

# Linear logic and recurrent neural networks

Huiyi Hu, Daniel Murfet

January 15, 2017

## 1 Introduction

In recent years deep neural networks have proven to be proficient at learning hierarchical vector representations of natural data, such as images [?] and language [?]. According to Leibniz [?] reasoning is the algebraic manipulation of such representations or “symbols” and it is therefore natural to look for systems that augment neural networks with the ability to perform these kinds of manipulations. There are by now numerous proposals for how to achieve this, including Neural Turing Machines [?, ?], stack-augmented recurrent neural networks [?] and others [?, ?, ?]. In this paper we investigate a new conceptual approach, based on the Curry-Howard isomorphism and the semantics of linear logic.

The Curry-Howard isomorphism [?] gives a bijection between a prototypical system of formal reasoning (Gentzen’s natural deduction) and a prototypical algorithmic system (Church’s simply-typed lambda calculus). Since the beginning of the AI field in the 70’s, lambda calculus and derivative languages like LISP have played an important role. While these methods may have been overshadowed in recent decades by the rise of the connectionist approach, now called deep learning, it has long been argued that a proper synthesis of the two approaches would be necessary in order to achieve reasoning in a connectionist system [?]. Arguably, the main obstacle to this synthesis is the discrepancy between the discrete symbolic nature of natural deduction, or equivalently lambda calculus, and the differentiable nature of neural networks. One way to overcome this obstacle is to augment a neural network not directly with symbolic logic, but rather with a *differentiable model* of logic; and the most natural way to construct such a model is not to work directly with simply-typed lambda calculus but rather with a refinement due to Girard known as linear logic [?] which has a canonical model (called a denotational semantics) in differentiable maps between vector spaces [?, ?, ?].

We demonstrate these ideas concretely in the form of an RNN controller which predicts at time  $t$  a linear combination of denotations of linear logic programs (or proofs) to be applied to vectors which are themselves predicted from the hidden state of the RNN at time  $t-1$  and the input symbol at time  $t$ . The output of the linear logic program is a linear operator, which is applied to the hidden state at time  $t-1$  to generate a vector which

is then added to the usual update equation of the RNN. This architecture borrows its basic features from various sources, including the second-order RNN [?], neural machine translation system [?], matrix-vector models in natural language processing [?] and the multiplicative RNN [?]. We apply it to the old problem [?] of approximating a function  $f : A^* \rightarrow A^*$  of sequences in some alphabet  $A$ , following recent work on neural networks augmented with algorithmic elements [?, ?, ?].

## 2 Architecture

### 2.1 The second-order RNN

As usual we use “weight” as a synonym for variable, or more precisely, the variables which we will vary during gradient descent. At time  $t \geq 1$  we denote the hidden state by  $h^{(t)}$  and the input by  $x^{(t)}$ . We denote by  $\sigma$  the function

$$\begin{aligned}\sigma : \mathbb{R}^k &\longrightarrow \mathbb{R}^k \\ \sigma(X)_i &= \frac{1}{2}(X_i + |X_i|).\end{aligned}$$

An RNN is defined by its *evolution equation* which expresses  $h^{(t+1)}$  as a function of  $h^{(t)}$ ,  $x^{(t)}$  and its *output equation* which gives the output  $y^{(t)}$  as a function of  $h^{(t)}$ .

**Definition 2.1.** A standard Elman-style RNN [1] is defined by weight matrices  $H, U, B$  and the evolution equation

$$(2.1) \quad h^{(t+1)} = \sigma(Hh^{(t)} + Ux^{(t+1)} + B)$$

where  $h^{(0)}$  is some specified initial state. The outputs are defined by

$$(2.2) \quad y^{(t)} = \sigma(Lh^{(t)} + M).$$

Suppose the input vector space is  $\mathcal{I}$  and the hidden state space is  $\mathcal{H}$ , so that  $x^{(t)} \in \mathcal{I}$  and  $h^{(t)} \in \mathcal{H}$  for all  $t$ . The value of the RNN on a sequence  $\mathbf{x} = (x^{(1)}, \dots, x^{(T)}) \in \mathcal{I}^{\oplus T}$  is computed by applying (2.1) for  $0 \leq t \leq T-1$ , accumulating the values  $\mathbf{y} = (y^{(1)}, \dots, y^{(T)})$  from each time step, and finally applying a fully-connected layer and the softmax function  $\tau$  obtain the output sequence  $o^{(t)} = \tau(Ey^{(t)} + F)$ . Each  $o^{(t)}$  has the property that its components in our chosen basis add to 1. We view such a vector as a probability distribution over the basis, sampling from which gives the output of the RNN on  $\mathbf{x}$ .

We consider a generalisation of the second-order RNN [11, 12, 9, 10] and the similar multiplicative RNN [7]. In addition to the hidden-to-hidden matrix  $H$  and the input-to-hidden matrix  $U$  of the traditional RNN, the second-order RNN learns a matrix  $V$  that

maps inputs to linear operators on the hidden state. More precisely, a vector is added to the evolution equation whose  $i$ th coordinate is given by the formula

$$(2.3) \quad \sum_j \sum_k V_i^{jk} h_j^{(t)} x_k^{(t+1)}$$

where  $V$  is interpreted as a tensor in

$$(2.4) \quad \mathcal{I}^* \otimes \mathcal{H}^* \otimes \mathcal{H} \cong \text{Hom}_{\mathbb{R}}(\mathcal{I}, \text{End}_{\mathbb{R}}(\mathcal{H})).$$

Identifying the tensor  $V$  with a linear map from the input space  $\mathcal{I}$  to linear operators on  $\mathcal{H}$ , we have that  $V(x^{(t+1)})(h^{(t)})$  is the vector whose  $i$ th coordinate is (2.3).

**Definition 2.2.** The second-order RNN [11, 12] is defined by weights  $H, U, B, V$  and

$$(2.5) \quad h^{(t+1)} = \sigma(V(x^{(t+1)})(h^{(t)}) + Hh^{(t)} + Ux^{(t+1)} + B),$$

with  $y^{(t)}, o^{(t)}$  as before.

The problem with second-order RNNs is that they may be difficult to train if the state space is large, since  $\dim(\text{End}_{\mathbb{R}}(\mathcal{H})) = \dim(\mathcal{H})^2$ . The *multiplicative RNN* is introduced in [7] as a more tractable model. The central idea is the same, but an auxiliary space  $\mathcal{K}$  is introduced and the RNN learns three linear maps

$$(2.6) \quad V : \mathcal{I} \longrightarrow \text{End}_{\mathbb{R}}(\mathcal{K}), \quad I : \mathcal{K} \longrightarrow \mathcal{H}, \quad J : \mathcal{H} \longrightarrow \mathcal{K}.$$

Moreover  $V$  factors through the subspace of diagonalisable matrices in some chosen basis, so it is defined by  $\dim(\mathcal{K})$  free parameters. The additional term in the evolution equation (2.5) is changed to

$$(2.7) \quad I(V(x^{(t+1)})(Jh^{(t)})).$$

The multiplicative RNN has been applied to character-level text generation [7] and sentiment analysis [6]. This is not to be confused with the *multiplicative integration* RNN of [5] which adds a term  $Hh^{(t)} \odot Ux^{(t+1)}$  to the evolution equation.

**Remark 2.3.** The second-order RNN transforms input symbols into linear operators on its hidden state. Observe that consecutive symbols  $x = x^{(t)}$  and  $x' = x^{(t+1)}$  in the input become composed operators on the hidden state, since (ignoring the non-linearity)

$$\begin{aligned} V(x')(h^{(t)}) &= V(x')(V(x)(h^{(t-1)}) + \dots) \\ &= \{V(x') \circ V(x)\}(h^{(t-1)}) + \dots. \end{aligned}$$

The relevance of this “compositionality” to NLP is remarked on in [6].

## 2.2 The linear logic RNN

At each time step the second-order RNN generates from its input a linear operator which is applied to the hidden state. The Linear Logic Recurrent Neural Network (LLRNN) generalises this by generating at each time step a sequence of *commands* (which are linear combinations of denotations of linear logic programs) which are applied to a sequence of *input operators* by a fixed *master algorithm*. The commands depend only on the previous hidden state, while the input operators depend only on the inputs  $x^{(t+1)}, x^{(t)}, \dots, x^{(t-s+2)}$  for some *window length*  $s$ .

**Definition 2.4.** The Linear Logic Recurrent Neural Network (LLRNN) is determined by the following data:

- types  $P_1, \dots, P_r$  called the *command types*;
- types  $B_1, \dots, B_s$ , called the *input types*;
- a type  $A$  with  $\llbracket A \rrbracket = \mathcal{H}$ , the hidden state space;
- a proof master of the sequent

$$(2.8) \quad !P_1, \dots, !P_r, !(B_1 \multimap B_1), \dots, !(B_s \multimap B_s) \vdash A \multimap A.$$

- Finite-dimensional subspaces

$$(2.9) \quad \mathcal{P}_i \subseteq \llbracket P_i \rrbracket, \quad \mathcal{B}_j \subseteq \llbracket B_j \multimap B_j \rrbracket = \text{End}_{\mathbb{R}}(\llbracket B_j \rrbracket)$$

for  $1 \leq i \leq r$  and  $1 \leq j \leq s$ . We refer to the elements of  $\mathcal{P}_i$  as *command vectors* and elements of  $\mathcal{B}$  as *input operators*.

From this data we obtain the smooth function

$$(2.10) \quad \llbracket \text{master} \rrbracket_{nl} : \mathcal{P}_1 \times \dots \times \mathcal{P}_r \times \mathcal{B}_1 \times \dots \times \mathcal{B}_s \longrightarrow \text{End}_{\mathbb{R}}(\mathcal{H})$$

by restriction. At each time step the RNN will generate from the hidden state a sequence of command vectors, and from the current and previous inputs a sequence of input operators. This is done via weight matrices  $\{W_i\}_{1 \leq i \leq r}$  and  $\{V_j\}_{1 \leq j \leq s}$  which determine linear maps

$$W_i : \mathcal{H} \longrightarrow \mathcal{P}_i, \quad V_j : \mathcal{I} \longrightarrow \mathcal{B}_j.$$

Finally, the evolution equation is given by

$$(2.11) \quad h^{(t+1)} = \sigma\left(Z + Hh^{(t)} + Ux^{(t+1)} + B\right)$$

where

$$Z = \llbracket \text{master} \rrbracket_{nl}\left(W_1(h^{(t)}), \dots, W_r(h^{(t)}), V_1(x^{(t+1)}), V_2(x^{(t)}), \dots, V_s(x^{(t-s+2)}))\right)(h^{(t)}).$$

**Example 2.5.** With  $r = 0$  and  $s = 1$  set  $B_1 = A$  take master to be the proof of

$$!(A \multimap A) \vdash A \multimap A$$

which is just the dereliction map, so that with  $\mathcal{B}_1 = \text{End}_{\mathbb{R}}(\mathcal{H})$  we have  $\llbracket \text{master} \rrbracket_{nl}$  equal to the identity map on  $\text{End}_{\mathbb{R}}(\mathcal{H})$ , and consequently

$$Z = V_1(x^{(t+1)})(h^{(t)}).$$

That is, the LLRNN with these settings is just the second-order RNN.

**Example 2.6.** We can emulate the stack-augmented RNN [?] with appropriate choices of the defining data of the LLRNN. Decompose the hidden state space as  $\mathcal{H} = \mathcal{H}' \oplus \mathcal{S}$  where the “memory stack”

$$\mathcal{S} = k[z]/z^{n+1} = k \oplus kz \oplus \cdots \oplus kz^n$$

is equipped with the operators

$$\begin{aligned} p_+ : \mathcal{S} &\longrightarrow \mathcal{S}, & p_+(v) &= zv, \\ p_- : \mathcal{S} &\longrightarrow \mathcal{S}, & p_-(v) &= \frac{\partial}{\partial z}(v), \end{aligned}$$

referred to as *push* and *pop* respectively. The master program is the proof of

$$(2.12) \quad !\mathbf{int}_A, !\mathbf{int}_A, !(S \multimap S), !(S \multimap S) \vdash A \multimap A$$

with the property that  $\llbracket S \rrbracket = \mathcal{S}$  and

$$\llbracket \text{master} \rrbracket_{nl}(\llbracket m \rrbracket, \llbracket n \rrbracket, \alpha, \beta) = J \circ \alpha^m \circ \beta^n \circ I,$$

where  $I : \mathcal{H} \longrightarrow \mathcal{S}$  and  $J : \mathcal{S} \longrightarrow \mathcal{H}$  are the projections and injections respectively into the direct sum. We have  $r = s = 2$  and  $P_i = \mathbf{int}_A, B_j = A$ . We either restrict  $\llbracket \text{master} \rrbracket_{nl}$  to a subset of its input space, or fix the weight matrices  $V_j$  so that the image of every basis element under  $V_1$  (resp.  $V_2$ ) is  $p_+$  (resp.  $p_-$ ). Finally, take  $\mathcal{B}_j = \text{End}_{\mathbb{R}}(\mathcal{H})$  and for some fixed integer  $L$  take

$$\mathcal{P}_1 = \mathcal{P}_2 = \text{span}(\llbracket 0 \rrbracket, \dots, \llbracket L \rrbracket) \subseteq \llbracket \mathbf{int}_A \rrbracket.$$

Roughly speaking, the LLRNN generates at each time step two integers  $m, n$  and writing  $h^{(t)} = (h_0^{(t)}, s^{(t)})^T$ , we have

$$Z = \begin{pmatrix} 0 & p_+^m p_-^n(s^{(t)}) \end{pmatrix}^T.$$

Note that the linear logic part of the update equation is only responsible for “rotating” the stack, adding an element to the top of the stack can be accomplished using the  $Hh^{(t)}$  term.

**Remark 2.7.** The complexity of an algorithm is rooted in the kind of *iteration* that it employs; for an exposition in the context of linear logic see [13, §]. The most elementary forms of iteration are simply repetitions, as in the proofs of type  $\mathbf{int}_\alpha$ ,  $\mathbf{bint}_\alpha$  above. Our proposal is that in order to enable an RNN controller to learn functions of a certain complexity, one should provide it with access to iterators of a “similar level” of complexity.

**Example 2.8.** If  $\Lambda = \{\underline{1}, \underline{2}\}$  and  $\Pi = \{\underline{010}\}$  then

$$\begin{aligned} Z = & I(p_1^{(t)}V(x^{(t)})(Jh^{(t)})) + I(p_2^{(t)}V(x^{(t)})^2(Jh^{(t)})) \\ & + I(q_{010}^{(t)}\{V(x^{(t)}) \circ V(x^{(t+1)}) \circ V(x^{(t)})\}(Jh^{(t)})). \end{aligned}$$

The first summand here is the same as (2.7), which shows that the linear logic RNN with  $\Lambda = \{\underline{1}\}$ ,  $\Pi = \emptyset$  and the weight  $p_1^{(t)}$  fixed to the identity, is the multiplicative RNN of [7].

**Remark 2.9.** To avoid too many weights and thus slow training, we follow the multiplicative RNN in assuming that  $V$  maps inputs to diagonal matrices. But the whole point of adding the ability to call programs from  $\Pi$  is that binary integers intertwine  $V(x^{(t)})$  and  $V(x^{(t+1)})$  in some useful way, which is compromised if these two matrices commute with one another.

So in practice we have two auxiliary spaces  $\mathcal{K}_1, \mathcal{K}_2$  and two sets of weights  $V_i, I_i, J_i$ . The index  $i = 1$  corresponds to  $\Lambda$ , and  $V_1 : \mathcal{I} \rightarrow \text{End}_{\mathbb{R}}(\mathcal{K}_1)$  factors through the subspace of matrices diagonal in a chosen basis. The map  $V_2$  is not restricted in this way, but to compensate we take  $\dim(\mathcal{K}_2) < \dim(\mathcal{K}_1)$ .

**Remark 2.10.** A more powerful system would substitute a neural theorem prover along the lines of [22, 24] in place of the library of functional programs  $\Pi, \Lambda$ . At each time step the RNN controller would predict a continuous vector, which when fed into the neural theorem prover as a set of parameters, generates a symbolic program whose denotation is then coupled back into the RNN [23].

**Remark 2.11.** Although intuitionistic linear logic is a weak system in terms of expressive strength the difficulty of training the RNN with polynomial functions of high degree makes it infeasible to make use of even this much strength. It would, however, be very interesting to define differentiable semantics for second-order logic and try to use this to improve the generalisation ability of neural networks.

**Remark 2.12.** Restricting ourselves to the  $\mathbf{int}_\alpha$  augmentation for the moment, let us discuss what kind of functions  $f : \{0, 1\}^* \rightarrow \{0, 1\}^*$  we would expect that it would be easier for the LLRNN to learn (all other things being equal) compared to an ordinary RNN. The LLRNN has the ability to predict, using its internal state, how many times to iterate the input-derived operator that is to be applied to the hidden state. Thus a function  $f$  which, for instance, sends 0s to 1 unless they are preceded by a 0 in which case they go to 0, should be helped.

## .1 Cuts

Let  $\Lambda = (\lambda_1, \dots, \lambda_l)$  be a sequence of proofs of  $\mathbf{int}_\alpha$  and let  $\Pi = (\pi_1, \dots, \pi_k)$  be proofs of  $\mathbf{bint}_\alpha$ . When we say that we allow the RNN controller to call this library of functional programs, what we actually mean is that it can predict at each time step diagonal matrices  $\{p_\lambda^{(t)}\}_{\lambda \in \Lambda}$  and  $\{q_\pi^{(t)}\}_{\pi \in \Pi}$  that determine which programs are used.

**Definition .13.** The *linear logic* RNN coupled to  $\Lambda, \Pi$  has weight matrices  $H, U, B$  as in the standard RNN, plus new weights  $V, P, Q, C, D, I, J$  and an evolution equation:

$$\begin{aligned} p_\lambda^{(t+1)} &= \sigma(P_\lambda h^{(t)} + C_\lambda), \\ q_\pi^{(t+1)} &= \sigma(Q_\pi h^{(t)} + D_\pi), \\ h^{(t+1)} &= \sigma\left(Z + Hh^{(t)} + Ux^{(t+1)} + B\right) \end{aligned}$$

where

$$Z = \sum_{\lambda \in \Lambda} I p_\lambda^{(t+1)} \llbracket \lambda \rrbracket (Vx^{(t+1)})(Jh^{(t)}) + \sum_{\pi \in \Pi} I q_\pi^{(t+1)} \llbracket \pi \rrbracket (Vx^{(t)}, Vx^{(t+1)})(Jh^{(t)}).$$

The weight matrices  $V, I, J$  are interpreted as linear maps as in (2.6). The operators  $p_\lambda^{(t)}$  and  $q_\pi^{(t)}$  are diagonal matrices in  $\text{End}_{\mathbb{R}}(\mathcal{K})$ . The outputs  $y^{(t)}$  of the linear logic RNN are defined by the same equation (2.2) as before.

Note that the RNN evolution equation for  $h^{(t+1)}$  involves both the current input  $x^{(t+1)}$  and the previous input  $x^{(t)}$ . We do this in order to allow the RNN to learn local transformations of an input sequence.

## .2 Differential linear logic

In functional programming languages like differential  $\lambda$ -calculus [30] and differential linear logic [25] it is possible to differentiate programs with respect to their inputs, even if the programs are higher-order (that is, take functions as input and return them as output).

This is not as exotic as it might sound: it is a generalisation of automatic differentiation [29] which is widely used to compute derivatives of complicated real-valued functions, for example in the backpropagation algorithms of machine learning packages like TensorFlow [27, §4.1]. The idea is to augment every computation in a code fragment so that derivatives with respect to a chosen variable are computed along with the main result. In TensorFlow this is done by adding special nodes in the dataflow graph.<sup>1</sup> The idea of differential  $\lambda$ -calculus is similar, but more complex [31]. It would be interesting to explore augmenting the dataflow graph of TensorFlow directly with terms of differential linear logic, in a way that generalises the coupling between semantics and RNNs in this paper.

---

<sup>1</sup>See the “gradients” method of `tensorflow/python/ops/gradients.py` in TensorFlow 0.10

### .3 Linear logic

In order to improve the expressive strength of an RNN, that is, allow it to approximate more complex functions from sequences of input symbols to sequences of output symbols, a natural idea is to provide a library of functional programs that can be “called” at each time step. The main problem lies in arranging for these programs to be differentiable, so that gradients can be backpropagated through them at training time.

One natural approach is to have an abstract, symbolic *source* language in which the relevant functional programs can be written, together with a concrete, differentiable *target* category in which the programs can be interpreted. This interpretation should respect the structure of the functional programs: for example, composition or plugging of programs should be interpreted as composition of differentiable maps. The purpose of the source language is to ensure expressive strength, that is, to provide a library of differentiable maps sufficiently rich to encode nontrivial symbolic transformations. The purpose of the target category is to facilitate a straightforward coupling into the RNN.

A candidate for the source language is linear logic [16] which has denotational semantics in differentiable maps [25]. We use the particular semantics developed in [13, 14, 26]. The second-order RNN encourages us to view symbols as operators on a vector space, and this aligns very well with the intuition provided by linear logic.

The *types* of intuitionistic linear logic (ILL) are built from atomic variables  $\alpha, \beta, \dots$  by applying binary connectives  $\otimes, \multimap$  (tensor, Hom) and a unary connective  $!$  (called bang). For example, for any type  $A$  we have the types

$$\begin{aligned}\mathbf{int}_A &= !(A \multimap A) \multimap (A \multimap A), \\ \mathbf{bint}_A &= !(A \multimap A) \multimap ( !(A \multimap A) \multimap (A \multimap A) ),\end{aligned}$$

respectively of integers on  $A$ , and binary integers on  $A$ . The semantics will map  $A$  to a vector space, and  $A \multimap A$  to the space of linear operators on that vector space.

A program  $\pi$  of type  $M$ , written  $\pi : M$ , is built from atomic pieces by rules that we will not explain here. The simplest rule is composition or *plugging*: given programs  $M : \alpha \multimap \beta$  and  $N : \beta \multimap \gamma$  we have a combined program  $NM : \alpha \multimap \gamma$ . A program of type  $!A \multimap B$  takes inputs of type  $A$  and returns outputs of type  $B$ . The same is true of a program of type  $A \multimap B$  but in addition it satisfies a *linearity constraint* in that it uses its input *precisely once*. So the  $!$  connective marks inputs that may be reused freely (or not used at all). All inputs and outputs are themselves programs of the appropriate type.

**Example .14.** A program  $\pi$  of type  $\mathbf{int}_A$  takes an input  $M$  of type  $A \multimap A$  that may be used non-linearly, and outputs a program of type  $A \multimap A$ . For each integer  $n \geq 0$  there is a program  $\underline{n} : \mathbf{int}_A$  called the  *$n$ th Church numeral*, with the property that it maps an input  $M$  to the  $n$ -fold repetition  $M^n$ .

The Church numeral is an *iterator* (something like a *for* loop in imperative programming) which takes a program  $M$  and iterates it  $n$  times.



**Example .15.** A program  $\pi$  of type  $\mathbf{bint}_\alpha$  maps an input of type  $\alpha \multimap \alpha$  to a program of type  $\mathbf{int}_\alpha = !(\alpha \multimap \alpha) \multimap (\alpha \multimap \alpha)$ . So, given two programs  $M, N$  of type  $\alpha \multimap \alpha$ , we can apply the output  $\pi(M)$  to  $N$ , to get a program of type  $\pi(M)(N) : \alpha \multimap \alpha$ .

Given a binary sequence  $S \in \{0, 1\}^*$  there is a program  $\underline{S} : \mathbf{bint}_\alpha$  with the property that it maps an input pair  $(M, N)$  to the program which intertwines copies of  $M, N$  according to  $S$ . For example,  $\underline{01101}$  sends  $(M, N)$  to the program  $MNNMN$ .

A *denotational semantics* of this language selects for each type  $A$  some mathematical object  $\llbracket A \rrbracket$  and for a program of type  $A \multimap B$  a transformation  $\llbracket A \rrbracket \longrightarrow \llbracket B \rrbracket$  [21]. The precise kind of mathematical object, and transformations between them, is determined by the kind of category we take as the target of the semantics. In this paper we use the Sweedler semantics in vector spaces of [13, 14] and the polynomial semantics which refines it. We will not give the details, since we only need a small number of examples.

Starting with an atomic type  $\alpha$  whose denotation  $\llbracket A \rrbracket = \mathcal{K}$  is some finite-dimensional vector space over  $\mathbb{R}$ , we have

$$\llbracket A \multimap A \rrbracket = \mathrm{Hom}_{\mathbb{R}}(\llbracket A \rrbracket, \llbracket A \rrbracket) = \mathrm{End}_{\mathbb{R}}(\mathcal{K}).$$

**Example .16.** The denotation of a proof  $\pi$  of type  $\mathbf{int}_\alpha$  is a function

$$\llbracket \pi \rrbracket : \mathrm{End}_{\mathbb{R}}(\mathcal{K}) \longrightarrow \mathrm{End}_{\mathbb{R}}(\mathcal{K})$$

which is a polynomial function of the entries of the input matrix. For example,

$$\llbracket \underline{n} \rrbracket(X) = X^n.$$

To be clear, writing  $\mathcal{K} = \mathbb{R}^{n_H}$  and identifying  $\mathrm{End}_{\mathbb{R}}(\mathcal{K})$  with the space of matrices  $M_{n_H}(\mathbb{R})$ , the denotation of any proof of  $\mathbf{int}_\alpha$  will be a function  $\llbracket \pi \rrbracket(X)_{ij} : M_{n_H}(\mathbb{R}) \longrightarrow \mathbb{R}$  which is a polynomial in the entries  $X_{kl}$  for  $1 \leq i, j \leq n_H$ .

**Example .17.** The denotation of a proof  $\pi$  of type  $\mathbf{bint}_\alpha$  is a polynomial function

$$\llbracket \pi \rrbracket : \mathrm{End}_{\mathbb{R}}(\mathcal{K}) \times \mathrm{End}_{\mathbb{R}}(\mathcal{K}) \longrightarrow \mathrm{End}_{\mathbb{R}}(\mathcal{K}).$$

For every binary sequence  $S \in \{0, 1\}^*$  the denotation  $\llbracket \underline{S} \rrbracket$  sends a pair of matrices  $X, Y$  to the product described by  $S$ , reading  $X$  for 0 and  $Y$  for 1. For example,

$$\llbracket \underline{01101} \rrbracket(X, Y) = XY YXY.$$

The idea of the second-order RNN is to view the input and output symbols as linear operators on a vector space  $\mathcal{K}$  (which, following the idea of the multiplicative RNN, we take as possibly different from the hidden state space  $\mathcal{H}$ ). Linear logic gives us a source language in which to create programs, whose denotations are differentiable maps taking linear operators as input and providing linear operators as output. This map from input to output may not itself be linear, as the above examples show. It is these differentiable shadows of symbolic transformations in linear logic that we will couple to the RNN.

Let us conclude with an example involving more complicated types:

**Example .18.** There is a program

$$\underline{subst} : \mathbf{bint}_\alpha \multimap \left( !\mathbf{bint}_\alpha \multimap \left( !\mathbf{bint}_\alpha \multimap \mathbf{bint}_\alpha \right) \right)$$

which given three input sequences  $\underline{S}, \underline{T}, \underline{U}$  uses the pattern of the sequence  $S$  to intertwine copies of  $T, U$ . A program, say  $\underline{T}$ , of type  $\mathbf{bint}_\alpha$  may be “promoted” to a program of type  $!\mathbf{bint}_\alpha$ , denoted  $[\underline{T}]$ . Then, for example,

$$\underline{subst}(\underline{010}, [\underline{T}], [\underline{U}]) = \underline{T} \cdot \underline{U} \cdot \underline{T} : \mathbf{bint}_\alpha$$

where  $\cdot$  means concatenation.

## A Background on CPUs

Recall that an assembly program for an ordinary CPU looks like

```
LOAD R1, A
ADD R3, R1, R2
STORE C, R3
```

Where  $R1, R2, R3$  stand for the first three registers of the CPU and  $A, B, C$  are numbers representing addresses in memory. Thus series of instructions will result in the CPU fetching a number from memory location  $A$  and storing it in  $R1$ , adding this number to a previously stored number in  $R2$  with the result being stored in  $R3$ , and finally writing that register out to the memory address  $C$ . In the analogy between a CPU and a vanilla RNN we think of the number read from  $A$  as the current input  $x^{(t)}$  and the previously stored value in  $R2$  as (part of) the internal state  $h^{(t-1)}$ .

Recent work [?, ?] on coupling memory to neural networks takes as its starting point the first of the above instructions **LOAD R1, A** and makes it “differentiable” by having the RNN controller predict at time  $t$  both the memory address  $A$  and the register  $R3$  to write to (in this case for example, as a mask on the vector giving the internal state  $h^{(t+1)}$ ). The same differentiable interpretation can be given of the **STORE** command. This is done by adding suitable terms to the update equation (2.1).

In contrast our focus is on the third command, the **ADD**. We increase the expressive power of the update equation by allowing it to predict at time  $t$  an operation  $p^{(t)}$  (a vector specifying a point in a space of “programs”) which is to be performed on the input and internal state. In order to preserve our analogy with CPU instructions even without **LOAD** and **STORE**, we could imagine a CPU with a command

```
ADD R3, A, R2
```

which at once reads the number from address  $A$ , adds it to the stored value in  $R2$  and writes the result to  $R3$ . Note that without a **LOAD** instruction, the only way a value could have gotten into  $R2$  in a previous time step is as the result of another **ADD**.

**Remark A.1.** Although the architecture takes some inspiration from normal CPUs, there is an important distinction: on a normal CPU the program is given as a series of instructions prior to the beginning of execution. In contrast, in the model we have described each command is *predicted* at runtime from the current internal state. Perhaps we can understand the process intuitively as follows: we are co-learning a part of  $H$ , call it  $H_0$ , which generates some part of the internal state  $h_0^{(1)}, h_0^{(2)}, \dots$  giving a path through the state space on which the weight matrix  $P$  picks out the right program to run at each time step. The overall algorithm is distributed amongst the weights of  $H_0$  and  $P$ .

This also suggests an alternative algorithm: we do not predict  $p^{(t)}$  at each time step, rather we have some fixed number of time steps  $T$  and matrices of weights  $p^{(1)}, \dots, p^{(T)}$  which are learned by gradient descent.

## References

- [1] J. Elman, *Finding structure in time*, Cognitive science, 14(2):179-211, 1990.
- [2] A. Graves, G. Wayne and I. Danihelka, *Neural Turing machines*, arXiv preprint arXiv:1410.5401 (2014).
- [3] A. Joulin and T. Mikolov, *Inferring algorithmic patterns with stack-augmented recurrent nets*, Advances in Neural Information Processing Systems, 2015.
- [4] A. Graves, *Hybrid computing using a neural network with dynamic external memory*, Nature 538.7626 (2016): 471–476.
- [5] Y. Wu, S. Zhang, Y. Zhang, Y. Bengio and R. R. Salakhutdinov, *On multiplicative integration with recurrent neural networks*, In Advances In Neural Information Processing Systems, pp. 2856-2864. 2016.
- [6] O. Irsoy and C. Cardie, *Modeling compositionality with multiplicative recurrent neural networks*, arXiv preprint arXiv:1412.6577 (2014).
- [7] I. Sutskever, J. Martens and G. E. Hinton, *Generating text with recurrent neural networks* Proceedings of the 28th International Conference on Machine Learning (ICML-11). 2011.
- [8] I. Sutskever, O. Vinyals and Q. V. Le, *Sequence to sequence learning with neural networks*, Advances in neural information processing systems, 2014.
- [9] M. W. Goudreau, C. L. Giles, S. T. Chakradhar and D. Chen, *First-order versus second-order single-layer recurrent neural networks*, IEEE Transactions on Neural Networks, 5(3), 511–513, 1994.

- [10] C. L. Giles, D. Chen, C. B. Miller, H. H. Chen, G. Z. Sun, Y. C. Lee, *Second-order recurrent neural networks for grammatical inference*, In Neural Networks, 1991., IJCNN-91-Seattle International Joint Conference on (Vol. 2, pp. 273-281). IEEE.
- [11] C. L. Giles, G. Z. Sun, H. H. Chen, Y. C. Lee, D. Chen, *Higher order recurrent networks and grammatical inference*, In NIPS (pp. 380-387) 1989.
- [12] J. B. Pollack, *The induction of dynamical recognizers*, Machine Learning, 7(2-3), 227-252 (1991).
- [13] D. Murfet, *Logic and linear algebra: an introduction*, preprint (2014) [arXiv:1407.2650].
- [14] D. Murfet, *On Sweedler's cofree cocommutative coalgebra*, J. Pure and Applied Algebra **219** (2015) 5289–5304.
- [15] A. Graves, G. Wayne and I. Danihelka, *Neural turing machines*, preprint (2014) [arXiv:1410.5401].
- [16] J.-Y. Girard, *Linear Logic*, Theoretical Computer Science **50** (1987), 1–102.
- [17] A. Joulin and T. Mikolov, *Inferring algorithmic patterns with stack-augmented recurrent nets*, Advances in Neural Information Processing Systems, 2015.
- [18] E. Grefenstette, et al, *Learning to transduce with unbounded memory*, Advances in Neural Information Processing Systems, 2015.
- [19] J. Weston, C. Sumit and B. Antoine, *Memory networks*, preprint (2014) [arXiv:1410.3916].
- [20] W. Zaremba, et al., *Learning Simple Algorithms from Examples*, preprint (2015) [arXiv:1511.07275].
- [21] P-A. Mellies, *Categorical semantics of linear logic*, in : Interactive models of computation and program behaviour, Panoramas et Synthèses 27, Société Mathématique de France, 2009.
- [22] A. A. Alemi, F. Chollet, G. Irving, C. Szegedy and J. Urban, *DeepMath-Deep Sequence Models for Premise Selection*, arXiv preprint arXiv:1606.04442.
- [23] M. Balog, A. L. Gaunt, M. Brockschmidt, S. Nowozin and D. Tarlow, *DeepCoder: Learning to Write Programs*, arXiv preprint arXiv:1611.01989.
- [24] T. Rocktäschel and S. Riedel, *Learning Knowledge Base Inference with Neural Theorem Provers*, In NAACL Workshop on Automated Knowledge Base Construction (AKBC) 2016.

- [25] T. Ehrhard, *An introduction to Differential Linear Logic: proof-nets, models and antiderivatives*, [\[arXiv:1606.01642\]](#) (2016).
- [26] J. Clift and D. Murfet, *The Sweedler semantics is a natural model of differential linear logic*, preprint.
- [27] Abadi, Martn, et al. *TensorFlow: A system for large-scale machine learning* arXiv preprint arXiv:1605.08695 (2016).
- [28] Abadi, Martn, et al. *Tensorflow: Large-scale machine learning on heterogeneous distributed systems*, arXiv preprint arXiv:1603.04467 (2016).
- [29] A. Griewank and A. Walther, *Evaluating derivatives: principles and techniques of algorithmic differentiation*, Siam (2008).
- [30] T. Ehrhard and L. Regnier, *The differential  $\lambda$ -calculus*, Theoretical Computer Science 309, pp. 1–41, (2003).
- [31] O. Manzyuk, *A simply typed  $\lambda$ -calculus of forward automatic differentiation*, In Mathematical Foundations of Programming Semantics Twenty-eighth Annual Conference, pages 259–73, Bath, UK, June 69 2012. [\[URL\]](#).