

Analog Computation Via Neural Networks

Hava T. Siegelmann

Department of Computer Science
Rutgers University
New Brunswick, NJ 08903

Eduardo D. Sontag

Department of Mathematics
Rutgers University
New Brunswick, NJ 08903

Abstract

We pursue a particular approach to analog computation, based on dynamical systems of the type used in neural networks research. Our systems have a fixed structure, invariant in time, corresponding to an unchanging number of "neurons". If allowed exponential time for computation, they turn out to have unbounded power. However, under polynomial-time constraints there are limits on their capabilities, though being more powerful than Turing Machines. (A similar but more restricted model was shown to be polynomial-time equivalent to classical digital computation in the previous work [17].) We note that these networks are not likely to solve polynomially NP-hard problems, as the equality "P = NP" in our model implies the almost complete collapse of the standard polynomial hierarchy. In contrast to classical computational models, the models studied here exhibit at least some robustness with respect to noise and implementation errors.

1 Introduction

"Neural networks" have attracted much attention lately as models of analog computation. Such nets consist of a finite number of simple processors, each of which computes a scalar—real-valued, not binary—function of an integrated input. This scalar function, or "activation," is meant to reflect the graded response of biological neurons to the net sum of excitatory and inhibitory inputs affecting them. The existence of feedback loops in the interconnection graph gives rise to a dynamical system. In this paper, we introduce a mathematical model for such recurrent neural networks, and we study their computational abilities.

1.1 Main Results

We focus on recurrent neural networks. In these networks, the activation of each processor is updated

according to a certain type of piecewise affine function of the activations (x_j) and inputs (u_j) at the previous instant, with real coefficients—also called weights—(a_{ij}, b_{ij}, c_i). Each processor i , ($i = 1, \dots, N$) updates its state by an equation of the type

$$x_i(t+1) = \sigma \left(\sum_{j=1}^N a_{ij} x_j(t) + \sum_{j=1}^M b_{ij} u_j(t) + c_i \right) \quad (1)$$

where N is the number of processors and M is the number of external input signals. The function σ is the simplest possible "sigmoid," namely the saturated-linear function:

$$\sigma(x) := \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } 0 \leq x \leq 1 \\ 1 & \text{if } x > 1. \end{cases} \quad (2)$$

We will give later a precise definition of language acceptance for these computational models.

We prove that neural networks can recognize in polynomial time the same class of languages as those recognized by Turing Machines that consult sparse oracles in polynomial time (the class P/poly); they can recognize all languages, including of course non-computable ones, in exponential time. Furthermore, we show that almost every language requires exponential recognition time. (For simplicity, we give our main results in terms of recognition; it is also possible to provide a more general version regarding the computation of more general functions.)

The proofs of the above results will be consequences of the following equivalence. For functions $T: \mathbb{N} \rightarrow \mathbb{N}$ and $S: \mathbb{N} \rightarrow \mathbb{N}$, let $\text{NET}_R(T)$ be the class of all functions computed by neural networks in time $T(n)$ —that is, recognition of strings of length n is in time at most $T(n)$ —and let $\text{CIRCUIT}(S)$ the class of functions computed by non-uniform families of circuits of size $S(n)$ —that is, circuits for input vectors of length n have size at most $S(n)$. We show that if F is so that $F(n) \geq n$, then $\text{NET}_R(F(n)) \subseteq \text{CIRCUIT}(\text{Poly}(F(n)))$ and

CIRCUIT $(F(n)) \subseteq \text{NET}_R(\text{Poly}(F(n)))$. This equivalence will allow us to make use of results from the theory of (nonuniform) circuit complexity.

We show that if one allows multiplications in addition to only linear operations in each neuron, that is, if one considers instead what are often called high order neural nets, the computational power does not increase. Even further, and perhaps more surprising, no increase in computational power (up to polynomial time) can be achieved by letting the activation function be not necessarily the simple saturated linear one in equation (2), but any function which satisfies certain reasonable assumptions. Also, no increase results even if the activation functions are not necessarily identical in the different processors.

One might ask about using such analog models, maybe high order nets, to "solve" NP-hard problems in polynomial time. We introduce a nondeterministic model and show that the equality $P = NP$ in the nets model is very not likely as it would imply the collapse of the polynomial hierarchy to Σ_2 .

The models used here have a weak property of "robustness" to noise and to implementation error, in the sense that small enough changes in the network would not affect the computation. The robustness includes changes in the precise form of the activation function, in the weights of the network, and even an error in the update. In classical models of (digital) computation, this type of robustness can not even be properly defined.

A Previous Related Result

In our previous work [17], we showed that if one restricts to nets all whose interconnection weights are rational numbers, which we call rational nets, then one obtains a model of computation that is polynomially related to Turing Machines. More precisely, given any multi-tape Turing Machine, one can simulate it in real time by some network with rational weights, and of course the converse simulation in polynomial time is obvious. Here we are interested in the case when weights are arbitrary real numbers. (It turns out that, as far as the results given here, the existence of just one irrational weight is all that is needed.)

1.2 The Model

The model we work with is that described by an iteration equation such as (1). For notational simplicity, we often summarize this equation, writing " $x^+(t)$ " instead of " $x(t+1)$ " and then dropping arguments t ; we

also write this in vector form, as

$$x^+ = \sigma(Ax + Bu + c) \quad (3)$$

where x is now a vector of size $N =$ number of processors, u is a vector of size $M =$ number of inputs, c is an N -vector, and A and B are, respectively, real matrices of sizes $N \times N$ and $N \times M$. (Now σ denotes application of σ into each coordinate of x .) Of course, one can drop the vector c from this description at the cost of adding a coordinate $x_0 \equiv 1$, but it is often useful to have c explicitly, and this allows us to take initial states to be $x = 0$, which corresponds to the intuitive idea that the system is at rest before the first input appears.

As part of the description, we assume that we have singled out a subset of the N processors, say x_{i_1}, \dots, x_{i_p} ; these are the p output processors, and they are used to communicate the outputs of the network to the environment. Thus a net is specified by the data (A, B, c) together with a subset of its nodes.

In our further development, both input and output channels will be forced to carry only binary data. Input and output are streams, that is, one input letter is transferred at each time (via M binary lines) and one output letter is produced at a time (and appears in the output via p binary lines). As opposed to the I/O, the computations inside the network will in general involve continuous real values.

We call a system defined by equations such as (3) simply a network or processor network. In the neural network literature, these are called recurrent first-order neural nets. We show later that considering higher-order nets, those in which multiplications of activations and/or inputs are allowed, does not result in any gain in computational capabilities (up to polynomial time).

The Finite Structure

We should emphasize from the outset that our networks are built up of finitely many processors, whose number does not increase with the length of the input. There is a small number of input channels (just two in our main result), into which inputs get presented sequentially. We assume that the structure of the network, including the values of the interconnection weights, does not change in time but rather remains constant. What changes in time are the activation values, or outputs of each processor, which are used in the next iteration. (A synchronous update model is used.) In this sense our model is very "uniform" in contrast with certain models used in the past, including those used in [9] or in the cellular automata

literature, which allow the number of units to increase over time and often even the structure to change depending on the length of inputs being presented.

The Meaning of Real Weights

One may ask about the meaning of real weights. In response, we recall that our intention is to model systems in which certain real numbers —corresponding to values of resistances, capacitances, physical constants, and so forth— may not be directly measurable, indeed may not even be computable real numbers, but they affect the “macroscopic” behavior of the system. For instance, imagine a spring/mass system. The dynamical behavior of this system is influenced by several real valued constants, such as stiffness and friction coefficients. On any finite time interval, one could replace these constants by rational numbers, and the same qualitative behavior is observed, but the long-term characteristics of the system depend on the true values. We take this use of real numbers as a basic feature of analog computation. (Another characteristic would be the use of differential as opposed to difference equations, but technical difficulties make that further study harder, and we will defer it to future work.)

What is interesting is to find a class of such systems which on the one hand is rich enough to exhibit behavior that is not captured by digital computation, while still being amenable to useful theoretical analysis, and in particular so that the imposition of resource constraints results in nontrivial reduction of computational power. That this is in accordance with models currently used in neural net studies, is especially attractive.

The remainder of this paper is organized as follows: Section 2 includes the basic definitions of networks and circuits, and states the main theorem regarding the relationships between these two models. Sections 3 and 4 contain the proof of this theorem: Section 3 shows that $\text{CIRCUIT}(F(n)) \subseteq \text{NET}_R(\text{Poly}(F(n)))$, and section 4 proves that $\text{NET}_R(F(n)) \subseteq \text{CIRCUIT}(\text{Poly}(F(n)))$. Section 5 states some corollaries for neural networks which follow from the above relation with circuits. We also define there a notion of nondeterministic network. In section 6, we show that our model does not gain power if one lets each neuron compute a polynomial function —rather than just affine combinations— of the activations of all the neurons and the external inputs, or by allowing more general activation functions than the piecewise linear one.

We now turn to precise definitions.

2 Basic Definitions

As we discussed above, we consider synchronous networks which can be represented as dynamical systems whose state at each instant is a real vector $x(t) \in \mathbb{R}^N$. The i th coordinate of this vector represents the activation value of the i th processor at time t . In matrix form, the equations are as in (3), for suitable matrices A, B and vector c .

Given a system of equations such as (3), an initial state $x(1)$, and an infinite input sequence $u = u(1), u(2), \dots$, we can define iteratively the state $x(t)$ at time t , for each integer $t \geq 1$, as the value obtained by recursively solving the equations. This gives rise, in turn, to a sequence of output values, by restricting attention to the output processors; we refer to this sequence as the “output produced by the input u ” starting from the given initial state.

2.1 Recognizing Languages

To define what we mean by a net recognizing a language $L \subseteq \{0, 1\}^+$, we must first define a formal network, a network which adheres to a rigid encoding of its input and output. We proceed as in [17] and define formal nets with two binary input lines. The first of these is a data line, and it is used to carry a binary input signal; when no signal is present, it defaults to zero. The second is the validation line, and it indicates when the data line is active; it takes the value “1” while the input is present there and “0” thereafter. We use “ D ” and “ V ” to denote the contents of these two lines, respectively, so

$$u(t) = (D(t), V(t)) \in \{0, 1\}^2$$

for each t . We always take the initial state $x(1)$ to be zero and to be an equilibrium state, that is, $\sigma(A0 + B0 + c) = 0$. We assume that there are two output processors, which also take the role of data and validation lines and are denoted $O_d(t), O_v(t)$ respectively.

(The convention of using two input lines allows us to have all external signals be binary; of course many other conventions are possible and would give rise to the same results, for instance, one could use a three-valued input, say with values $\{-1, 0, 1\}$, where “0” indicates that no signal is present, and ± 1 are the two possible binary input values.)

We now encode each word $\alpha = \alpha_1 \dots \alpha_k \in \{0, 1\}^+$ as follows. Let $u_\alpha(t) = (V_\alpha(t), D_\alpha(t))$, $t = 1, \dots$, where

$$V_\alpha(t) = \begin{cases} 1 & \text{if } t = 1, \dots, k \\ 0 & \text{otherwise,} \end{cases}$$

and

$$D_\alpha(t) = \begin{cases} \alpha_k & \text{if } t = 1, \dots, k \\ 0 & \text{otherwise.} \end{cases}$$

Given a formal net \mathcal{N} , with two inputs as above, we say that a word α is classified in time τ , if the following property holds: the output sequence

$$y(t) = (O_d(t), O_v(t))$$

produced by u_α when starting from $x(1) = 0$ has the form $O_d = \underbrace{0 \dots 0}_{\tau-1} \eta_\alpha 000 \dots$, $O_v = \underbrace{0 \dots 0}_{\tau-1} 1000 \dots$,

where $\eta_\alpha = 0$ or 1 .

Let $T : \mathbb{N} \rightarrow \mathbb{N}$ be a function on natural numbers. We say that the language $L \subseteq \{0, 1\}^+$ is recognized in time T by the formal net \mathcal{N} provided that each word $\alpha \in \{0, 1\}^+$ is classified in time $\tau \leq T(|\alpha|)$, and η_α equals 1 when $\alpha \in L$ and is $= 0$ otherwise.

2.2 Circuit Families

We briefly recall some of the basic definitions of non-uniform families of circuits. A Boolean circuit is a directed acyclic graph. Its nodes of in-degree 0 are called input nodes, while the rest are called gates and are labeled by one of the Boolean functions AND, OR, or NOT (the first two seen as functions of many variables, the last one as a unary function). One of the nodes, which has no outgoing edges, is designated as the output node. The size of the circuit is the total number of gates. Adding if necessary extra gates, we assume that nodes are arranged into levels $0, 1, \dots, d$, where the input nodes are at level zero, the output node is at level d , and each node only has incoming edges from the previous level. The depth of the circuit is d , and its width is the maximum size of each level. Each gate computes the corresponding Boolean function of the values from the previous level, and the value obtained is considered as an input to be used by the successive level; in this fashion each circuit computes a Boolean function of the inputs.

A family of circuits \mathcal{C} is a set of circuits $\{c_n, n \in \mathbb{N}\}$. These have sizes $S_{\mathcal{C}}(n)$, depth $D_{\mathcal{C}}(n)$, and width $W_{\mathcal{C}}(n)$, $n = 1, 2, \dots$, which are assumed to be monotone nondecreasing functions. If $L \subseteq \{0, 1\}^+$, we say that the language L is computed by the family \mathcal{C} if the characteristic function of $L \cap \{0, 1\}^n$ is computed by c_n , for each $n \in \mathbb{N}$.

The qualifier "nonuniform" serves as a reminder that there is no requirement that circuit families be recursively described. It is this lack of classical computability that makes circuits a possible model of resource-bounded "computing," as emphasized in [14].

We will show that recurrent neural networks, although more "uniform" in the sense that they have an unchanging physical structure, share exactly the same power.

If L is recognized by the formal net \mathcal{N} in time T , we write $\phi_{\mathcal{N}} = L$ and $T_{\mathcal{N}} = T$. If L is computed by the family of circuits \mathcal{C} , we write $\phi_{\mathcal{C}} = L$. We are interested in comparing the functions $T_{\mathcal{N}}$ and $S_{\mathcal{C}}$ for formal nets and circuits so that $\phi_{\mathcal{N}} = \phi_{\mathcal{C}}$.

2.3 Statement Of Result

Recall that $\text{NET}_R(T(n))$ is the class of languages recognized by formal networks (with real weights) in time $T(n)$ and that $\text{CIRCUIT}(S(n))$ is the class of languages recognized by (non-uniform) families of circuits of size $S(n)$.

Theorem 1 Let F be so that $F(n) \geq n$. Then, $\text{NET}_R(F(n)) \subseteq \text{CIRCUIT}(\text{Poly}(F(n)))$, and $\text{CIRCUIT}(F(n)) \subseteq \text{NET}_R(\text{Poly}(F(n)))$. ■

More precisely, we prove the following two facts. For each function $F(n) \geq n$:

- $\text{CIRCUIT}(F(n)) \subseteq \text{NET}_R(nF^2(n))$.
- $\text{NET}_R(F(n)) \subseteq \text{CIRCUIT}(F^3(n))$.

3 Circuit Families Are Simulated By Networks

We start by reducing circuit families to networks. The proof will construct a fixed, "universal" net, having roughly $N = 1000$ processors, which, through the setting of a particular real weight which encodes an entire circuit family, can simulate that family.

Theorem 2 There exists a positive integer N such that the following property holds: For each circuit family \mathcal{C} of size $S_{\mathcal{C}}(n)$ there exists an N -processor formal network $\mathcal{N} = \mathcal{N}(\mathcal{C})$ so that $\phi_{\mathcal{N}} = \phi_{\mathcal{C}}$ and $T_{\mathcal{N}}(n) = O(nS_{\mathcal{C}}^2(n))$.

The proof is provided in the remainder of this section.

3.1 The circuit Encoding

Given a circuit c_i —with size s_i and width w_i , we encode it as a sequence over the alphabet $\{0, 2, 4, 6\}$

called $en[c_i]$. This sequence corresponds to a concatenated encoding of its gates provided in a bottom-up manner.

We encode a non-uniform family of circuits, \mathcal{C} as an infinite sequence

$$e(\mathcal{C}) = 8 \overline{en}[c_1] 8 \overline{en}[c_2] 8 \overline{en}[c_3] \dots, \quad (4)$$

where $\overline{en}[c_i]$ is the encoding of c_i in the reversed order. This allows for a decoding procedure that reveals the gates in a bottom up manner, thus allowing for a quick simulation of the circuit.

Let $\hat{\mathcal{C}}$ be the interpretation of formula (4) in base 9. That is,

$$\hat{\mathcal{C}} = 8 \overline{en}[c_1] 8 \overline{en}[c_2] 8 \overline{en}[c_3] \dots |_9 \equiv \sum_{i=1}^{\infty} \frac{r_i}{9^i}, \quad (5)$$

where r_i is the i th bit of $e(\mathcal{C})$.

A Possible Encoding

Given a circuit c —with size s , width w , and w_i gates in the i th level—we encode it as a finite sequence over the alphabet $\{0, 2, 4, 6\}$, as follows:

- The encoding of each level i starts with the letter 6. Levels are encoded successively, starting with the bottom level and ending with the top one.
- At each level, gates are encoded successively. The encoding of a gate g consists of three parts—a starting symbol, a 2-digit code for the gate type, and a code to indicate which gate feeds into it:
 - It starts with the letter 0.
 - A two digit sequence $\{42, 44, 22\}$ denotes the type of the gate, $\{AND, OR, NOT\}$ respectively.
 - If gate g is in level i , then the input to g is represented as a sequence in $\{2, 4\}^{w_{i-1}}$, such that the j th position in the sequence is 4 if and only if the j th gate of the $(i-1)$ th level feeds into gate g .

The encoding of a gate g in level i is of length $(w_{i-1} + 3)$. The length of the encoding of a circuit c is $l(c) \equiv |en(c)| = O(sw)$.

Cantor Like Set Encoding

The number $\hat{\mathcal{C}}$ which encodes a family of circuits, or one that is a suffix of such an encoding, is a number between $[0, 1]$. However, not every value in $[0, 1]$

appears. The set of possible values is not continuous and has “holes”. Such a set of values “with holes” is a Cantor set. Its self-similar structure means that bit (base 9) shifts preserve the “holes.”

The advantage of this approach is that there is never a need to distinguish among two very close numbers in order to read the desired circuit out of the encoding; the circuit can be then retrieved with finite-precision operations employing a finite number of neurons.

In the proof, we exhibit a network having as one of its weights $\hat{\mathcal{C}}$ (all other weights are rational numbers), which upon receiving an input α of size n , computes (i.e. retrieves) $en[c_i]$. Then the network simulates the operation of circuit c_i on the input α step by step.

3.2 A Circuit Retrieval

Lemma 3.1 For each (non-uniform) family of circuits \mathcal{C} there exists a 16-processor network $\mathcal{N}_R(\mathcal{C})$ with one input line such that, starting from the zero initial state and given the input signal

$$u(1) = \underbrace{11 \dots 1}_n 00 \dots |_2 = 1 - 2^{-n},$$

$$u(t) = 0 \quad \text{for } t > 1,$$

$\mathcal{N}_R(\mathcal{C})$ outputs

$$x_r = \underbrace{000 \dots 0}_{2n+2 \sum_{i=1}^n l(c_i)+4} \overline{en}[c_n] 000 \dots$$

Proof. Let $\Sigma = \{0, 2, 4, 6, 8\}$. Denote by \mathcal{C}_9 the “Cantor 9-set,” which consists of all those real numbers q which admit an expansion of the form

$$q = \sum_{i=1}^{\infty} \frac{\alpha_i}{9^i} \quad (6)$$

with each $\alpha_i \in \Sigma$. Let $\Lambda : \mathbb{R} \rightarrow [0, 1]$ be the function

$$\Lambda[x] := \begin{cases} 0 & \text{if } x < 0 \\ 9x - [9x] & \text{if } 0 \leq x \leq 1 \\ 1 & \text{if } x > 1. \end{cases} \quad (7)$$

Let $\Xi : \mathbb{R} \rightarrow [0, 1]$ be the function

$$\Xi[x] := \begin{cases} 0 & \text{if } x < 0 \\ 2 \lfloor \frac{9x}{2} \rfloor & \text{if } 0 \leq x \leq 1 \\ 1 & \text{if } x > 1. \end{cases} \quad (8)$$

Note that, for each

$$q = \sum_{i=1}^{\infty} \alpha_i / 9^i \in \mathcal{C}_9,$$

we may think of $\Xi[q]$ as the "select left" operation, since

$$\Xi[q] = \alpha_1,$$

and of $\Lambda[q]$ as the "shift left" operation, since

$$\Lambda[q] = \sum_{i=1}^{\infty} \alpha_{i+1}/9^i \in C_9.$$

For each $i \geq 0$, $q \in C_9$, $\Xi[\Lambda^i[q]] = \alpha_{i+1}$. The following procedure summarizes the task to be performed by the network constructed below, which in turn satisfies the requirements of the lemma.

```

Procedure Retrieval( $\hat{C}, n$ )
Variables counter,  $y$ ,  $z$ 
Begin
  counter  $\leftarrow 0$ ,  $y \leftarrow 0$ ,  $z \leftarrow \hat{C}$ ,
  While counter  $< n$ 
    Parbegin
       $z \leftarrow \Lambda[z]$ 
      if  $\Xi[z] = 8$  then increment counter
    Parend,
    While  $\Xi[z] < 8$ 
      Parbegin
         $z \leftarrow \Lambda[z]$ 
         $y \leftarrow \frac{1}{9}(y + \Xi[z])$ 
      Parend,
    Return( $y$ )
End

```

The functions Λ and Ξ can not be programmed within the neural network model due to their discontinuity. However, we can program the functions $\tilde{\Lambda}, \tilde{\Xi}$, which coincide with Λ, Ξ respectively on C_9 :

$$\tilde{\Lambda}[q] = \sum_{j=0}^8 (-1)^j \sigma(9q - j), \quad (9)$$

and

$$\tilde{\Xi}[q] = 2 \sum_{j=0}^3 \sigma(9q - (2j + 1)). \quad (10)$$

It is easy to provide a network of 16 processors that executes the above procedure for the substitute functions $\tilde{\Lambda}$ and $\tilde{\Xi}$. ■

3.3 Circuit Simulation By A Network

Let $\alpha \in \{0, 1\}^n$ be a binary sequence. Denote by $en[\alpha]$ the sequence $\in \{2, 4\}^n$ that substitutes $(2\alpha_i + 2)$ for each α_i , and by $\widehat{en}[\alpha]$ the interpretation of $en[\alpha]$ in base 9, that is, $en[\alpha]_9$. We next construct a "universal net" for interpreting circuits.

Lemma 3.2 There exists a network \mathcal{N}_s , such that for each circuit c and binary sequence α , starting from the zero initial state and applying the input signal

$$u_1 = \widehat{en}[c]00 \dots \quad u_2 = \widehat{en}[\alpha]00 \dots,$$

\mathcal{N}_s outputs

$$x_0 = \underbrace{00 \dots 0}_T y 00 \dots \quad x_v = \underbrace{00 \dots 0}_T 100 \dots,$$

where y is the response of circuit c on the input α , and $T = O(l(c) + |\alpha|)$.

Proof. It is easy to verify that, given any circuit description with gates ordered bottom up, there is a three-tape Turing Machine which can simulate the given circuit in time $O(l(c) + |\alpha|)$. Indeed, we proved in ([17]) that if M is a k -tape Turing Machine with s states which computes in time T a function f on binary input strings, then there exists a rational network \mathcal{N} , which consists of $9^k s + s + 28k + 2$ processors, that computes the same function f in time $O(T)$. Closer counting shows that less than 1000 processors suffice. ■

Remark 3.3 If the lemma would only require an estimate of a polynomial number of processors, as opposed to the more precise estimate that we obtain, the proof would have been immediate from the consideration of the circuit value problem (CVP). This is the problem of recognizing the set of all pairs $\langle x, y \rangle$, where $x \in \{0, 1\}^+$, and y encodes a circuit with $|x|$ input lines which outputs 1 on input x . It is known that $CVP \in P$ ([3] volume I, pg 110). □

3.4 The General Proof

Proof of Theorem 2.

Let \mathcal{C} be a circuit family. We construct the required formal network as a composition of the following three networks:

- An input network, \mathcal{N}_I , which receives the input

$$u_1 = \alpha 00 \dots$$

$$u_2 = \underbrace{11 \dots 1}_{|\alpha|} 00 \dots,$$

and computes $\widehat{en}[\alpha]$ and $u_2|_2$, for each $\alpha \in \{0, 1\}^+$. This network is trivial to implement.

- A retrieval network, $\mathcal{N}_R(c)$, as described in Lemma 3.1, which receives $u_2|_2$ from \mathcal{N}_I , and computes $\widehat{en}[c|_{|\alpha|}]$. (Note that during the encoding operation, network \mathcal{N}_I produces an output of zero, and $\mathcal{N}_R(c)$ remains in its initial state 0.)

- A simulation network, \mathcal{N}_S , as stated in Lemma 3.2, which receives $\widehat{\text{en}}[c_{|\alpha|}]$ and $\widehat{\text{en}}[\alpha]$, and computes

$$\begin{aligned} x_0 &= \underbrace{00 \cdots 0}_T \phi_c(\alpha) 00 \cdots \\ x_v &= \underbrace{00 \cdots 0}_T 100 \cdots \end{aligned}$$

Notice that out of the above three networks, only \mathcal{N}_R depends on the specific family of circuits \mathcal{C} . Moreover, all weights can be taken to be rational numbers, except for the one weight that encodes the entire circuit family.

The time complexity to compute the response of \mathcal{C} to the input α is dominated by that of retrieving the circuit description. Thus, the complexity is of order $T = O\left(\sum_{i=1}^{|\alpha|} l(c_i)\right)$. We remarked that the length of the encoding $l(c_i)$ is of order $O(W_C(i)S_C(i))$, which is itself $O(S_C^2(i))$. Since $S_C(i) \leq S_C(i+1)$ for $i = 1, 2, \dots$, we achieve the claimed bound $T = O(|\alpha| S_C^2(|\alpha|))$.

Remark 3.4 In case of bounded fan-in, the “standard encoding” of circuit c_n is of length $l(c_n) = O(S_C(n) \log(S_C(n)))$. The total running time of the algorithm is then $O(n S_C(n) \log(S_C(n)))$. \square

4 Networks Are Simulated By Circuit Families

We next state the reverse simulation, of nets by nonuniform families of circuits.

Theorem 3 Let \mathcal{N} be a formal network that computes in time $T : \mathbb{N} \rightarrow \mathbb{N}$. There exists a non-uniform family of circuits $\mathcal{C}(\mathcal{N})$ of size $O(T^3)$, depth $O(T \log(T))$, and width $O(T^2)$, that accepts the same language as \mathcal{N} does. \blacksquare

The proof is given in the next two subsections. In the first part, we replace a single formal network by a family of formal networks with small rational weights. (This is unrelated to the standard fact for threshold gates that weights can be taken to have $n \log n$ bits.) In the second part, we simulate such a family of formal networks by circuits.

4.1 Linear Precision Suffices

Define a processor to be a designated output processor if its activation value is used as an output of the

network (i.e. it is an output processor) and is not fed into any other processor. A formal network, for which its two output processors are designated, is called an output designated network. Its processors, which are not the designated output processors, are called internal processors.

For the next result, we introduce the notion of a q -truncation net. This is a processor network in which the update equations take the form

$$x_i^{\dagger} = q\text{-Truncation} \left[\sigma \left(\sum_{j=1}^N a_{ij} x_j + \sum_{j=1}^M b_{ij} u_j + c_i \right) \right],$$

where q -Truncation means the operation of truncating after q bits.

Lemma 4.1 Let \mathcal{N} be an output designated network. If \mathcal{N} computes in time T , there exists a family of $T(n)$ -Truncation output designated networks $\mathcal{N}_1(n)$ such that

- For each n , $\mathcal{N}_1(n)$ has the same number of processors and input and output channels as \mathcal{N} does.
- The weights feeding into the internal processors of $\mathcal{N}_1(n)$ are like those of \mathcal{N} , but truncated after $O(T(n))$ bits.
- For each designated output processor in \mathcal{N} , if this processor computes $x_i^{\dagger} = \sigma(f)$, where f is a linear function of processors and inputs, then the respective processor in $\mathcal{N}_1(n)$ computes $\sigma(2\tilde{f} - .5)$, where \tilde{f} is the same as the linear function f but applied instead to the processors of $\mathcal{N}_1(n)$ and with weights truncated at $O(T(n))$ bits.
- The respective output processors of \mathcal{N} and $\mathcal{N}_1(n)$ have the same activation values at all times $t \leq T(n)$.

Proof. We first measure the difference (error) between the activations of the corresponding internal processors of $\mathcal{N}_1(n)$ and \mathcal{N} at time $t \leq T(n)$. This calculation is analogous to that of the chop error in floating point computation, [2]. We use the following notation: The network has N processors and M input lines. We denote by L the value $(N + M + 1)$, and by W a bound on the sum of weights. We denote by $\tilde{x}_i(t)$, \tilde{a}_{ij} , \tilde{b}_{ij} , and \tilde{c}_i the respective activation values of processors, and weights of $\mathcal{N}_1(n)$. The errors considered include: $\delta_w \in (0, 1)$ and $\delta_p > 0$ are the truncation errors at weights and processors, respectively; and $\epsilon_t > 0$ is the largest accumulated error at time t in processors of $\mathcal{N}_1(n)$. Network $\mathcal{N}_1(n)$ computes at each step

$$\tilde{x}_i^{\dagger} = q\text{-Truncation} \left[\sigma \left(\sum_{j=1}^N \tilde{a}_{ij} \tilde{x}_j + \sum_{i=1}^M \tilde{b}_{ij} u_j + \tilde{c}_i \right) \right].$$

We assume inductively on t that for all internal processors i, j , $|\bar{x}_i(t) - x_i(t)| \leq \epsilon_t$. Using the global Lipschitz property $|\sigma(a) - \sigma(b)| \leq |a - b|$, it follows that

$$\begin{aligned} \epsilon_t &\leq N(W' + \delta_w)\epsilon_{t-1} + (N + M + 1)\delta_w + \delta_p \\ &\leq LW\epsilon_{t-1} + L\delta_w + \delta_p. \end{aligned}$$

Therefore,

$$\epsilon_t \leq \sum_{i=0}^{t-1} (LW)^i (L\delta_w + \delta_p) \leq (LW)^t (L\delta_w + \delta_p).$$

We now analyze the behavior of the output processors. We need to prove that $\sigma(2\bar{f} - .5) = 0, 1$ when $\sigma(f) = 0, 1$ respectively. That is, $f \leq 0 \implies \bar{f} < \frac{1}{4}$ and $f \geq 1 \implies \bar{f} > \frac{3}{4}$. This happens if $|f - \bar{f}| < \frac{1}{4}$. Arguing as earlier, the condition $\epsilon_t < \frac{1}{4}$ suffices. This is translated into the requirement $(L\delta_w + \delta_p) \leq \frac{1}{4}(LW)^{-t}$. If both δ_w and δ_p are bounded by $\frac{1}{8}(LW)^{-(t-1)}$, this inequality holds. This happens when the weights and the processor activations are truncated after $O(t \log(LW))$ bits. As L and W are constants, we conclude as desired that a sufficient truncation for a computation of length T is $O(T)$. ■

4.2 The Network Simulation by a Circuit

Lemma 4.2 Let \mathcal{N}_1 be a family of $T(n)$ -Truncation output designated networks, where all networks $\mathcal{N}_1(n)$ consist of N processors and the weights are all rational numbers with $O(T)$ bits. Then, there exists a circuit family \mathcal{C} of size $O(T^3)$, depth $O(T \log(T))$, and width $O(T^2)$, so that c_n accepts the same language as $\mathcal{N}_1(n)$ does on $\{0, 1\}^n$.

This proof is omitted due to space limitations. The crucial points are as follows. First, we change the input convention from feeding serially into the network $\mathcal{N}_1(n)$ via the two input lines data and validation (where the validation line includes n consecutive 1's) to n bits that are fed simultaneously into the circuit c_n . Then, we associate with each processor p a subcircuit $sc(p)$. Each processor $p \in \mathcal{N}_1(n)$ computes a truncated sum of up to $N + 2$ numbers, N of which are multiplications of two T -bit numbers. Hardwiring the weights, we can say that each processor computes a sum of $(TN + 2)$ $(2T)$ -bit numbers. Using the carry-look-ahead method, [16], the summation can be computed via a subcircuit of depth $O(\log(TN))$, width $O(T^2N)$, and size $O(T^2N)$. (This depth is of the same order as the lower bound of similar tasks, see [5], [7].)

The proof of Theorem 3 follows immediately from Lemma 4.1 and Lemma 4.2.

5 Corollaries

Let NET-P and NET-EXP be the classes of languages accepted by formal networks in polynomial time and exponential time, respectively. Let CIRCUIT-P and CIRCUIT-EXP be the classes of languages accepted by families of circuits in polynomial and exponential size, respectively.

Corollary 5.1 NET-P = CIRCUIT-P and NET-EXP = CIRCUIT-EXP.

The class CIRCUIT-P is often called "P/poly" and coincides with the class of languages recognized by Turing Machine "with advice sequences" in polynomial time. From [3], volume I, Theorem 5.11, pg 122 (originally, [13]), we conclude as follows:

Corollary 5.2 NET-EXP includes all possible binary languages. Furthermore, most Boolean functions require exponential time complexity.

The concept of a nondeterministic circuit family is usually defined by means of an extra input, whose role is that of an oracle. Similarly, we define a nondeterministic network to be a network having an extra binary input line, the Guess line, in addition to the Data and Validation lines. A language L accepted by a nondeterministic formal network \mathcal{N} in time B is defined as $L = \{\alpha \mid \exists \text{ a guess } \gamma, \phi_{\mathcal{N}}(\alpha, \gamma) = 1, T_{\mathcal{N}}(\alpha, \gamma) \leq B(|\alpha|)\}$.

It is easy to see that Corollary (5.1), stated for the deterministic case, holds for the nondeterministic case as well. That is, if we define NET-NP to be the class of languages accepted by nondeterministic formal networks in polynomial time, and CIRCUIT-NP to be the class of languages accepted by nondeterministic non-uniform families of circuits of polynomial size, then:

Corollary 5.3 NET-NP = CIRCUIT-NP. □

Since $NPC \subseteq \text{NET-NP}$ (one may simulate a nondeterministic Turing Machine by a nondeterministic network with rational weights), the equality NET-NP = NET-P implies $NPC \subseteq \text{CIRCUIT-P} = P/\text{poly}$. Thus, from [10] we conclude: If NET-NP = NET-P then the polynomial hierarchy collapses to Σ_2 .

The above result says that a theory of computation similar to that which arises in the classical case of Turing machine computation is also possible for our model of analog computation. In particular, even though the two models have very different power, the question of knowing if the verification of solutions to problems is really easier than finding solutions, at the core of

modern computational complexity, has a precise corresponding version in our setup, and its solution will be closely related to that of the classical case. Of course, it follows from this that it is quite likely that NET-NP is strictly more powerful than NET-P.

6 Equivalence of Different Dynamical Systems

We consider dynamical systems –which we will call generalized processor networks– with far less restrictive structure than the recurrent neural network model which was described in equation (1). We show that these networks are not more powerful, up to polynomial time slowdown, than the previously considered model.

Let N, M, p be natural numbers. A generalized processor network is a dynamical system that consists of N processors x_1, x_2, \dots, x_N , and receives its input $u_1(t), u_2(t), \dots, u_M(t)$ via M input lines. A subset of the N processors, say x_{i_1}, \dots, x_{i_p} , is the set of output processors of the system, used to communicate the output of the system to the environment. In vector form, a generalized processor network D updates its processors via the dynamic equation

$$x^+ = f(x, u),$$

where x is the current state of the network (a vector), u is an external input (also possibly a vector), and f is a composition of functions: $f = \psi \circ \pi$, where $\pi : \mathbb{R}^{N+M} \mapsto \mathbb{R}^N$ is some vector polynomial in $N+M$ variables with real coefficients, and $\psi : \mathbb{R}^N \mapsto \mathbb{R}^N$ is any vector function which has a bounded range and is locally Lipschitz. (Thus, the composite function $f = \psi \circ \pi$ again satisfies the same properties.)

We also assume, as part of the definition of generalized processor network, that, at least for binary inputs of the type considered in the definition of “formal networks,” given in section 2.1, D outputs “soft” binary information. That is, there exist two constants α, β , satisfying $\alpha < \beta$ and called the decision thresholds, so that each output neuron of D outputs a stream of numbers each of which is either smaller than α or larger than β . We interpret the outputs of each output neuron y as a binary value:

$$\text{binary}(y) = \begin{cases} 0 & \text{if } y \leq \alpha \\ 1 & \text{if } y \geq \beta. \end{cases}$$

In the usual model we studied earlier, the values are always binary, but we allow more generality to show that even if one allows more general analog values, no

increase in computational power is attained, at least up to polynomial time.

A neural network is a special case of a generalized processor network, in which all coordinates of the function ψ compute the same sigmoidal function, and the polynomial π is a first order polynomial, that is, an affine function.

Let $T : \mathbb{N} \mapsto \mathbb{N}$ be a function from integers into integers. We say that a generalized processor network D computes in time T if for every input of size $n \in \mathbb{N}$, D completes its output in no more than $T(n)$ steps.

Theorem 4 Let D be a generalized processor network which computes via a function $f = \psi \circ \pi$, where the function $T(n)$ -truncation(f) is in $P/poly$. Then there exists a neural network N_D which recognizes the same language as D and which does so with at most a polynomial time slowdown. Furthermore, if $T(n)$ -truncation(ψ) $\in P$ the weights utilized by N_D are of the same type as the coefficients of the polynomial π (rational or real, respectively).

The proof of this theorem is omitted due to space limitations. Briefly, we prove “linear precision suffices” to the generalized network similarly to the proof in subsection 4.1:

Lemma 6.1 Assume D computes in time T , with decision thresholds α, β . Then, there is a constant c such that the function

$$q(n) = cT(n)$$

satisfies the following property. For each positive integer n , let $Q = q(n)$. Then, Q -Truncation(D) computes the same function as D on inputs of length at most n , with decision thresholds

$$\alpha' = \alpha + \frac{\beta - \alpha}{3} \quad \text{and} \quad \beta' = \beta - \frac{\beta - \alpha}{3}.$$

Then, we show that if $T(n)$ -truncation(f) is in P , one can simulate D via a neural network with rational weights. If, however, $T(n)$ -truncation(f) is in $P/poly$, real weights are required. In both cases, no more than polynomial slow down in the computation occurs while simulating.

Corollary 6.2 Adding flexibility to the neural network model, described in Equation (1), does not add power to the model, except for a possible polynomial time speed up. This flexibility includes:

- Using a high order polynomial π rather than an affine function.

- Using other ψ functions rather than the saturation we used earlier, including the possibility of having different functions in different neurons.
- Allowing for the output to be “soft binary” rather than pure binary.

Note that networks with high order polynomials have appeared especially in the language recognition literature (see e.g. [8] and references there). We emphasize the relationship between these models: Let N_1 be neural network (of any order), which recognizes a language L in polynomial time. Then there is a first order network N_2 which recognizes the same language L in polynomial time.

Remark 6.3 The networks that we consider are mildly “robust to noise and to implementation error” in the sense that small enough perturbations in weights or the sigmoid activation function do not affect the computation, as long as “soft binary” outputs are considered. Given any time T , there is some ϵ_T so that an error of ϵ_T would not affect the computation up to a time T . \square

Acknowledgment

This Research was partially supported by US Air Force Grant AFOSR-91-0343. We wish to thank Richard Beigel for suggesting that we study the generalization of our results to the high-order net case treated in section 6. We are also grateful to Robert Solovay for his useful comments during the early steps of this research.

References

- [1] Alspector J., R.B. Allen, “A neuromorphic VLSI learning system,” in *Advanced Research in VLSI: Proceedings of the 1987 Stanford Conference*, (P. Loséleben ed.), MIT Press, Cambridge, MA, 1987: 313-349.
- [2] Atkinson K.E., *An Introduction to Numerical Analysis*, Wiley, New York, 1989.
- [3] Balcazar J.L., J. Diaz, J. Gabarro, *Structural Complexity*, Springer-Verlag, Berlin, 1988.
- [4] Blum L., M. Shub, and S. Smale, “On a theory of computation and complexity over the real numbers: NP completeness, recursive functions, and universal machines,” *Bull. A.M.S.* 21(1989): 1-46.
- [5] Chandra A.K., L. Stockmeyer, U. Vishkin, “Constant depth reducibility,” *SIAM J. Computing* 13(1984): 423-439.
- [6] Eberhardt S.P., T. Daud, D. A. Kerns, T. X. Brown, and A. P. Thakoor, “Competitive neural architecture for hardware solution to the assignment problem,” *Neural Networks* 4(1989): 431-442.
- [7] Furst M., J.B.Saxe, M. Sipser “Parity, circuits, and the polynomial-time hierarchy,” *Proc. 22nd IEEE Symp. Foundations of Comp. Sci.*, 1981: 260-270.
- [8] Giles, C.L., D. Chen, C.B. Miller, H.H. Chen, G.Z. Sun and Y.C. Lee, “Second-Order Recurrent Neural Networks for Grammatical Inference,” *Proceedings of the International Joint Conference on Neural Networks*, vol. 2 (1991), pp. 273-281.
- [9] Hong J.W., “On Connectionist Models,” *Comm. on Pure and Applied Mathematics* 41(1988): 1039-1050.
- [10] Karp R.M., R. Lipton, “Turing Machines that take advice,” *Enseign. Math.* 28(1982): 191-209.
- [11] Maass W., G. Schnitger, and E.D. Sontag, “On the computational power of sigmoid versus Boolean threshold circuits,” *Proc. 32nd IEEE Symp. Foundations of Comp. Sci.*, 1991: 767-776.
- [12] MacLennan B.J., “Continuous symbol systems: The logic of connectionism,” in D.S. Levine and M. Aparicio IV (eds.), *Neural Networks for Knowledge Representation and Inference*, Lawrence Erlbaum, Hillsdale, NJ, 1992.
- [13] Muller D.E., “Complexity in electronic switching circuits,” *IRE Trans. Electronic Comp.* 5(1956): 15-19.
- [14] Parberry I., “Knowledge, understanding, and computational complexity,” *Technical Report CRPDC-92-2*, Center for Research in Parallel and Distributed Computing, Department of Computer Sciences, University of North Texas, Feb. 1992.
- [15] Penrose R., *The Emperor’s New Mind*, Oxford University Press, Oxford, 1989.
- [16] Savage J.E. *The Complexity of Computing*, New York, Wiley, 1976.
- [17] Siegelmann H.T., E.D. Sontag, “On the computational power of neural nets,” in *Proc. Fifth ACM Workshop on Computational Learning Theory*, Pittsburgh, July 1992: 440-449.
- [18] Vergis A., K. Steiglitz, B. Dickinson, “The complexity of analog computation,” in *Math. and Computers in Simulation* 28(1986): 91-113.
- [19] Wolpert D., “A computationally universal field computer which is purely linear,” *Los Alamos National Laboratory report LA-UR-91-2937*.