x)\|_\infty \le \frac{D^2}{2} 4^{-(2^k-2)} = 8D^2\, 4^{-2^k}, \qquad \|g_1^k(x)\|_\infty \le D^2,$$

which matches the claimed bounds for $\ell = 1$.

*Case $\ell \geq 2$.* By Lemma 24, for all admissible $x$ and all $t \in \mathbb{N}_0$,

$$\left\|(\mathcal{R}_D^\ell \mathcal{D}x)[t] - f^\ell(x)\right\|_\infty \leq \frac{D^{2^\ell}}{2} 4^{-t}, \qquad \left\|(\mathcal{R}_D^\ell \mathcal{D}x)[t]\right\|_\infty \leq D^{2^\ell}.$$

Setting $t = 2^k - 2$ and using $g_\ell^k(x) = (\mathcal{R}_D^\ell \mathcal{D}x)[2^k - 2]$, this yields

$$\|g_\ell^k(x) - f^\ell(x)\|_\infty \leq \frac{D^{2^\ell}}{2} 4^{-(2^k-2)} = 8D^{2^\ell} 4^{-2^k}, \qquad \|g_\ell^k(x)\|_\infty \leq D^{2^\ell}.$$

The two cases together prove the corollary. $\qquad \square$

We now demonstrate that each composite function $F^\ell$ can be approximated by applying an affine output map to the hidden-state sequence generated by $\mathcal{K}_{D,L}^\pi$ (see Definition 25).

**Lemma 28.** *Let $L \in \mathbb{N}$ and $D \geq 1$. For each $\ell \in \{1, \ldots, L\}$, there exist matrices $A^\ell$ and vectors $b^\ell$ such that, for all $x \in [-D, D]$, and all $k \geq \lceil \log(L) \rceil + 2$,*

$$\left\|A^\ell (\mathcal{K}_{D,L}^\pi \mathcal{D}x)[2^k - 2] + b^\ell - F^\ell(x)\right\|_\infty \leq 8 \cdot 2^\ell D^{2^\ell} 4^{-\frac{1}{2\ell}2^k} =: \varepsilon_{\ell,k}.$$

*Proof.* Let $L \in \mathbb{N}$ and $D \geq 1$, and define $g_\ell^k$ as in (32) for $\ell \in \{1, \ldots, L\}$ and $k \geq 2$. Fix $\ell \in \{1, \ldots, L\}$ arbitrarily. By Corollary 17, there exist matrices $A^\ell$, vectors $b^\ell$, and integers $\tilde{k}_1, \ldots, \tilde{k}_\ell \leq \lceil \log \ell \rceil$ such that

$$A^\ell (\mathcal{K}_{D,L}^\pi \mathcal{D}x)[2^k - 2] + b^\ell = \left(g_\ell^{k-\tilde{k}_\ell} \circ \cdots \circ g_1^{k-\tilde{k}_1}\right)(x), \text{ for all } k \geq \lceil \log(L) \rceil + 2, \text{ and all } x \in [-D, D]. \quad (51)$$

We now verify that $A^\ell$ and $b^\ell$ satisfy the desired properties. To this end, define for each $\ell' \in \{1, \ldots, \ell\}$ and $k \geq \lceil \log(L) \rceil + 2$, the mappings

$$G_{\ell'}^k := g_{\ell'}^{k-\tilde{k}_{\ell'}} \circ \cdots \circ g_1^{k-\tilde{k}_1}.$$

Fix $x \in [-D, D]$ and $k \geq \lceil \log(L) \rceil + 2$. We show by induction on $\ell' \in \{1, \ldots, \ell\}$ that

$$\|G_{\ell'}^k(x)\|_\infty \leq D^{2^{\ell'}}, \tag{52}$$

and

$$\|G_{\ell'}^k(x) - F^{\ell'}(x)\|_\infty \leq 8D^{2^{\ell'}} 4^{-2^{k-\lceil\log(\ell)\rceil}} \left(\sum_{i=0}^{\ell'-1} 2^i\right). \tag{53}$$

The base case, corresponding to $\ell' = 1$, follows directly from Corollary 27. Indeed, since $G_1^k = g_1^{k-\tilde{k}_1}$, we have for all $x \in [-D, D]$,

$$\|G_1^k(x)\|_\infty = \|g_1^{k-\tilde{k}_1}(x)\|_\infty \leq D^2,$$
$$\|G_1^k(x) - F^1(x)\|_\infty = \|g_1^{k-\tilde{k}_1}(x) - f^1(x)\|_\infty \leq 8D^2 4^{-2^{k-\tilde{k}_1}} \leq 8D^2 4^{-2^{k-\lceil\log(\ell)\rceil}}.$$

Next, assume that (52) and (53) hold for some $\ell' \in \{1, \ldots, \ell-1\}$, and consider $\ell' + 1$. We have

$$\|G_{\ell'+1}^k(x)\|_\infty = \|g_{\ell'+1}^{k-\tilde{k}_{\ell'+1}}(G_{\ell'}^k(x))\|_\infty \leq \sup_{\|z\|_\infty \leq D^{2^{\ell'}}} \|g_{\ell'+1}^{k-\tilde{k}_{\ell'+1}}(z)\|_\infty \leq D^{2^{\ell'+1}},$$

where the last inequality again follows from Corollary 27. This shows that (52) holds for $\ell' + 1$ as well.

26

Next, we have

$$\left\|G_{\ell'+1}^k(x) - F^{\ell'+1}(x)\right\|_\infty = \left\|g_{\ell'+1}^{k-\widetilde{k}_{\ell'+1}}(G_{\ell'}^k(x)) - f^{\ell'+1}(F^{\ell'}(x))\right\|_\infty$$

$$\leq \left\|g_{\ell'+1}^{k-\widetilde{k}_{\ell'+1}}(G_{\ell'}^k(x)) - f^{\ell'+1}(G_{\ell'}^k(x))\right\|_\infty + \left\|f^{\ell'+1}(G_{\ell'}^k(x)) - f^{\ell'+1}(F^{\ell'}(x))\right\|_\infty$$

$$\overset{(52)}{\leq} \sup_{\|z\|_\infty \leq D^{2^{\ell'}}} \left\|g_{\ell'+1}^{k-\widetilde{k}_{\ell'+1}}(z) - f^{\ell'+1}(z)\right\|_\infty + \left\|f^{\ell'+1}(G_{\ell'}^k(x)) - f^{\ell'+1}(F^{\ell'}(x))\right\|_\infty$$

$$\overset{\text{Cor. } 27}{\leq} 8D^{2^{\ell'+1}}4^{-2^{k-\widetilde{k}_{\ell'+1}}} + \left\|f^{\ell'+1}(G_{\ell'}^k(x)) - f^{\ell'+1}(F^{\ell'}(x))\right\|_\infty$$

$$\overset{\text{Lem. } 21}{\leq} 8D^{2^{\ell'+1}}4^{-2^{k-\widetilde{k}_{\ell'+1}}} + 2D^{2^{\ell'}}\left\|G_{\ell'}^k(x) - F^{\ell'}(x)\right\|_\infty$$

$$\overset{(53)}{\leq} 8D^{2^{\ell'+1}}4^{-2^{k-\widetilde{k}_{\ell'+1}}} + 2D^{2^{\ell'}}\left(8 \cdot D^{2^{\ell'}}4^{-2^{k-\lceil\log(\ell)\rceil}}\left(\sum_{i=0}^{\ell'-1} 2^i\right)\right)$$

$$\overset{(*)}{\leq} 8D^{2^{\ell'+1}}4^{-2^{k-\lceil\log(\ell)\rceil}} + 8 \cdot D^{2^{\ell'+1}}4^{-2^{k-\lceil\log(\ell)\rceil}}\left(\sum_{i=0}^{\ell'-1} 2^{i+1}\right)$$

$$= 8D^{2^{\ell'+1}}4^{-2^{k-\lceil\log(\ell)\rceil}}\left(\sum_{i=0}^{\ell'} 2^i\right),$$

where in $(*)$ we used $\widetilde{k}_{\ell'+1} \leq \lceil\log(\ell)\rceil$. This proves that (53) holds for $\ell' + 1$ as well. In particular, using (53) with $\ell' = \ell$, we can thus compute

$$\left\|A^\ell\left(\mathcal{K}_{D,L}^\pi \mathcal{D}x\right)[2^k - 2] + b^\ell - F^\ell(x)\right\|_\infty \overset{(51)}{=} \left\|G_\ell^k(x) - F^\ell(x)\right\|_\infty$$

$$\overset{(53)}{\leq} 8 \cdot D^{2^\ell}4^{-2^{k-\lceil\log(\ell)\rceil}}\left(\sum_{i=0}^{\ell-1} 2^i\right)$$

$$= 8 \cdot D^{2^\ell}4^{-2^{k-\lceil\log(\ell)\rceil}}(2^\ell - 1)$$

$$\leq 8 \cdot 2^\ell \cdot D^{2^\ell}4^{-2^{k-\log(\ell)-1}}$$

$$= 8 \cdot 2^\ell \cdot D^{2^\ell}4^{-\frac{1}{2\ell}2^k}.$$

Since $\ell \in \{1, \ldots, L\}$, $k \geq \lceil\log(L)\rceil + 2$, and $x \in [-D, D]$ were arbitrary, this completes the proof. $\quad\square$

Having constructed the hidden-state operator $\mathcal{K}_{D,L}^\pi$ in Definition 25 and established its approximation properties in Lemma 28, we now specify the affine output map that extracts the desired polynomial terms from its hidden-state sequence.

**Lemma 29.** *Let $L \in \mathbb{N}$ and $D \geq 1$. There exist matrices $A_o^\pi$ and vectors $b_o^\pi$, such that the affine output mapping $\mathcal{Q}_{D,L}^\pi : h \to A_o^\pi h + b_o^\pi$ satisfies, for all $x \in [-D, D]$ and all $k \geq \lceil\log(L)\rceil + 2$,*

$$((\mathcal{Q}_{D,L}^\pi \mathcal{K}_{D,L}^\pi \mathcal{D}x)[2^k - 2])_1 = x,$$

*and, for every $\ell \in \{1, \ldots, L\}$ and $j \in \{1, \ldots, 2^{\ell-1}\}$,*

$$|((\mathcal{Q}_{D,L}^\pi \mathcal{K}_{D,L}^\pi x)[2^k - 2])_{2^{\ell-1}+j} - x^{2^{\ell-1}+j}| \leq \varepsilon_{\ell,k}.$$

*Here, $\varepsilon_{\ell,k}$ denotes the approximation error defined in Lemma 28, and $\mathcal{K}_{D,L}^\pi$ is the hidden-state operator introduced in Definition 25.*

27

*Proof.* Fix $L \in \mathbb{N}$ and $D \geq 1$. Let $A^\ell, b^\ell$ be as in Lemma 28 and note that $A^\ell \in \mathbb{R}^{(2^{\ell-1}+1) \times \mathcal{M}_{\mathrm{hid}}(\mathcal{K}^\pi_{D,L})}$ and $b^\ell \in \mathbb{R}^{2^{\ell-1}+1}$, for $\ell \in \{1, \ldots, L\}$. Define the matrices and vectors

$$\widetilde{A^1_o} := \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} A^1, \qquad\qquad \widetilde{b^1_o} := \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} b^1,$$

$$\widetilde{A^\ell_o} := \begin{pmatrix} \mathbb{I}_{2^{\ell-1}} & 0 \end{pmatrix} A^\ell, \qquad \widetilde{b^\ell_o} := \begin{pmatrix} \mathbb{I}_{2^{\ell-1}} & 0 \end{pmatrix} b^\ell, \quad \text{for } \ell \in \{2, \ldots, L\},$$

and the output mapping

$$\mathcal{Q}^\pi_{D,L}(h) := A^\pi_o h + b^\pi_o, \text{ with } A^\pi_o := \begin{pmatrix} \widetilde{A^1_o} \\ \widetilde{A^2_o} \\ \vdots \\ \widetilde{A^L_o} \end{pmatrix} \text{ and } b^\pi_o := \begin{pmatrix} \widetilde{b^1_o} \\ \widetilde{b^2_o} \\ \vdots \\ \widetilde{b^L_o} \end{pmatrix}.$$

We now verify that $\mathcal{Q}^\pi_{D,L}$ satisfies the claimed properties. To this end, fix $x \in [-D, D]$ and $k \geq \lceil \log(L) \rceil + 2$ arbitrarily and note that, by (51) in the proof of Lemma 28,

$$A^1(\mathcal{K}^\pi_{D,L}\mathcal{D}x)[2^k - 2] + b^1 = g^{k - \widetilde{k}_1}_1(x)$$

and, by (50), $g^{k - \widetilde{k}_1}_1(x) = (\mathcal{R}^1_D \mathcal{D}x)[2^{k - \widetilde{k}_1} - 2]$. Combining the two identities yields

$$((\mathcal{Q}^\pi_{D,L}\mathcal{K}^\pi_{D,L}\mathcal{D}x)[2^k - 2])_1 = \left( A^1(\mathcal{K}^\pi_{D,L}\mathcal{D}x)[2^k - 2] + b^1 \right)_2 = \left( (\mathcal{R}^1_D \mathcal{D}x)\left[2^{k - \widetilde{k}_1} - 2\right] \right)_2 \overset{\text{Lem. } 23}{=} x.$$

For the second entry we have

$$\left| ((\mathcal{Q}^\pi_{D,L}\mathcal{K}^\pi_{D,L}\mathcal{D}x)[2^k - 2])_2 - x^{2^0 + 1} \right| = \left| \left( A^1(\mathcal{K}^\pi_{D,L}\mathcal{D}x)[2^k - 2] + b^1 \right)_1 - F^1(x)_1 \right| \overset{\text{Lem. } 28}{\leq} \varepsilon_{1,k}.$$

For $\ell \in \{2, \ldots, L\}$ and $j \in \{1, \ldots, 2^{\ell-1}\}$, consider the $(2^{\ell-1}+j)$-th coordinate of $(\mathcal{Q}^\pi_{D,L}\mathcal{K}^\pi_{D,L}\mathcal{D}x)[2^k - 2]$. Since $\widetilde{A^1_o}$ has two rows and, for $i = 2, \ldots, \ell - 1$, each block $\widetilde{A^i_o}$ has $2^{i-1}$ rows, the total number of preceding rows is

$$2 + \sum_{i=2}^{\ell-1} 2^{i-1} = 2^{\ell-1}.$$

Consequently, the $(2^{\ell-1}+j)$-th coordinate of the global output corresponds to the $j$-th coordinate of $A^\ell(\mathcal{K}^\pi_{D,L}\mathcal{D}x)[2^k - 2] + b^\ell$. By Lemma 28, we obtain

$$\left| (\mathcal{Q}^\pi_{D,L}\mathcal{K}^\pi_{D,L}\mathcal{D}x)[2^k - 2]_{2^{\ell-1}+j} - x^{2^{\ell-1}+j} \right| = \left| (A^\ell(\mathcal{K}^\pi_{D,L}\mathcal{D}x)[2^k - 2] + b^\ell)_j - F^\ell(x)_j \right| \leq \varepsilon_{\ell,k}.$$

This completes the proof. $\qquad\square$

We now consolidate the preceding results into a single statement describing the full RNN and its approximation properties.

**Theorem 30.** *Let $L \in \mathbb{N}$ and $D \geq 1$. Then, there exists an RNN $\mathcal{R}^\pi_{D,L} = \mathcal{Q}^\pi_{D,L}\mathcal{K}^\pi_{D,L}$ such that*

$$\mathcal{M}_{\mathrm{hid}}(\mathcal{R}^\pi_{D,L}) \leq 40 \cdot 2^L, \tag{56}$$

*and, for all $x \in [-D, D]$ and all $k \geq \lceil \log L \rceil + 2$,*

$$((\mathcal{R}^\pi_{D,L}\mathcal{D}x)[2^k - 2])_1 = x, \tag{57}$$

*and, for every $\ell \in \{1, \ldots, L\}$ and $j \in \{1, \ldots, 2^{\ell-1}\}$,*

$$|((\mathcal{R}^\pi_{D,L}x)[2^k - 2])_{2^{\ell-1}+j} - x^{2^{\ell-1}+j}| \leq \varepsilon_{\ell,k}, \tag{58}$$

*where $\varepsilon_{\ell,k} = 8 \cdot 2^\ell D^{2^\ell} 4^{-\frac{1}{2\ell}2^k}$.*

*Proof.* Fix $L \in \mathbb{N}$ and $D \geq 1$. All results in this section apply for these values. Let $\mathcal{K}^\pi_{D,L}$ be as in Definition 25 and $\mathcal{Q}^\pi_{D,L}$ as in Lemma 29. Now, (56) follows from Lemma 26 and (57) and (58) are by Lemma 29. $\qquad\square$

# 5 Approximation of polynomials

It is now clear that any monomial can be approximated by an RNN according to Theorem 30. This construction, however, produces meaningful outputs only at discrete time steps $t = 2^k - 2$. The following modification yields an RNN that generates a continuous sequence of outputs by holding the most recent valid approximation, rather than remaining idle between update times. Specifically, the output of the modified RNN remains constant between successive update times and is clipped to a prescribed range $[-B, B]$ to ensure boundedness; this clipping is operationally irrelevant, since $B > 0$ may be chosen arbitrarily large.

**Lemma 31.** *Fix $B > 0$ and let $\mathcal{R} = \mathcal{QK}$ be an RNN with $\mathcal{M}_{\mathrm{out}}(\mathcal{R}) = 1$. Then there exists an RNN $\mathcal{R}'$ such that, for all $x \in \mathbb{R}$,*

$$(\mathcal{R}'\mathcal{D}x)[2^k - 1 + \ell] = \mathcal{C}((\mathcal{R}\mathcal{D}x)[2^k - 2], -B, B), \quad k \geq 2, \ell \in \{0, \ldots, 2^k - 1\},$$

*with the clipping operator*

$$\mathcal{C}(y, A, B) := \begin{cases} A, & \text{if } y \leq A \\ B, & \text{if } y \geq B \\ y, & \text{otherwise.} \end{cases}$$

*Furthermore, the hidden state dimension satisfies*

$$\mathcal{M}_{\mathrm{hid}}(\mathcal{R}') = \mathcal{M}_{\mathrm{hid}}(\mathcal{R}) + 11.$$

*Proof.* Let the weights of $\mathcal{R}$ be $A_x \in \mathbb{R}^{m \times d}$, $A_h \in \mathbb{R}^{m \times m}$, $b_h \in \mathbb{R}^m$, $A_o \in \mathbb{R}^{1 \times m}$, and $b_o \in \mathbb{R}$. The modified RNN $\mathcal{R}'$ has hidden state dimension $m + 11$ and corresponding weights $A'_x \in \mathbb{R}^{(m+11) \times d}$, $A'_h \in \mathbb{R}^{(m+11) \times (m+11)}$, $b'_h \in \mathbb{R}^{m+11}$, $A'_o \in \mathbb{R}^{1 \times (m+11)}$, $b'_o \in \mathbb{R}$, which are given by

$$A'_h = \left( \begin{array}{c|c|c|c|c} A_h & 0 & 0 & 0 & 0 \\ \hline A_o & 0 & 0 & 0 & B\widehat{A}_o \\ \hline -A_o & 0 & 0 & 0 & B\widehat{A}_o \\ \hline A_o & 0 & 0 & 0 & 0 \\ \hline -A_o & 0 & 0 & 0 & 0 \\ \hline 0 & \mathbb{I}_2 & -\mathbb{I}_2 & \mathbb{I}_2 & -B\mathbf{1}_2\widehat{A}_o \\ \hline 0 & 0 & 0 & 0 & \widehat{A} \end{array} \right), \quad b'_h = \left( \begin{array}{c} b_h \\ \hline b_o - B \\ \hline -b_o - B \\ \hline b_o - B \\ \hline -b_o - B \\ \hline 0 \\ \hline \widehat{b} \end{array} \right), \quad A'_x = \left( \begin{array}{c} A_x \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \end{array} \right),$$

$$A'_o = \left( \begin{array}{c|cc|cc|cc|c} 0 & 1 & -1 & -1 & 1 & 1 & -1 & 0 \end{array} \right), \quad b'_o = 0,$$

where $\widehat{A}_o := \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \end{pmatrix}$ and $\widehat{A}$ and $\widehat{b}$ are as in (14). All zero blocks are of the dimensions required for blockwise compatibility. Fix an arbitrary input $x \in \mathbb{R}^d$ and consider the hidden state sequence $h'[\cdot]$ of $\mathcal{R}'$ generated by the corresponding input sequence $(\mathcal{D}x)[\cdot]$. We partition $h'[\cdot]$ according to the block structure of $A'_h$ as follows

$$h'[t] = \left( \begin{array}{c} h'_1[t] \\ \hline h'_2[t] \\ \hline h'_3[t] \\ \hline h'_4[t] \\ \hline h'_5[t] \end{array} \right) = \rho \left( A'_h h'[t-1] + A'_x ((Dx)[t]) + b'_h \right), \, t \in \mathbb{N}_0,$$

and analyze the dynamics of each block individually. Observe that $h'_5[\cdot]$ evolves according to

$$h'_5[t] = \rho(\widehat{A}h'_5[t-1] + \widehat{b}), \, t \in \mathbb{N}_0, \text{ with } h'_5[-1] = 0.$$

By Lemma 13, this implies

$$\widehat{A}_o h'_5[t] = \widehat{\delta}[t+2], \quad \text{for all } t \in \mathbb{N}_0. \tag{59}$$

The first block $h'_1[t]$ follows the recursion

$$h'_1[t] = \rho\left(A_h h'_1[t-1] + A_x((\mathcal{D}x)[t]) + b_h\right), \ t \in \mathbb{N}_0, \quad \text{with } h'_1[-1] = 0,$$

and hence coincides with the hidden state of the base network, $h'_1[t] = (\mathcal{K}\mathcal{D}x)[t]$. Next, the second block $h'_2[t]$ evolves according to

$$h'_2[t] = \rho\begin{pmatrix} A_o h'_1[t-1] + B\widehat{A}_o h'_5[t-1] + b_o - B \\ -A_o h'_1[t-1] + B\widehat{A}_o h'_5[t-1] - b_o - B \end{pmatrix}$$

$$= \rho\begin{pmatrix} (A_o(\mathcal{K}\mathcal{D}x)[t-1] + b_o) - B(1 - \widehat{A}_o h'_5[t-1]) \\ -(A_o(\mathcal{K}\mathcal{D}x)[t-1] + b_o) - B(1 - \widehat{A}_o h'_5[t-1]) \end{pmatrix}$$

$$= \rho\begin{pmatrix} (\mathcal{R}\mathcal{D}x)[t-1] - B(1 - \widehat{\delta}[t+1]) \\ -(\mathcal{R}\mathcal{D}x)[t-1] - B(1 - \widehat{\delta}[t+1]) \end{pmatrix}.$$

Similarly, we get

$$h'_3[t] = \rho\begin{pmatrix} A_o h'_1[t-1] + b_o - B \\ -A_o h'_1[t-1] - b_o - B \end{pmatrix}$$

$$= \rho\begin{pmatrix} (\mathcal{R}\mathcal{D}x)[t-1] - B \\ -(\mathcal{R}\mathcal{D}x)[t-1] - B \end{pmatrix} \tag{60}$$

$$= \begin{pmatrix} ((\mathcal{R}\mathcal{D}x)[t-1] - B)\, \mathbb{1}_{\{(\mathcal{R}\mathcal{D}x)[t-1] \geq B\}} \\ (-(\mathcal{R}\mathcal{D}x)[t-1] - B)\, \mathbb{1}_{\{-(\mathcal{R}\mathcal{D}x)[t-1] \geq B\}} \end{pmatrix}$$

$$= \begin{pmatrix} (\rho\left((\mathcal{R}\mathcal{D}x)[t-1]\right) - B)\, \mathbb{1}_{\{(\mathcal{R}\mathcal{D}x)[t-1] \geq B\}} \\ (\rho\left(-(\mathcal{R}\mathcal{D}x)[t-1]\right) - B)\, \mathbb{1}_{\{-(\mathcal{R}\mathcal{D}x)[t-1] \geq B\}} \end{pmatrix}. \tag{61}$$

Next, we show that

$$h'_2[t] - h'_3[t] = \begin{pmatrix} \mathcal{C}((\mathcal{R}\mathcal{D}x)[t-1], 0, B) \\ \mathcal{C}(-(\mathcal{R}\mathcal{D}x)[t-1], 0, B) \end{pmatrix} \widehat{\delta}[t+1]. \tag{62}$$

This follows by a case distinction. If $\widehat{\delta}[t+1] = 0$, then by (60),

$$h'_2[t] - h'_3[t] = \rho\begin{pmatrix} (\mathcal{R}\mathcal{D}x)[t-1] - B \\ -(\mathcal{R}\mathcal{D}x)[t-1] - B \end{pmatrix} - \rho\begin{pmatrix} (\mathcal{R}\mathcal{D}x)[t-1] - B \\ -(\mathcal{R}\mathcal{D}x)[t-1] - B \end{pmatrix} = 0.$$

For $\widehat{\delta}[t+1] = 1$, we obtain from (61),

$$h'_2[t] - h'_3[t] = \rho\begin{pmatrix} (\mathcal{R}\mathcal{D}x)[t-1] \\ -(\mathcal{R}\mathcal{D}x)[t-1] \end{pmatrix} - \begin{pmatrix} (\rho\left((\mathcal{R}\mathcal{D}x)[t-1]\right) - B)\, \mathbb{1}_{\{(\mathcal{R}\mathcal{D}x)[t-1] \geq B\}} \\ (\rho\left(-(\mathcal{R}\mathcal{D}x)[t-1]\right) - B)\, \mathbb{1}_{\{-(\mathcal{R}\mathcal{D}x)[t-1] \geq B\}} \end{pmatrix}$$

$$= \begin{pmatrix} \mathcal{C}((\mathcal{R}\mathcal{D}x)[t-1], 0, B) \\ \mathcal{C}(-(\mathcal{R}\mathcal{D}x)[t-1], 0, B) \end{pmatrix},$$

thereby establishing (62).

Hence, $h'_4[t], t \in \mathbb{N}_0$, satisfies the recurrence

$$h'_4[t] = \rho\left(h'_2[t-1] - h'_3[t-1] + h'_4[t-1] - B\mathbf{1}_2\widehat{A}_o h'_5[t-1]\right)$$

$$\overset{(62),(59)}{=} \rho\left(\begin{pmatrix} \mathcal{C}((\mathcal{R}\mathcal{D}x)[t-2], 0, B) \\ \mathcal{C}(-(\mathcal{R}\mathcal{D}x)[t-2], 0, B) \end{pmatrix}\widehat{\delta}[t] + h'_4[t-1] - B\mathbf{1}_2\widehat{\delta}[t+1]\right), \text{with } h'_4[-1] = 0. \tag{63}$$

30

Note that $h'_4[t] = 0$, for $t \in \{0, 1, 2, 3\}$. We now prove by nested induction that

$$h'_4[2^k - 1] = 0, \qquad\qquad\qquad\qquad\text{for } k \in \mathbb{N}, \ k \geq 2, \qquad (64)$$

$$\text{and} \quad h'_4[2^k + \ell] = \begin{pmatrix} \mathcal{C}((\mathcal{RD}x)[2^k - 2], 0, B) \\ \mathcal{C}(-(\mathcal{RD}x)[2^k - 2], 0, B) \end{pmatrix}, \qquad \text{for } \ell \in \{0, \ldots, 2^k - 2\}. \qquad (65)$$

We have already concluded that (64) holds for $k = 2$. Assume now that (64) is valid for some $k \geq 2$, and compute

$$h'_4[2^k] = \rho\left( \begin{pmatrix} \mathcal{C}((\mathcal{RD}x)[2^k - 2], 0, B) \\ \mathcal{C}(-(\mathcal{RD}x)[2^k - 2], 0, B) \end{pmatrix} \widehat{\delta}[2^k] + h'_4[2^k - 1] - B\mathbf{1}_2\widehat{\delta}[2^k + 1] \right)$$

$$= \begin{pmatrix} \mathcal{C}((\mathcal{RD}x)[2^k - 2], 0, B) \\ \mathcal{C}(-(\mathcal{RD}x)[2^k - 2], 0, B) \end{pmatrix}.$$

This verifies (65) for $\ell = 0$. Next, assume that (65) holds for some $\ell \in \{0, \ldots, 2^k - 3\}$, and compute

$$h'_4[2^k + \ell + 1] = \rho\left( \begin{pmatrix} \mathcal{C}((\mathcal{RD}x)[2^k + \ell - 1], 0, B) \\ \mathcal{C}(-(\mathcal{RD}x)[2^k + \ell - 1], 0, B) \end{pmatrix} \widehat{\delta}[2^k + \ell + 1] + h'_4[2^k + \ell] - B\mathbf{1}_2\widehat{\delta}[2^k + \ell + 2] \right)$$

$$= \rho(h'_4[2^k + \ell]) = \begin{pmatrix} \mathcal{C}((\mathcal{RD}x)[2^k - 2], 0, B) \\ \mathcal{C}(-(\mathcal{RD}x)[2^k - 2], 0, B) \end{pmatrix}.$$

This completes the induction over $\ell$ and establishes (65). To complete the induction step for (64), we use (63) to get

$$h'_4[2^{k+1} - 1] = \rho\left( \begin{pmatrix} \mathcal{C}((\mathcal{RD}x)[2^{k+1} - 3], 0, B) \\ \mathcal{C}(-(\mathcal{RD}x)[2^{k+1} - 3], 0, B) \end{pmatrix} \widehat{\delta}[2^{k+1} - 1] + h'_4[2^{k+1} - 2] - B\mathbf{1}_2\widehat{\delta}[2^{k+1}] \right)$$

$$= \rho\left( h'_4[2^k + 2^k - 2] - B\mathbf{1}_2 \right)$$

$$\overset{(65)}{=} \rho\left( \begin{pmatrix} \mathcal{C}((\mathcal{RD}x)[2^k - 2], 0, B) \\ \mathcal{C}(-(\mathcal{RD}x)[2^k - 2], 0, B) \end{pmatrix} - B\mathbf{1}_2 \right) = 0.$$

This establishes (64) for $k+1$ and thereby completes the overall nested induction. The network output is given by

$$(\mathcal{R}'\mathcal{D}x)[t] = \begin{pmatrix} 1 & -1 \end{pmatrix} h'_2[t] + \begin{pmatrix} -1 & 1 \end{pmatrix} h'_3[t] + \begin{pmatrix} 1 & -1 \end{pmatrix} h'_4[t]$$

$$= \begin{pmatrix} 1 & -1 \end{pmatrix} (h'_2[t] - h'_3[t] + h'_4[t])$$

$$\overset{(62)}{=} \begin{pmatrix} 1 & -1 \end{pmatrix} \left( \begin{pmatrix} \mathcal{C}((\mathcal{RD}x)[t - 1], 0, B) \\ \mathcal{C}(-(\mathcal{RD}x)[t - 1], 0, B) \end{pmatrix} \widehat{\delta}[t + 1] + h'_4[t] \right).$$

By (64), (65), and the fact that $\widehat{\delta}[t + 1] = 1$, for $t = 2^k - 1$ with $k \geq 2$, and $\widehat{\delta}[t + 1] = 0$ otherwise, we find

$$(\mathcal{R}'\mathcal{D}x)[2^k - 1 + \ell] = \begin{pmatrix} 1 & -1 \end{pmatrix} \begin{pmatrix} \mathcal{C}((\mathcal{RD}x)[2^k - 2], 0, B) \\ \mathcal{C}(-(\mathcal{RD}x)[2^k - 2], 0, B) \end{pmatrix}$$

$$= \mathcal{C}((\mathcal{RD}x)[2^k - 2], -B, B), \qquad \ell \in \{0, \ldots, 2^k - 1\},$$

as claimed. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \square$

To describe the RNN behavior uniformly across all time indices, we next introduce a convenient re-indexing of the dyadic intervals arising in Lemma 31.

**Lemma 32.** *For $t \in \mathbb{N}$, let $\widetilde{k}(t) := \lfloor \log(t+1) \rfloor$ and $\widetilde{\ell}(t) := t + 1 - 2^{\widetilde{k}(t)}$. Then,*

$$t = 2^{\widetilde{k}(t)} - 1 + \widetilde{\ell}(t) \quad \text{and} \quad \widetilde{\ell}(t) \in \{0, \dots, 2^{\widetilde{k}(t)} - 1\}.$$

*Proof.* Arbitrarily fix $t \in \mathbb{N}$ and let $\widetilde{k} := \widetilde{k}(t)$ and $\widetilde{\ell} := \widetilde{\ell}(t)$. Then,

$$2^{\widetilde{k}} - 1 + \widetilde{\ell} = 2^{\widetilde{k}} - 1 + t + 1 - 2^{\widetilde{k}} = t.$$

Moreover, we have

$$\widetilde{\ell} = t + 1 - 2^{\lfloor \log(t+1) \rfloor} \geq t + 1 - 2^{\log(t+1)} = t + 1 - t - 1 = 0.$$

Furthermore,

$$
\begin{aligned}
\widetilde{\ell} &= t + 1 - 2^{\lfloor \log(t+1) \rfloor} \\
&= 2 \cdot 2^{\log(t+1)-1} - 2^{\lfloor \log(t+1) \rfloor} \\
&< 2 \cdot 2^{\lfloor \log(t+1) \rfloor} - 2^{\lfloor \log(t+1) \rfloor} \\
&= 2^{\lfloor \log(t+1) \rfloor} = 2^{\widetilde{k}},
\end{aligned}
$$

so $\widetilde{\ell} \in \{0, ..., 2^{\widetilde{k}} - 1\}$. $\qquad\square$

We now combine the constructions from the previous sections to approximate arbitrary polynomials by a single RNN. The following theorem summarizes the result.

**Theorem 33.** *Let $N \in \mathbb{N}$, $D \geq 1$, and $a_0, \dots, a_N \in \mathbb{R}$. There exists an RNN $\mathcal{R}_a^\pi$ such that, for all $x \in [-D, D]$ and all $t \geq 16 \log(N)$,*

$$\left| (\mathcal{R}_a^\pi \mathcal{D}x)[t] - \sum_{i=0}^{N} a_i x^i \right| \leq \|a\|_1 C_1 4^{-C_2 t},$$

*with $C_1 = 16ND^{2N}$, $C_2 = \frac{1}{4\lceil \log(N) \rceil}$, and $\|a\|_1 = \sum_{i=0}^{N} |a_i|$. Furthermore, $\mathcal{M}_{\mathrm{hid}}(\mathcal{R}_a^\pi) \leq 80N + 11$.*

*Proof.* We may assume without loss of generality that $N \geq 2$. Indeed, if $N \leq 1$, then by Lemma 22 the identity RNN $\mathcal{R}^{\mathrm{Id}}$ satisfies $(\mathcal{R}^{\mathrm{Id}} \mathcal{D}x)[t] = x$, for all $t \in \mathbb{N}_0$. Applying Lemma 11 with the affine readout $h \mapsto a_1 h + a_0$ (with $a_1 = 0$ when $N = 0$) yields an RNN that realizes the polynomial $a_1 x + a_0$ exactly. Let $L := \lceil \log(N) \rceil \in \mathbb{N}$ and extend the coefficient sequence by setting $a_{N+1} = \cdots = a_{2^L} = 0$. We now modify the RNN $\mathcal{R}_{D,L}^\pi$ from Theorem 30 using Lemma 11 with $A = \left( a_1, \dots, a_{2^L} \right)$ and $b = a_0$ to obtain the RNN $\widetilde{\mathcal{R}_a^\pi}$ realizing

$$(\widetilde{\mathcal{R}_a^\pi} \mathcal{D}x)[t] = a_0 + \sum_{i=1}^{2^L} a_i \left( (\mathcal{R}_{D,L}^\pi \mathcal{D}x)[t] \right)_i.$$

We bound, for $k \geq \lceil \log(L) \rceil + 2$,

$$
\begin{aligned}
\left| (\widetilde{\mathcal{R}_a^\pi} \mathcal{D}x)[2^k - 2] - \sum_{i=0}^{N} a_i x^i \right| &\leq \sum_{i=2}^{2^L} |a_i| \left| \left( (\mathcal{R}_{D,L}^\pi \mathcal{D}x)[2^k - 2] \right)_i - x^i \right| \\
&= \sum_{\ell=1}^{L} \sum_{j=1}^{2^{\ell-1}} |a_{2^{\ell-1}+j}| \left| \left( (\mathcal{R}_{D,L}^\pi \mathcal{D}x)[2^k - 2] \right)_{2^{\ell-1}+j} - x^{2^{\ell-1}+j} \right| \\
&\overset{\text{Thm. 30}}{\leq} \sum_{\ell=1}^{L} \varepsilon_{\ell,k} \sum_{j=1}^{2^{\ell-1}} |a_{2^{\ell-1}+j}| \\
&\leq \varepsilon_{L,k} \sum_{\ell=1}^{L} \sum_{j=1}^{2^{\ell-1}} |a_{2^{\ell-1}+j}| \\
&= \|a\|_1 8 \cdot 2^L D^{2^L} 4^{-\frac{1}{2L} 2^k},
\end{aligned}
\tag{66}
$$

where $\varepsilon_{\ell,k}$ was defined in Lemma 28. Next, set $B = \max_{x\in[-D,D]}\left|\sum_{i=0}^{N}a_ix^i\right|$ and apply Lemma 31 with this $B$ and $\mathcal{R} = \widetilde{\mathcal{R}_a^\pi}$ to obtain an RNN $\mathcal{R}_a^\pi$ satisfying, for $k \geq 2$,

$$(\mathcal{R}_a^\pi \mathcal{D}x)[2^k - 1 + \ell] = \mathcal{C}((\widetilde{\mathcal{R}_a^\pi}\mathcal{D}x)[2^k - 2], -B, B), \qquad \ell \in \{0,\ldots,2^k - 1\}. \tag{67}$$

Combining this with Lemma 32, we compute

$$\left|(\mathcal{R}_a^\pi \mathcal{D}x)[t] - \sum_{i=0}^{N}a_ix^i\right| = \left|(\mathcal{R}_a^\pi \mathcal{D}x)[2^{\widetilde{k}(t)} - 1 + \widetilde{\ell}(t)] - \sum_{i=0}^{N}a_ix^i\right|$$

$$\overset{(67)}{=} \left|\mathcal{C}((\widetilde{\mathcal{R}_a^\pi}\mathcal{D}x)[2^{\widetilde{k}(t)} - 2], -B, B) - \sum_{i=0}^{N}a_ix^i\right|$$

$$\leq \left|(\widetilde{\mathcal{R}_a^\pi}\mathcal{D}x)[2^{\widetilde{k}(t)} - 2] - \sum_{i=0}^{N}a_ix^i\right|.$$

Since $t \geq 16\log(N)$, we have

$$\widetilde{k}(t) = \lfloor \log(t+1) \rfloor \geq \lfloor \log(16\log(N)) \rfloor = \lfloor \log(2\log(N)) \rfloor + 3$$
$$\geq \lceil \log(\log(N) + \log(N)) \rceil + 2 \geq \lceil \log(\lceil \log(N) \rceil) \rceil + 2 = \lceil \log(L) \rceil + 2,$$

and therefore (66) applies. Hence,

$$\left|(\mathcal{R}_a^\pi \mathcal{D}x)[t] - \sum_{i=0}^{N}a_ix^i\right| \leq \|a\|_1 8 \cdot 2^L D^{2^L} 4^{-\frac{1}{2L}2^{\widetilde{k}(t)}}. \tag{68}$$

We further upper bound the RHS of (68) according to

$$\|a\|_1 8 \cdot 2^L D^{2^L} 4^{-\frac{1}{2L}2^{\widetilde{k}(t)}}$$

$$= \|a\|_1 8 \cdot 2^{\lceil \log(N) \rceil} D^{2^{\lceil \log(N) \rceil}} 4^{-\frac{1}{2\lceil \log(N) \rceil}2^{\widetilde{k}(t)}}$$

$$\leq \|a\|_1 8 \cdot 2N D^{2N} 4^{-\frac{1}{2\lceil \log(N) \rceil}2^{\widetilde{k}(t)}}$$

$$= \|a\|_1 16N D^{2N} 4^{-\frac{1}{2\lceil \log(N) \rceil}2^{\lfloor \log(t+1) \rfloor}}$$

$$\leq \|a\|_1 16N D^{2N} 4^{-\frac{1}{2\lceil \log(N) \rceil}2^{\log(t+1)-1}}$$

$$\leq \|a\|_1 16N D^{2N} 4^{-\frac{t}{4\lceil \log(N) \rceil}}.$$

The proof is finalized by noting that

$$\mathcal{M}_{\mathrm{hid}}(\mathcal{R}_a^\pi) \overset{\text{Lem. 31}}{=} \mathcal{M}_{\mathrm{hid}}(\widetilde{\mathcal{R}_a^\pi}) + 11 = \mathcal{M}_{\mathrm{hid}}(\mathcal{R}_{D,L}^\pi) + 11 \overset{\text{Lem. 26}}{\leq} 40 \cdot 2^L + 11 \leq 80N + 11.$$

$\square$

# 6 Conclusion and Outlook

This paper introduced a new recurrent neural network approximation paradigm based on a reversal of the classical quantifier structure. Specifically, for a fixed target function, a single RNN with fixed topology and fixed weights achieves any prescribed error tolerance by running for sufficiently many time steps. Applied to univariate polynomials, this approach yields RNNs whose size is independent of the error tolerance and whose hidden-state dimension grows linearly with the polynomial degree.

Several directions for further work follow from this new approximation paradigm. An extension to multivariate polynomials is a natural next step and appears readily achievable. Extensions to more general smooth or continuous functions appear to be more involved, yet the methods developed here provide a reasonable starting point. It would also be interesting to examine whether the new paradigm can lead to metric-entropy optimality in the approximation of function classes. Finally, the analysis of RNNs with quantized or reduced-precision weights constitutes another interesting research direction.

## Acknowledgments

## References

[1] G. Cybenko, "Approximation by superpositions of a sigmoidal function," *Mathematics of Control, Signals and Systems*, vol. 2, no. 4, pp. 303–314, 1989, ISSN: 1435-568X. DOI: 10.1007/BF02551274

[2] K.-I. Funahashi, "On the approximate realization of continuous mappings by neural networks," *Neural Networks*, vol. 2, no. 3, pp. 183–192, 1989, ISSN: 0893-6080. DOI: 10.1016/0893-6080(89)90003-8

[3] K. Hornik, M. Stinchcombe, and H. White, "Multilayer feedforward networks are universal approximators," *Neural Networks*, vol. 2, no. 5, pp. 359–366, 1989, ISSN: 0893-6080. DOI: 10.1016/0893-6080(89)90020-8

[4] 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, ISSN: 1557-9654. DOI: 10.1109/18.256500

[5] A. R. Barron, "Approximation and estimation bounds for artificial neural networks," *Machine Learning*, vol. 14, no. 1, pp. 115–133, Jan. 1994, ISSN: 1573-0565. DOI: 10.1007/BF00993164

[6] D. Yarotsky, "Error bounds for approximations with deep ReLU networks," *Neural Networks*, vol. 94, pp. 103–114, Oct. 2017, ISSN: 0893-6080. DOI: 10.1016/j.neunet.2017.07.002

[7] M. Telgarsky, "Neural networks and rational functions," in *Proceedings of the 34th International Conference on Machine Learning*, PMLR, Jul. 2017, pp. 3387–3393. Accessed: Sep. 9, 2025.

[8] J. Schmidt-Hieber, "The Kolmogorov–Arnold representation theorem revisited," *Neural Networks*, vol. 137, pp. 119–126, May 2021, ISSN: 0893-6080. DOI: 10.1016/j.neunet.2021.01.020 Accessed: Jun. 24, 2024.

[9] J. W. Siegel, "Optimal approximation rates for deep ReLU neural networks on Sobolev and Besov spaces," *Journal of Machine Learning Research*, vol. 24, no. 357, pp. 1–52, 2023.

[10] D. Elbrächter, D. Perekrestenko, P. Grohs, and H. Bölcskei, "Deep neural network approximation theory," *IEEE Transactions on Information Theory*, vol. 67, no. 5, pp. 2581–2623, May 2021, ISSN: 0018-9448, 1557-9654. DOI: 10.1109/TIT.2021.3062161

[11] J. L. Elman, "Finding structure in time," *Cognitive Science*, vol. 14, no. 2, pp. 179–211, 1990. DOI: 10.1207/s15516709cog1402_1

[12] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016. [Online]. Available: www.deeplearningbook.org

[13] J. Bohn and M. Feischl, "Recurrent neural networks as optimal mesh refinement strategies," *Computers & Mathematics with Applications*, vol. 97, pp. 61–76, Sep. 2021, ISSN: 08981221. DOI: 10.1016/j.camwa.2021.05.018 Accessed: Sep. 9, 2025.