

# Linear logic and recurrent neural networks

Huiyi Hu, Daniel Murfet

February 20, 2017

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background</b>	<b>2</b>
2.1	Recurrent Neural Network . . . . .	3
2.2	Second-order Recurrent Neural Network . . . . .	3
2.3	Neural Turing Machine . . . . .	5
2.4	Linear Logic . . . . .	7
2.4.1	Types . . . . .	7
2.4.2	Terms . . . . .	8
2.4.3	Semantics . . . . .	9
<b>3</b>	<b>The Linear Logic Recurrent Neural Network</b>	<b>11</b>
3.1	Examples . . . . .	13
3.2	NTM as an LLRNN . . . . .	14
<b>4</b>	<b>Generalisations of the NTM</b>	<b>16</b>
4.1	Pattern NTM . . . . .	16
4.2	Multiple pattern NTM . . . . .	19
4.3	Polynomial step NTM . . . . .	21
4.3.1	Dihedral polynomial step NTM . . . . .	22
<b>5</b>	<b>Unpolished remarks</b>	<b>23</b>
5.1	Decaying pattern NTM . . . . .	25

## 1 Introduction

Deep neural networks have proven to be proficient at learning hierarchical vector representations of natural data, including images and language [9]. If we follow Leibniz [6] in

believing that reasoning is the algebraic manipulation of such representations or “symbols” then it is natural to look for systems which augment the capacity of neural networks for performing these kinds of manipulations. This is a foundational problem in the field of artificial intelligence which remains far from solved, despite renewed attention in recent years [11, 25, 5, 13, 26, 27, 31, 4, 3]. Arguably one obstacle is that some of the underlying mathematical ideas are still unclear; to help remedy this we investigate a new approach, based on the Curry-Howard correspondence and linear logic.

The Curry-Howard correspondence [8] is a bijection between a prototypical system of formal reasoning (Gentzen’s natural deduction) and a prototypical algorithmic system (Church’s simply-typed lambda calculus). For this reason lambda calculus and derivative languages such as LISP have played an important role in the symbolic approach to reasoning in the modern field of artificial intelligence. While these methods may have been overshadowed in recent decades by the rise of the connectionist approach, now called deep learning, it has been argued that a synthesis of the two approaches is necessary to achieve general reasoning in a connectionist system [40]. 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 [24] which has a canonical model (called a denotational semantics) in differentiable maps between vector spaces [33, 7, 34].

In this paper we realise these ideas in the form of a Recurrent Neural Network (RNN) controller [10] augmented with smooth denotations of linear logic programs. We call the resulting system the *Linear Logic Recurrent Neural Network* or LLRNN. This architecture is inspired by many papers in the neural network literature, most notably the second-order RNN [25], multiplicative RNN [16] and Neural Turing Machine [11]. Conceptually, all of these models augment the basic RNN by generating, at each time step, a *linear operator* from the previous hidden state and the current input. This linear operator is then applied to the previous hidden state to generate the new hidden state. The LLRNN generalises these models by using algorithms from linear logic to generate the linear operator.

## 2 Background

It is our hope that the paper is accessible both for readers coming from the deep learning community, and for readers coming from the linear logic community. To this end we will make some introductory remarks on both subjects.

## 2.1 Recurrent Neural Network

We begin with a review of the ordinary Recurrent Neural Network (RNN). A good text-book introduction to deep learning in general and RNNs is [43]. As usual we use “weight” as a synonym for variable, or more precisely, the variables which we will vary during gradient descent. We denote by  $\sigma$  the function

$$\sigma : \mathbb{R}^k \longrightarrow \mathbb{R}^k$$

$$\sigma(\mathbf{x})_i = \frac{1}{2}(x_i + |x_i|)$$

and by  $\zeta$  the *softmax* function

$$\zeta : \mathbb{R}^k \longrightarrow \mathbb{R}^k ,$$

$$\zeta(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^k e^{x_j}} .$$

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  $h^{(t)}$ . At time  $t \geq 1$  we denote the hidden state by  $h^{(t)}$  and the input by  $x^{(t)}$ .

**Definition 2.1.** A standard Elman-style RNN [10] 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

$$(2.2) \quad y^{(t)} = \sigma(W_y h^{(t)} + B_y) .$$

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  $o^{(t)} = \zeta(W_o y^{(t)} + B_o)$ . 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  $\mathbf{x}$ .

## 2.2 Second-order Recurrent Neural Network

We consider a generalisation of the second-order RNN [20, 21, 18, 19] and the similar multiplicative RNN [16]. 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 [20, 21] 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.

**Remark 2.3.** 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 [16] 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 [16] and sentiment analysis [15]. This is not to be confused with the *multiplicative integration* RNN of [14] which adds a term  $Hh^{(t)} \odot Ux^{(t+1)}$  to the evolution equation.

**Remark 2.4.** 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 [15].

## 2.3 Neural Turing Machine

The Neural Turing Machine (NTM) [11] is an RNN augmented with a second hidden state space that functions like an external memory. This memory is read from and written to using attention vectors which are manipulated based on the hidden state of the RNN. To describe the precise form of this manipulation, it is important to know that the memory state can be viewed as a sequence of vectors arranged at the vertices of a regular polygon. An *attention vector* is by definition a probability distribution over the vertices of this polygon. There is an natural action of rotation  $R$  on both memory states and attention vectors, and the read and write attention vectors are manipulated by predicting, at each time step, which power of the linear operator  $R$  to apply.

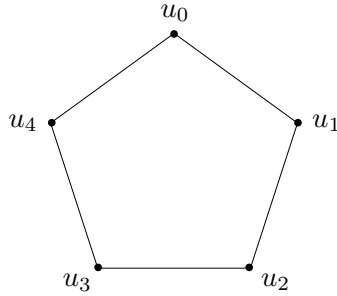
To explain, we first introduce the *address space*

$$\mathcal{W} = (\mathbb{Z}/N\mathbb{Z})^{\mathbb{R}} \cong \mathbb{R}\bar{0} \oplus \cdots \oplus \mathbb{R}\overline{N-1}$$

and given a *memory coefficient space*  $\mathcal{V}$  the *memory space* is the tensor product

$$\mathcal{S} = \mathcal{W}^* \otimes \mathcal{V} \cong \text{Hom}_{\mathbb{R}}(\mathcal{W}, \mathcal{V}).$$

A memory state  $M$  is a linear map  $M : \mathcal{W} \longrightarrow \mathcal{V}$  which we view as an arrangement of the vectors  $u_j = M(\bar{j})$  on the vertices of the regular  $N$ -gon:



Sometimes we refer to this structure as a *memory ring*. In the notation of [11] we may take  $\mathcal{W} = \mathbb{R}^N$  and  $\mathcal{V} = \mathbb{R}^M$  so that elements of  $\mathcal{S}$  are identified with  $N \times M$  matrices.

The space of memory states has a natural action of the cyclic group of order  $N$ , which acts by rotation on  $\mathcal{W}$  and therefore on  $\mathcal{S}$ . The rotated state assigns the vector  $u_{j+1}$  to the position  $j$  (we index the positions in the clockwise order) with  $j+1$  computed modulo  $N$ . Said differently, let  $R : \mathcal{W} \longrightarrow \mathcal{W}$  be the linear map defined by

$$R(\bar{a}) = \overline{a+1},$$

with dual map  $R^* : \mathcal{W}^* \longrightarrow \mathcal{W}^*$ . The rotation of the memory state  $M$  is the composition  $M \circ R$ , or what is the same thing,  $R^*(M)$ . Thus, applying the  $j$ -th power of  $R^*$  to a memory state has the effect of rotating the state  $j$  times around the regular polygon.

The internal state space of the RNN underlying the NTM is decomposed as

$$(2.8) \quad \mathcal{H} = \mathcal{H}_0 \oplus \mathcal{W} \oplus \mathcal{W}^* \oplus \mathcal{I}$$

and the state at time  $t$  is written

$$h^{(t)} = (h_0^{(t)}, r^{(t)}, w^{(t)}, M^{(t)}) \in \mathcal{H}.$$

The components of this vector are interpreted as the *controller internal state*  $h_0^{(t)}$ , *read address weights*  $r^{(t)}$ , *write address weights*  $w^{(t)}$  and *memory contents*  $M^{(t)}$ . For convenience we will refer to the read address weight as the *read address*, even though the vector represents a distribution over the available addresses rather than a specific address. At each time step the RNN generates from its hidden state a distribution  $s$  over rotations of the write address, a distribution  $q$  over rotations of the read address, an operator  $e$  to be applied to each coefficient vector in the memory state (this encodes “forgetting”) and a vector  $a$  to be written to memory, via the formulas

$$\begin{aligned} s^{(t+1)} &= \zeta(W_s h^{(t)} + B_s) \in \mathcal{W}^*, \\ q^{(t+1)} &= \zeta(W_q h^{(t)} + B_q) \in \mathcal{W}^*, \\ e^{(t+1)} &= \sigma(W_e h^{(t)} + B_e) \in \text{End}_{\mathbb{R}}(\mathcal{V}), \\ a^{(t+1)} &= \sigma(W_u h^{(t)} + B_u) \in \mathcal{V}. \end{aligned}$$

In practice we will assume that  $e^{(t+1)}$  is diagonal with respect to a chosen basis of  $\mathcal{V}$ . The update equation for the addresses [11, Eq. (8)] and memory [11, Eq. (3),(4)] are

$$(2.9) \quad w^{(t+1)} = \sum_{i=0}^{N-1} s^{(t+1)}(\bar{i}) \cdot (R^*)^i(w^{(t)}),$$

$$(2.10) \quad r^{(t+1)} = \sum_{i=0}^{N-1} q^{(t+1)}(\bar{i}) \cdot R^i(r^{(t)}),$$

$$(2.11) \quad M^{(t+1)} = (1_{\mathcal{V}} - e^{(t+1)}) \circ M^{(t)} + w^{(t)} \otimes a^{(t+1)}.$$

Finally, the evolution equation is

$$(2.12) \quad h_0^{(t+1)} = \sigma\left(M^{(t)}(r^{(t)}) + H_0 h_0^{(t)} + U_0 x^{(t+1)} + B_0\right).$$

where

$$\begin{aligned} H_0 &: \mathcal{H}_0 \longrightarrow \mathcal{H}_0, \\ U_0 &: \mathcal{I} \longrightarrow \mathcal{H}_0, \end{aligned}$$

and  $B_0 \in \mathcal{H}_0$  are weight matrices. We assume for simplicity that the memory coefficient space is the same as the hidden state space, that is,  $\mathcal{V} = \mathcal{H}$ . Following [11, §3.2] we

write  $e$  for the *erase vector* and  $a$  for the *add vector*. These equations define the NTM with location-based addressing, with an RNN controller and memory of size  $N$  with a single read and write head.<sup>1</sup>

**Remark 2.5.** To connect the above more explicitly to the notation of [11, Eq. (8)] observe that (2.10) can be rewritten, writing  $r[j]$  for the component of  $r \in \mathcal{W}$  for  $\bar{j}$ , as

$$(2.13) \quad r^{(t+1)}[i] = \sum_j q^{(t+1)}(\bar{j}) r^{(t)}[i - j].$$

## 2.4 Linear Logic

Linear logic is the universal way of describing algorithms which manipulate tensors and linear maps using the basic constructions of linear algebra (such as composition and tensor contraction) and iterations of such constructions. At the level of *syntax* these algorithms may be thought of as programs written in a functional programming language much like the simply-typed lambda calculus. This programming language has a functorial *semantics* denoted here by  $\llbracket - \rrbracket$  which realises the abstract algorithms as concrete functions.

In this section we give a very brief introduction to linear logic; for a longer survey see [22, 34]. We divide the introduction into three parts, covering the *types*, *terms* (meaning programs) and *semantics* of intuitionistic linear logic.

### 2.4.1 Types

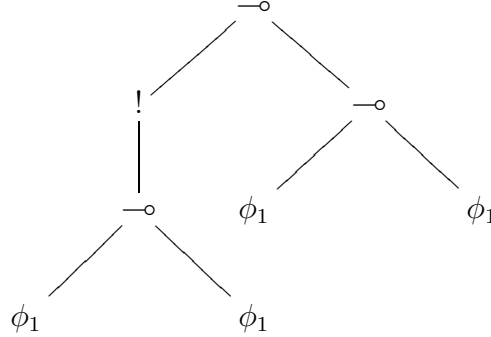
The *types* of intuitionistic linear logic are built from atomic types  $\phi_1, \phi_2, \dots$  via binary connectives  $\otimes, \multimap$  (tensor, Hom) and a unary connective  $!$  (called bang). Thus a type may be represented as a syntax tree, where the leaves are labelled by atomic types  $\phi_i$  and the other vertices are labelled by  $\otimes, \multimap, !$ . For example the type

$$!(\phi_1 \multimap \phi_1) \multimap (\phi_1 \multimap \phi_1)$$

---

<sup>1</sup>We omit the sharpening step [11, Eq. (9)] for simplicity. The original NTM paper [11] is not specific about how the output  $M^{(t)}(r^{(t)})$  of the read heads enters the RNN evolution equation; the form we have given above follows the construction of the differentiable neural computer [5, p. 7]. Note that the term appearing in the equation for  $h^{(t+1)}$  is the memory state at time  $t$  applied to the output of the read heads at time  $t$ . That is, the output of the read heads is incorporated into the hidden state of the RNN at the *next* time step; again, in this architectural choice we are following [5]. Similarly, the write address used to define  $M^{(t+1)}$  is  $w^{(t)}$  rather than  $w^{(t+1)}$ . In this it seems we actually *disagree* with [11, 5]. It makes the formulas cleaner, but maybe it is wrong (TODO).

corresponds to the tree



Given a type  $A$  we can construct two new types

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

called respectively the types of *integers* on  $A$  and *binary integers* on  $A$ , which will play an important role in what follows.

### 2.4.2 Terms

There is a countable collection of variables  $x_1^A, x_2^A, \dots$  of each type  $A$ . Variables are *terms* and a general term is built up from variables by a set *deduction rules* which we will not review here, but which we will try to illuminate with examples. As usual we write  $\pi : A$  to indicate that  $\pi$  is a term of type  $A$  and present terms as linear lambda terms, following [42]. Some examples of linear lambda terms are, for variables  $x : A, y : B$  and  $z : !A$ ,

$$(2.14) \quad (\lambda x. x) : A \multimap A,$$

$$(2.15) \quad (\lambda x. (\lambda y. (y x))) : A \multimap ((A \multimap A) \multimap A),$$

$$(2.16) \quad (\lambda z. \text{derelict}(z)) : !A \multimap A.$$

**Remark 2.6.** A term of type  $A \multimap B$  should be thought of as a program which takes inputs of type  $A$  and returns outputs of type  $B$ , under the constraint that the input is used *linearly* (that is, precisely once). For example, the term (2.14) just returns its input, while (2.15) takes two inputs  $x : A, y : A \multimap A$  and returns the result  $(y x)$  of applying the second program to the first. In both cases the inputs are used linearly.

On the other hand a term of type  $!A \multimap B$  is a program which takes inputs of type  $A$  and returns outputs of type  $B$ , with no linearity constraint (so it is like a normal program in the lambda calculus). Usually such programs are constructed by first making copies of the input and then feeding those copies into a program, like those in (2.14), (2.15), which is linear with respect to its inputs. To give such an example we need the *dereliction* operator  $\text{derelict}(-)$  appearing in (2.16) which takes an input which may be used arbitrarily ( $z$  in our example) and returns a restricted form of the input which may only be used linearly.



The simplest nontrivial example is the term (with  $q : !(A \multimap A), z : A$ )

$$(2.17) \quad \underline{2} = (\lambda q. (\text{copy } q \text{ as } q', q'' \text{ in } (\lambda z. (\text{derelict}(q'') (\text{derelict}(q') z))))))$$

which is of type

$$\underline{2} : \mathbf{int}_A = !(A \multimap A) \multimap (A \multimap A).$$

Intuitively, this program takes an input  $q : !(A \multimap A)$ , copies it to form  $q', q''$ , and returns

$$(\lambda z. (\text{derelict}(q'') (\text{derelict}(q') z)))$$

which is itself a program, sending  $z : A$  to what is intuitively the result of evaluating  $q$  twice on  $z$ . The correct way to make this precise is as shown above, using the dereliction operator. In more informal language, the overall program  $\underline{2}$  has the meaning  $q \mapsto q^2$ .

Every term in intuitionistic logic can be interpreted as a linear map between (usually infinite-dimensional) vector spaces, and it is these interpretations (or *denotations*, in the logic terminology) which will be used in the recurrent neural networks.

### 2.4.3 Semantics

The semantics  $\llbracket - \rrbracket$  of linear logic in the category  $\mathcal{V}$  of vector spaces [22, §5.1, §5.3] assigns to each type  $A$  a vector space  $\llbracket A \rrbracket$  and to each term  $\pi : A$  a vector  $\llbracket \pi \rrbracket \in \llbracket A \rrbracket$  in a way which is compositional with respect to the structure of types and terms. We refer to the objects  $\llbracket A \rrbracket, \llbracket \pi \rrbracket$  as the *denotations* of  $A, \pi$  respectively.

The semantics depends on a choice, for each atomic type  $\phi_i$ , of a vector space  $\llbracket \phi_i \rrbracket$ . This choice is made freely and once made, determines all denotations, according to the following inductive definitions:

$$\begin{aligned} \llbracket A \multimap B \rrbracket &= \text{Hom}_k(\llbracket A \rrbracket, \llbracket B \rrbracket), \\ \llbracket A \otimes B \rrbracket &= \llbracket A \rrbracket \otimes \llbracket B \rrbracket, \\ \llbracket !A \rrbracket &= !\llbracket A \rrbracket, \end{aligned}$$

where  $\text{Hom}_{\mathbb{R}}(W, V)$  is the space of linear maps from  $W$  to  $V$  (we will also use the notation  $\text{End}_{\mathbb{R}}(V) = \text{Hom}_{\mathbb{R}}(V, V)$ ). Here  $!V$  denotes the universal cocommutative counital coalgebra mapping to  $V$ , or what is the same, the space of distributions with finite support on (a finite-dimensional subspace of)  $V$  viewed as a smooth manifold. This is a bit complicated to describe, but it suffices here to know that there is a canonical (non-linear!) function

$$(2.18) \quad V \longrightarrow !V, \quad u \longmapsto |\emptyset\rangle_u$$

sending a vector  $u \in V$  to the Dirac distribution concentrated at  $u$ .

**Example 2.7.** The vector space

$$\begin{aligned} \llbracket !(\phi_1 \multimap \phi_1) \multimap (\phi_1 \multimap \phi_1) \rrbracket &= \text{Hom}_{\mathbb{R}}(\llbracket !(\phi_1 \multimap \phi_1) \rrbracket, \llbracket \phi_1 \multimap \phi_1 \rrbracket) \\ &= \text{Hom}_{\mathbb{R}}(!\text{End}_{\mathbb{R}}(\llbracket \phi_1 \rrbracket), \text{End}_{\mathbb{R}}(\llbracket \phi_1 \rrbracket)) \end{aligned}$$

is determined once  $\llbracket \phi_1 \rrbracket$  is. The construction can be understood in terms of the tree (??) by starting with the vector space  $\llbracket \phi_1 \rrbracket$  at the leaves and combining these vector spaces according to the operations at the vertices.

Having explained how the denotations of types are defined, we now turn to denotations of terms. For any term  $\pi$  of type  $A$  the denotation is a vector  $\llbracket \pi \rrbracket \in \llbracket A \rrbracket$ . In particular, if  $\rho$  is of type  $A \multimap B$  then its denotation is a vector

$$\llbracket \rho \rrbracket \in \llbracket A \multimap B \rrbracket = \text{Hom}_{\mathbb{R}}(\llbracket A \rrbracket, \llbracket B \rrbracket)$$

that is, the denotation is a linear map  $\llbracket A \rrbracket \longrightarrow \llbracket B \rrbracket$ . Similarly, if  $\gamma$  is of type  $!A \multimap B$  then its denotation is a linear map

$$\llbracket \gamma \rrbracket : \llbracket !A \rrbracket \longrightarrow \llbracket B \rrbracket.$$

In Remark 2.6 we explained how to think of both  $\rho, \gamma$  as programs taking inputs of type  $A$ , but where  $\rho$  is restricted to use this input linearly. In the semantics this resolves to the distinction between  $\llbracket A \rrbracket$  and  $\llbracket !A \rrbracket$ . We can still think of  $\llbracket \gamma \rrbracket$  as taking inputs of type  $A$  (i.e. vectors in  $\llbracket A \rrbracket$ ) using (2.18) to define the (non-linear) composite

$$(2.19) \quad \llbracket A \rrbracket \longrightarrow \llbracket !A \rrbracket \xrightarrow{\llbracket \gamma \rrbracket} \llbracket B \rrbracket$$

which feeds  $u \in \llbracket A \rrbracket$  to  $\llbracket \gamma \rrbracket$ . This composite is denoted  $\llbracket \gamma \rrbracket_{nl}$ , see [22, Definition 5.10].

We will not describe here the precise rules which construct the denotations of terms, but illustrate the idea with examples from [34, §3]. We fix a type  $A$  and set  $V = \llbracket A \rrbracket$ .

**Example 2.8.** The denotation of a term  $\pi$  of type  $\mathbf{int}_A$  is by Example 2.7 a vector

$$\llbracket \pi \rrbracket \in \llbracket \mathbf{int}_A \rrbracket = \text{Hom}_{\mathbb{R}}(!\text{End}_{\mathbb{R}}(V), \text{End}_{\mathbb{R}}(V)).$$

Thus  $\llbracket \pi \rrbracket$  is a linear map  $!\text{End}_{\mathbb{R}}(V) \longrightarrow \text{End}_{\mathbb{R}}(V)$ . For each integer  $n \geq 0$  there is a term  $\underline{n}$  of type  $\mathbf{int}_A$ , with  $\underline{2}$  being given already in (??). One calculates that for  $\alpha \in \text{End}_{\mathbb{R}}(V)$ ,

$$\llbracket \underline{n} \rrbracket |\emptyset\rangle_{\alpha} = \alpha^n.$$

Another way to express this is to say that the non-linear function

$$\llbracket \pi \rrbracket_{nl} : \text{End}_{\mathbb{R}}(V) \longrightarrow \text{End}_{\mathbb{R}}(V)$$

satisfies  $\llbracket \underline{n} \rrbracket_{nl}(\alpha) = \alpha^n$ .

**Example 2.9.** The denotation of a proof  $\pi$  of type  $\mathbf{bint}_A$  gives a function

$$\llbracket \pi \rrbracket_{nl} : \text{End}_{\mathbb{R}}(V) \times \text{End}_{\mathbb{R}}(V) \longrightarrow \text{End}_{\mathbb{R}}(V).$$

For every binary sequence  $S \in \{0, 1\}^*$  there is a corresponding term  $\underline{S}$  of type  $\mathbf{bint}_A$  [34, §3.2]. For example

$$\underline{001} = (\lambda q. (\lambda p. (\text{copy } q \text{ as } r, s \text{ in } (\lambda z. (\text{derelict}(p) (\text{derelict}(s) (\text{derelict}(r) z))))))).$$

The denotation  $\llbracket \underline{S} \rrbracket_{nl}$  sends a pair of linear operators  $\alpha, \beta$  to the product described by  $S$ , reading  $\alpha$  for 0 and  $\beta$  for 1, and reading in reverse order. For example,

$$\llbracket \underline{001} \rrbracket_{nl}(\alpha, \beta) = \beta \alpha \alpha.$$

### 3 The Linear Logic Recurrent Neural Network

There is a class of augmented RNNs with the following general form: at each time step a linear operator on the hidden state space

$$Z = Z(h^{(t)}, x^{(t+1)}) \in \text{End}_{\mathbb{R}}(\mathcal{H})$$

is generated from the previous hidden state and current input, and then this linear operator  $Z$  is applied to the hidden state to generate a vector  $Z(h^{(t)})$  which is inserted into the RNN evolution equation for  $h^{(t+1)}$ . Both the second-order RNN (??) and NTM (??) have this general form. What is the *generic* augmented RNN of this kind?

Obviously  $Z$  should be a (piece-wise) differentiable function of the hidden state and current input, but we consider a more restricted class of functions constructed as follows: in the first step, using feed-forward neural networks, we generate from  $h^{(t)}, x^{(t+1)}$  a family of vectors (keeping in mind that tensors and linear operators are also vectors). In the second step, these vectors are combined algorithmically using iteration and the basic operations of linear algebra to produce the linear operator  $Z$ . The vectors generated in the first step are called *command vectors*, *data vectors* or *input vectors*, depending on how they are constructed (see below for details). The algorithmic construction in the second step is a fixed linear logic program which we refer to as the *master algorithm*. The generic augmented RNN defined in this way is called the Linear Logic RNN (LLRNN).

In fact it is natural to consider two linear operators  $Z_{\text{in}}, Z_{\text{out}}$ , with the former appearing within the evolution equation *inside* the nonlinearity and the latter appearing *outside* the nonlinearity. Thus the evolution equation of the LLRNN model is

$$(3.1) \quad h^{(t+1)} = \sigma\left(Z_{\text{in}}(h^{(t)}) + Hh^{(t)} + Ux^{(t+1)} + B\right) + Z_{\text{out}}(h^{(t)})$$

where the linear operators  $Z_z \in \text{End}_{\mathbb{R}}(\mathcal{H})$  for  $z \in \{\text{in}, \text{out}\}$  depend on  $h^{(t)}, x^{(t+1)}$  via the aforementioned feed-forward networks

$$\begin{aligned} p_i^{(t+1)} &= \zeta(W_i^p h^{(t)} + B_i^p) \in \mathcal{P}_i, & 1 \leq i \leq r, \\ b_j^{(t+1)} &= \sigma(W_j^b h^{(t)} + B_j^b) \in \mathcal{B}_j, & 1 \leq j \leq s, \\ c^{(t+1)} &= W^c x^{(t+1)} \in \mathcal{C}. \end{aligned}$$

and the denotation of a term master<sup>z</sup> in linear logic, via the formula

$$(3.2) \quad Z_z = \llbracket \text{master}^z \rrbracket_{nl} \left( p_1^{(t+1)}, \dots, p_r^{(t+1)}, b_1^{(t+1)}, \dots, b_s^{(t+1)}, c^{(t+1)} \right).$$

Here we use that the denotation  $\llbracket \text{master}^z \rrbracket$  defines a differentiable function

$$(3.3) \quad \mathcal{P}_1 \times \dots \times \mathcal{P}_r \times \mathcal{B}_1 \times \dots \times \mathcal{B}_s \times \mathcal{C} \longrightarrow \text{End}_{\mathbb{R}}(\mathcal{H}).$$

We give the vectors generated in the first step names; we have

- A sequence of *command vectors*  $p_i^{(t+1)}$  generated via softmax  $\zeta$  from  $h^{(t)}$ ,
- a sequence of *data vectors*  $b_i^{(t+1)}$  generated via  $\sigma$  from  $h^{(t)}$ ,
- an *input vector*  $c^{(t+1)}$  generated via  $\sigma$  from  $x^{(t+1)}$ .

The idea is that the command vectors give distributions over finite-dimensional spaces of linear logic programs, while the data vectors and input vector give the inputs to be fed in some way into the selected programs; the exact details of this “feeding” are specified by the master algorithm.

**Definition 3.1.** 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 *data types*;
- a type  $C$ , called the *input type*;
- a type  $A$  with  $\llbracket A \rrbracket = \mathcal{H}$ , the hidden state space;
- two proofs  $\underline{\text{master}}^z$  for  $z \in \{\text{in}, \text{out}\}$  of the sequent

$$(3.4) \quad !P_1, \dots, !P_r, !B_1, \dots, !B_s, !C \vdash A \multimap A.$$

- Finite-dimensional subspaces for  $1 \leq i \leq r$  and  $1 \leq j \leq s$

$$(3.5) \quad \mathcal{P}_i \subseteq \llbracket P_i \rrbracket, \quad \mathcal{B}_j \subseteq \llbracket B_j \rrbracket, \quad \mathcal{C} \subseteq \llbracket C \rrbracket$$

which are spanned by the denotations of linear logic proofs. We refer to the elements of  $\mathcal{P}_i$  as *command vectors*, elements of  $\mathcal{B}_j$  as *data vectors* and elements of  $\mathcal{C}$  as *input vectors*.

The fundamental fact which makes this model reasonable is the following:

**Proposition 3.2.** *The functions  $\llbracket \underline{\text{master}}^z \rrbracket_{nl}$  are smooth.*

*Proof.* This follows from the hypothesis that  $\mathcal{P}_i, \mathcal{B}_j$  are generated by denotations of linear logic proofs, and the smoothness of these denotations [34].  $\square$

Moreover, the derivatives of  $\llbracket \underline{\text{master}}^z \rrbracket_{nl}$  can be computed symbolically using the cut-elimination algorithm of differential linear logic, and it is therefore feasible to implement a general LLRNN in a software package like TensorFlow.

**Remark 3.3.** It will often be convenient to write  $\Gamma = P_1, \dots, P_r$  for the command types and  $\Delta = B_1, \dots, B_s$  for the data types. A standard notational device in the linear logic literature is to then write  $!\Gamma$  for the list prepending  $!$  to all the types in the list  $\Gamma$ , that is,  $!\Gamma = !P_1, \dots, !P_r$ . With this notation,  $\underline{\text{master}}$  is a proof of  $!\Gamma, !\Delta, !V \vdash A \multimap A$ .

### 3.1 Examples

As defined the LLRNN is a class of models. In this section we explain how to implement various existing augmentations of RNNs, including a subset of the Neural Turing Machine, in the framework of the LLRNN. The guiding intuition from logic is that 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 [22, §7]. From this point of view, the purpose of augmenting an RNN with linear logic is to provide access to iterators of a complexity “similar” to the function that the neural network is trying to approximate.

**Example 3.4 (Second-order RNN).** With  $r = s = 0$  there are no command or data vectors. Take for the input type  $C = A \multimap A$  and

$$\mathcal{C} = \text{End}_{\mathbb{R}}(\mathcal{H}) = \llbracket A \multimap A \rrbracket.$$

so the master algorithm takes a single input, which is a linear operator on the hidden state, and returns such an operator. The only weight matrix involved in the LLRNN beyond the usual RNN is the matrix  $W^c$  which maps inputs to linear operators on  $\mathcal{H}$ .

We choose  $\underline{\text{master}}^{\text{in}}$  to be the proof of  $!(A \multimap A) \vdash A \multimap A$  which is given by dereliction, so that  $\llbracket \underline{\text{master}}^{\text{in}} \rrbracket_{nl}(\alpha) = \alpha$  and hence

$$(3.6) \quad Z_{\text{in}} = W^c(x^{(t+1)}).$$

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

The most elementary coupling of an RNN to linear logic adds the ability to raise linear operators (generated say from an input symbol, in the manner of the second-order RNN) to a power generated from the hidden state of the RNN.

**Example 3.5 (Higher-order RNN).** Consider the generalisation of the second-order RNN where the controller predicts at each time step an integer power of the linear operator  $W^c(x^{(t+1)})$  to apply to the hidden state. Suppose we allow powers in the range  $\{0, \dots, L\}$ . Then at each time step the RNN will generate a distribution  $p^{(t+1)}$  over  $\{0, \dots, L\}$  from the current hidden state  $h^{(t)}$  by the formula

$$p^{(t+1)} = \zeta(W_p h^{(t)} + B_p)$$

and therefore

$$(3.7) \quad Z_{\text{in}} = \sum_{i=0}^L p_i^{(t+1)} (W^c(x^{(t+1)}))^i.$$

The operation of taking a linear operator and raising it to the  $n$ th power is encoded by the proof  $\underline{n}$  of type  $\mathbf{int}_A$ . We can therefore represent the higher-order RNN as a LLRNN,

as follows. There is one command type  $\mathbf{int}_A$ , no data types, and input type  $A \multimap A$ . The spaces of command and input vectors are respectively

$$\begin{aligned}\mathcal{P}_1 &= \text{span}(\llbracket 0 \rrbracket, \dots, \llbracket L \rrbracket) \subseteq \llbracket \mathbf{int}_A \rrbracket, \\ \mathcal{C} &= \text{End}_{\mathbb{R}}(\mathcal{H}) = \llbracket A \multimap A \rrbracket.\end{aligned}$$

We omit  $\underline{\text{master}}^{\text{out}}$  and take  $\underline{\text{master}}^{\text{in}}$  to be the proof of

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

given in the term calculus of [42] by

$$\underline{\text{master}}^{\text{in}} = (\lambda n. (\lambda a. (\text{derelict}(n) a))).$$

This proof has the property that for  $n_i \geq 0$  and  $\alpha \in \text{End}_{\mathbb{R}}(\mathcal{H})$  and  $p_i \in \mathbb{R}$ ,

$$(3.8) \quad \llbracket \underline{\text{master}}^{\text{in}} \rrbracket_{nl} \left( \sum_i \lambda_i \llbracket n_i \rrbracket, \alpha \right) = \sum_i \lambda_i \alpha^{n_i}.$$

The coupling of this function to the RNN is via  $p_1^{(t+1)}$  which we identify with  $p^{(t+1)}$  above, and  $c^{(t+1)} = W^c(x^{(t+1)})$ . By construction the LLRNN with this configuration reproduces the  $Z$  of (3.7) and thus the higher-order RNN.

Since one doesn't need to know linear logic to understand how to raise a linear operator to an integer power, it is natural to wonder to what degree linear logic is actually necessary here. All of our models can be formulated without any mention of linear logic, and indeed we will generally present the functions  $\llbracket \underline{\text{master}} \rrbracket_{nl}$  rather than the underlying proof. Nonetheless, the construction of  $\underline{\text{master}}$  within linear logic constrains the model and makes conceptually clear the computational ideas involved; these ideas are not necessarily clear from the polynomial algebra that results from applying the denotation functor  $\llbracket - \rrbracket$ . The strongest example of this point of view is the role of iteration in the LLRNN approach to the Neural Turing Machine and its generalisations.

### 3.2 NTM as an LLRNN

We now explain how to present the NTM of Section 2.3 as an LLRNN. We first give the command, data and input types. Let  $V, W$  be types of linear logic with  $\llbracket V \rrbracket = \mathcal{V}$ ,  $\llbracket W \rrbracket = \mathcal{W}$ , and write  $W^\vee = W \multimap 1$  so that  $\llbracket W^\vee \rrbracket = \mathcal{W}^*$ . The command types  $\Gamma$  and data types  $\Delta$  are defined to be respectively

$$\Gamma = \mathbf{int}_W, \mathbf{int}_{W^\vee}, \quad \Delta = W \multimap W, W^\vee \multimap W^\vee, V \multimap V, V.$$

The master algorithms are the proofs of  $!\Gamma, !\Delta \vdash A \multimap A$  such that for inputs

$$\begin{aligned}\llbracket m \rrbracket &\in \llbracket \mathbf{int}_W \rrbracket, \quad \llbracket n \rrbracket \in \llbracket \mathbf{int}_{W^\vee} \rrbracket \\ \alpha &\in \llbracket W \multimap W \rrbracket = \text{End}_{\mathbb{R}}(\mathcal{W}) \\ \beta &\in \llbracket W^\vee \multimap W^\vee \rrbracket = \text{End}_{\mathbb{R}}(\mathcal{W}^*) \\ e &\in \llbracket V \multimap V \rrbracket = \text{End}_{\mathbb{R}}(\mathcal{V}) \\ a &\in \llbracket V \rrbracket = \mathcal{V}\end{aligned}$$

we have, for  $h = (h_0, r, w, M) \in \mathcal{H}$  decomposed according to (2.8),

$$\begin{aligned}\llbracket \underline{\text{master}}^{\text{in}} \rrbracket_{nl}(\llbracket m \rrbracket, \llbracket n \rrbracket, \alpha, \beta, e, a)(h) &= (M(r), 0, 0, 0), \\ \llbracket \underline{\text{master}}^{\text{out}} \rrbracket_{nl}(\llbracket m \rrbracket, \llbracket n \rrbracket, \alpha, \beta, e, a)(h) &= (0, \alpha^m(r), \beta^n(w), (1_{\mathcal{V}} - e) \circ M + w \otimes a).\end{aligned}$$

As before we take  $\mathcal{V} = \mathcal{H}$ . We take

$$\begin{aligned}\mathcal{P}_1 &= \text{span}(\llbracket 0 \rrbracket, \dots, \llbracket N-1 \rrbracket) \subseteq \llbracket \mathbf{int}_W \rrbracket, \\ \mathcal{P}_2 &= \text{span}(\llbracket 0 \rrbracket, \dots, \llbracket N-1 \rrbracket) \subseteq \llbracket \mathbf{int}_{W^\vee} \rrbracket, \\ \mathcal{B}_1 &= \llbracket W \multimap W \rrbracket = \text{End}_{\mathbb{R}}(\mathcal{W}), \\ \mathcal{B}_2 &= \llbracket W^\vee \multimap W^\vee \rrbracket = \text{End}_{\mathbb{R}}(\mathcal{W}^*), \\ \mathcal{B}_3 &= \llbracket V \multimap V \rrbracket = \text{End}_{\mathbb{R}}(\mathcal{V}) \\ \mathcal{B}_4 &= \llbracket V \rrbracket = \mathcal{V}.\end{aligned}$$

We restrict the functions  $\llbracket \underline{\text{master}}^z \rrbracket_{nl}$  to the subset of inputs where  $\alpha = R, \beta = R^*$ . The coupling of this restricted function to the RNN is via command vectors  $p_1^{(t+1)}, p_2^{(t+1)}$  giving distributions over the basis  $\{\llbracket i \rrbracket\}_{i=0}^{N-1}$  of  $\mathcal{P}_1, \mathcal{P}_2$  which we identify respectively with  $q^{(t+1)}$  (the distribution over powers of  $R$  used to manipulate the read address) and  $s^{(t+1)}$  (the distribution over powers of  $R^*$  used to manipulate the write address).

We assume the weight matrix  $H$  of (3.1) is the projection from  $\mathcal{H}$  to  $\mathcal{H}_0$  followed by the weight  $H_0$  above, and similarly for  $U, B$ . Then with the current notation the evolution equation (3.1) of the LLRNN reads

$$\begin{aligned}(h_0^{(t+1)}, r^{(t+1)}, w^{(t+1)}, M^{(t+1)}) &= \left( \sigma(M^{(t)}(r^{(t)}) + H_0 h_0^{(t)} + U_0 x^{(t+1)} + B_0), \right. \\ &\quad \sum_{i=0}^{N-1} (p_1^{(t+1)})_i R^i(r^{(t)}), \sum_{i=0}^{N-1} (p_2^{(t+1)})_i (R^*)^i(w^{(t)}), \\ &\quad \left. (1_{\mathcal{V}} - e^{(t+1)}) \circ M^{(t)} + w^{(t)} \otimes a \right),\end{aligned}$$

which agrees with the equations (2.9) – (2.12).

**Example 3.6 (Dihedral NTM).** The NTM manipulates its memory state via rotations of the regular  $N$ -gon. In this example we study the natural extension which allows access the full symmetry group, the dihedral group, by adding the reflection

$$\begin{aligned}R' : \mathcal{W} &\longrightarrow \mathcal{W}, \\ R'(\bar{a}) &= \overline{-a}.\end{aligned}$$

Note that  $R'$  and  $R$  do not commute. The command and data types are now

$$\begin{aligned}\Gamma &= \mathbf{bint}_W, \mathbf{bint}_{W^\vee}, \\ \Delta &= W \multimap W, W \multimap W, W^\vee \multimap W^\vee, W^\vee \multimap W^\vee, V \multimap V, V\end{aligned}$$

and  $\underline{\text{master}}^{\text{out}}$  is the proof such that for  $h \in \mathcal{H}$  and  $F, G \in \{0, 1\}^*$

$$\begin{aligned} & \llbracket \underline{\text{master}}^{\text{out}} \rrbracket_{nl}(\llbracket F \rrbracket, \llbracket G \rrbracket, \alpha_1, \alpha_2, \beta_1, \beta_2, e, a)(h) \\ &= \left( M(r), \llbracket F \rrbracket_{nl}(\alpha_1, \alpha_2)(r), \llbracket G \rrbracket_{nl}(\beta_1, \beta_2)(w), (1_{\mathcal{V}} - e) \circ M + w \otimes a \right). \end{aligned}$$

We then fix  $\alpha_1 = R, \alpha_2 = R'$  and  $\beta_1 = R^*, \beta_2 = (R')^*$ .

## 4 Generalisations of the NTM

Taking the linear logic point of view on the NTM from Section 3.2 as a starting point, we develop several generalisations. The main purpose of these models is to demonstrate how a high-level idea for extending the basic NTM model can be realised by first translating the idea into linear logic, and then using the LLRNN form of the NTM.

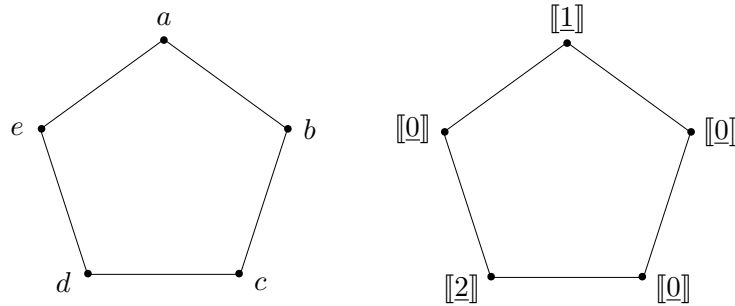
### 4.1 Pattern NTM

At each time step the NTM predicts distributions  $s^{(t)}, q^{(t)}$  over the possible rotations of the write and read addresses. A pattern of memory accesses is a *sequence* of such rotations, and since the “angle” of rotation is represented in the LLRNN as the denotation of a linear logic proof (with  $\llbracket n \rrbracket$  representing a clockwise angle of  $\frac{2\pi n}{5}$ ) the pattern may be stored as a sequence of vectors in  $\llbracket \text{int}_W \rrbracket$ . It is therefore natural to add a memory ring with coefficients in this vector space, which in principle should allow the NTM to learn more complex patterns of memory access. In this example we first present the construction using linear logic, and then afterward translate this into explicit equations involving matrices.

Let  $\mathbf{a} = \{a, b, \dots, z\}$  and let  $\mathcal{V}_1 = \oplus_{i \in \mathbf{a}} \mathbb{R}i$  be the free vector space on  $A$ . We also take a finite-dimensional subspace  $\mathcal{V}_2 \subseteq \llbracket \text{int}_W \rrbracket$ . For convenience we keep the same address space for both memory rings. So we have memory spaces for  $i \in \{1, 2\}$

$$\mathcal{S}_i = \mathcal{W}^* \otimes \mathcal{V}_i \cong \text{Hom}_{\mathbb{R}}(\mathcal{W}, \mathcal{V}_i).$$

An example of a combined state  $\mathcal{S} = \mathcal{S}_1 \oplus \mathcal{S}_2$  of these two memory spaces is





We consider a modification of the NTM in which the contents of the second ring control the motion of the read address of the first ring (rather than this being manipulated directly by the RNN controller). To demonstrate the dynamics in an example, suppose that the memory state is as shown at time  $t = t_0$ , that the read address for the first ring is focused at the zero position at time  $t_0 - 1$ , and that the read address for the second ring is focused at zero at time  $t_0$  and increases its focus position by one step over each time interval.

Then as time increases from  $t_0$  the repeating sequence

$$(4.1) \quad \llbracket 1 \rrbracket, \llbracket 0 \rrbracket, \llbracket 0 \rrbracket, \llbracket 0 \rrbracket, \llbracket 2 \rrbracket, \llbracket 0 \rrbracket, \dots$$

will be used to control the read address of the first memory ring. More precisely, the operator applied to the read address  $r_1 \in \mathcal{W}$  will be  $R^n$  with  $n$  varying over the repeating sequence 1, 0, 0, 2, 0. Finally, the sequence of vectors delivered to the controller will be

$$a, b, b, b, d, d, e, e, e, b, \dots$$

In more detail, as part of master there will be an evaluation program

$$\begin{array}{c} \underline{\text{eval}} \\ \vdots \\ \mathbf{int}_W, !(W \multimap W) \vdash W \multimap W \end{array}$$

which has the property that

$$\llbracket \underline{\text{eval}} \rrbracket \left( \sum_i \lambda_i \llbracket n_i \rrbracket, |\emptyset\rangle_\alpha \right) = \sum_i \lambda_i \alpha^{n_i}.$$

The first argument to this evaluation function will be the state  $M_2(r_2) = \sum_i \lambda_i \llbracket n_i \rrbracket$  of the second memory ring. We apply the output to the read address  $r_1$  of the first memory ring and specialise to  $\alpha = R$  so that the relevant expression is

$$\llbracket \underline{\text{eval}} \rrbracket (M_2(r_2), |\emptyset\rangle_R)(r_1) = \sum_i \lambda_i R^{n_i}(r_1).$$

The command and data types are

$$\begin{aligned} \Gamma &= \mathbf{int}_{W^\vee}, \mathbf{int}_W, \mathbf{int}_{W^\vee} \\ \Delta &= W \multimap W, W^\vee \multimap W^\vee, W \multimap W, W^\vee \multimap W^\vee, V \multimap V, V, \mathbf{int}_W \multimap \mathbf{int}_W, \mathbf{int}_W. \end{aligned}$$

The first type in  $\Gamma$  is connected to the write address for the first memory ring, while the second and third types relate respectively to the read and write address of the second memory ring. We do not include an input to the master algorithm for manipulating the read address of the first memory ring since this is purely under the control of the second memory ring. The inputs of type  $V \multimap V, V$  in  $\Delta$  are the erase and add vector for the

first memory ring, while the inputs  $\mathbf{int}_W \multimap \mathbf{int}_W, \mathbf{int}_W$  give the erase and add vectors for the second memory ring. As usual  $\underline{\text{master}}^z$  is a proof of  $!\Gamma, !\Delta \vdash A \multimap A$ .

Given the command and data vectors

$$g = (\llbracket \underline{n}_1 \rrbracket, \llbracket \underline{m}_2 \rrbracket, \llbracket \underline{n}_2 \rrbracket, \alpha_1, \beta_1, \alpha_2, \beta_2, e_1, a_1, e_2, a_2) \in \bigoplus_{i=1}^3 \mathcal{P}_i \oplus \bigoplus_{j=1}^8 \mathcal{B}_j$$

and state vector

$$(4.2) \quad h = (h_0, r_1, w_1, r_2, w_2, M_1, M_2) \in \mathcal{H} = \mathcal{H}_0 \oplus \mathcal{W} \oplus \mathcal{W}^* \oplus \mathcal{W} \oplus \mathcal{W}^* \oplus \mathcal{S}_1 \oplus \mathcal{S}_2$$

the value of the master functions are

$$\llbracket \underline{\text{master}}^{\text{in}} \rrbracket_{nl}(g)(h) = (Q(M_1(r_1)), 0, 0, 0, 0, 0, 0).$$

where  $Q : \mathcal{V}_1 \longrightarrow \mathcal{H}_0$  is a weight matrix, and

$$\begin{aligned} \llbracket \underline{\text{master}}^{\text{out}} \rrbracket_{nl}(g)(h) = & \left( 0, \llbracket \text{eval} \rrbracket(M_2(r_2), |\emptyset\rangle_{\alpha_1})(r_1), \beta_1^{n_1}(w_1), \alpha_2^{m_2}(r_2), \beta_2^{n_2}(w_2), \right. \\ & \left. (1_{\mathcal{V}_1} - e_1) \circ M_1 + \beta_1^{n_1}(w_1) \otimes a_1, (1_{\mathcal{V}_2} - e_2) \circ M_2 + w_2 \otimes a_2 \right). \end{aligned}$$

As above we set  $\alpha_1 = \alpha_2 = R$  and  $\beta_1 = \beta_2 = R^*$ . We take  $\mathcal{B}_8 = \mathcal{V}_2 \subseteq \llbracket \mathbf{int}_W \rrbracket$  and identify  $\mathcal{B}_7 \subseteq \llbracket \mathbf{int}_W \multimap \mathbf{int}_W \rrbracket$  with  $\text{End}_{\mathbb{R}}(\mathcal{V}_2)$  in an appropriate way.

**Remark 4.1.** The pattern NTM can be presented without linear logic as follows. Suppose we allow rotations given by powers of  $R$  in the set  $\{0, \dots, L\}$ . Then  $\mathcal{V}_2 = \mathbb{R}^{L+1}$  and if we take  $N = 5$  as in the above pictures, then  $\mathcal{W} = \mathbb{R}^5$ . The system is a modified RNN with state space (4.2), and equations (for  $1 \leq i \leq 2$  in all cases)

$$(4.3) \quad s_i^{(t+1)} = \zeta(W_{i,s}h^{(t)} + B_{i,s}) \in \mathcal{W}^* \cong \mathbb{R}^5$$

$$(4.4) \quad q_2^{(t+1)} = \zeta(W_q h^{(t)} + B_q) \in \mathcal{W}^* \cong \mathbb{R}^5,$$

$$(4.5) \quad e_i^{(t+1)} = \sigma(W_{i,e}h^{(t)} + B_{i,e}) \in \text{End}_{\mathbb{R}}(\mathcal{V}_i),$$

$$(4.6) \quad a_i^{(t+1)} = \sigma(W_{i,a}h^{(t)} + B_{i,a}) \in \mathcal{V}_i,$$

$$(4.7) \quad w_i^{(t+1)} = \sum_{i=0}^{N-1} s^{(t+1)}(\bar{i}) \cdot (R^*)^i(w_i^{(t)}),$$

$$(4.8) \quad r_1^{(t+1)} = \sum_{j=0}^L M_2^{(t)}(r_2^{(t)})_j \cdot R^j(r_1^{(t)}),$$

$$(4.9) \quad r_2^{(t+1)} = \sum_{i=0}^{N-1} q_2^{(t+1)}(\bar{i}) \cdot R^i(r_2^{(t)}),$$

$$(4.10) \quad M_i^{(t+1)} = (1_{\mathcal{V}_i} - e_i^{(t+1)}) \circ M_i^{(t)} + w_i^{(t)} \otimes a_i^{(t+1)}.$$

Finally, the evolution equation is

$$(4.11) \quad h_0^{(t+1)} = \sigma \left( Q(M_1^{(t)}(r_1^{(t)})) + H_0 h_0^{(t)} + U_0 x^{(t+1)} + B_0 \right).$$

The main difference between the ordinary NTM and the pattern NTM is in (4.8) where the content  $M_2(r_2) \in \mathbb{R}^{L+1}$  of the second memory ring is used to determine the coefficients of each rotation  $R^j$  of the first memory ring's read address.

## 4.2 Multiple pattern NTM

In the previous example the second memory ring stores patterns of memory access in the form of sequences of integers, or more precisely, proofs of linear logic of type  $\mathbf{int}_W$ . More complex patterns (this word being of course a euphemism for *algorithms*) can be encoded using more rings and higher types, in various ways. In this section we give an example in which, in addition to the original memory ring from the standard NTM, there are two rings with coefficients in  $\llbracket \mathbf{int}_W \rrbracket$  and a ring with coefficients in  $\llbracket \mathbf{bint}_W \rrbracket$ . The idea is that the second and third memory rings encode two patterns of memory access, and we allow the controller to switch between the two patterns at any time using the fourth ring. For ease of reference in the sequel we call this the *multiple pattern NTM*.

Before getting into the details, let us give an example of the kind of function this new model is designed to learn. Let  $\mathbf{a} = \{a, b, \dots, z\} \cup \{\circ, \bullet, S\}$  and  $\mathcal{J} = \mathcal{V}_1 = \bigoplus_{i \in \mathbf{a}} \mathbb{R}i$  where  $S$  is the *swap symbol* used to switch between two learned patterns,  $\circ$  is the *initial symbol* used by the controller to initialise its memory rings, and  $\bullet$  is the *terminal symbol* used to switch the controller to “output” mode. Suppose we are trying to approximate

$$(4.12) \quad \circ abcdeSabcdeSabcde \bullet \mapsto abcde aabbccdde abcde,$$

where we insert spaces for clarity in the output. Without being too precise: one pattern is the identity, the other pattern doubles every digit, and at occurrences of  $S$  we switch between the patterns. This algorithm is realised using linear logic as follows: if at some time the read address of the second and third rings are focused at positions containing  $\llbracket i \rrbracket, \llbracket j \rrbracket$  respectively and the read address of the fourth ring is focused at  $\llbracket S \rrbracket$  for some sequence  $S \in \{0, 1\}^*$  (notation of Example 2.9) then the read address of the first ring will be acted on by  $R^p$  where

$$p = ai + bj$$

and  $a, b$  are the number of times 0, 1 appear in  $S$ , respectively. Thus whenever we see  $\llbracket 0 \rrbracket$  on the fourth ring we will follow the pattern on the second ring, and if we see  $\llbracket 1 \rrbracket$  we will follow the pattern on the third ring. An entry  $\llbracket 00 \rrbracket$  on the fourth ring will run the pattern on the second ring but with every “angle of rotation” doubled.

The proof that encodes this logic is

$$\begin{array}{c} \text{feed} \\ \vdots \\ \mathbf{int}_W, \mathbf{int}_W, \mathbf{bint}_W \vdash \mathbf{int}_W \end{array}$$

with linear lambda term

$$\begin{aligned} \underline{\text{feed}} = & (\lambda p q r a. (\text{copy } a \text{ as } a', a'' \text{ in} \\ & ((r \text{ (promote } a' \text{ for } x \text{ in } (p \ x))) \\ & (\text{promote } a'' \text{ for } x' \text{ in } (q \ x')))))) . \end{aligned}$$

We will not spell out the full model here since it is essentially the same as the one given in the previous section, with the main difference being that the second entry in the tuple defining  $\llbracket \underline{\text{master}}^{\text{out}} \rrbracket_{nl}(g)(h)$  now looks like

$$(4.13) \quad \llbracket \underline{\text{feed}} \rrbracket (M_2(r_2), M_3(r_3), M_4(r_4))_{nl}(R)(r_1) .$$

To be more concrete, suppose  $\mathcal{V}_2 = \mathcal{V}_3 = \text{span}(\llbracket 0 \rrbracket, \dots, \llbracket L \rrbracket)$  and

$$\mathcal{V}_4 = \text{span}(\llbracket S_1 \rrbracket, \dots, \llbracket S_m \rrbracket)$$

where  $S_k \in \{0, 1\}^*$  for  $1 \leq k \leq m$ . We write  $a_k, b_k$  for, respectively, the number of 0's and 1's in  $S_k$ . At any given time the controller observes some linear combination of the contents of the rings

$$(4.14) \quad M_2(r_2) = \sum_{i=0}^L \lambda_i \llbracket i \rrbracket, \quad M_3(r_3) = \sum_{j=0}^L \mu_j \llbracket j \rrbracket, \quad M_4(r_4) = \sum_{k=0}^m \kappa_k \llbracket S_k \rrbracket .$$

In this case the term in (4.13) is

$$(4.15) \quad \sum_{i,j,k} \lambda_i \mu_j \kappa_k R^{a_k i + b_k j}(r_1)$$

so that the time evolution of the first read address will become (cf. (2.13))

$$(4.16) \quad r_1[u] \longleftarrow \sum_{i,j,k} \lambda_i \mu_j \kappa_k r_1[u - a_k i - b_k j] .$$

Returning to the example, set  $\mathcal{W}_i = (\mathbb{Z}/N_i\mathbb{Z})^{\mathbb{R}}$  with  $N_1 > 18$  and  $N_2 = N_3 = N_4 = 2$ . The particular values are not important, the memory rings need only be sufficiently large. One method (described anthropomorphically) for the controller to approximate (4.12) in the above system is to perform the following steps:

1. Conditioned on the initial symbol  $\circ$  (or a sequence of distinct initial symbols) the controller inserts  $\llbracket 1 \rrbracket$  at every position of the second memory ring and the repeating pattern  $\llbracket 0 \rrbracket, \llbracket 1 \rrbracket$  on both the third and fourth rings (note that in the first case these are of type  $\mathbf{int}_W$  and in the second case of type  $\mathbf{bint}_W$ ).
2. The controller writes the input sequence  $a, b, c, d, e, S, \dots, \bullet$  including the terminal symbol but not the initial symbol, onto the first memory ring.

3. The controller switches into output mode, with the read addresses of all memory rings sharply focused at the zero position. In particular, the fourth memory ring is focused at  $\llbracket 0 \rrbracket$  and the pattern on the second memory ring (the identity) is “active”.
4. The controller advances the read address of the second memory ring one position in each time step, yielding  $\llbracket 1 \rrbracket, \llbracket 1 \rrbracket, \dots$  and so advances the read address of the first memory ring one position per unit time, emitting  $V(a), V(b), V(c), V(d), V(e)$ .
5. The first encounter with the vector  $V(S)$  triggers the controller to
  - (i) advance the read address of the fourth memory ring so that the second pattern becomes “active”, and
  - (ii) begin advancing the read address of the third ring rather than the second.
6. The controller advances the read address of the third memory ring one position in each time step, yielding  $\llbracket 0 \rrbracket, \llbracket 1 \rrbracket, \llbracket 0 \rrbracket, \dots$  and so advances the read address of the first memory ring by zero positions (i.e. stays stationary), then by one position, then by zero, etc., emitting  $V(a), V(a), V(b), V(b), \dots$
7. Finally, when the controller emits the vector  $V(\bullet)$  we stop.

### 4.3 Polynomial step NTM

For our final example extends the pattern NTM of Section 4.1 with a modal symbol  $T$  which modifies the way the pattern is interpreted. The alphabet is therefore

$$(4.17) \quad \mathfrak{a} = \{a, b, \dots, z\} \cup \{\circ, \bullet, T\}$$

with  $\mathcal{S} = \mathcal{V}_1 = \bigoplus_{i \in \mathfrak{a}} \mathbb{R}i$  as before. The behaviour of the modal symbol  $T$  can be modelled in the LLRNN framework by adding a memory ring with coefficients in  $\llbracket !\mathbf{int}_W \multimap \mathbf{int}_W \rrbracket$ . Any polynomial  $a(x) \in \mathbb{N}[x]$  can be encoded by a term  $\underline{a}$  of  $!\mathbf{int}_W \multimap \mathbf{int}_W$ . We therefore have three memory rings with coefficient spaces

$$\mathcal{V}_1 = \bigoplus_{i \in \mathfrak{a}} \mathbb{R}i, \quad \mathcal{V}_2 \subseteq \llbracket \mathbf{int}_W \rrbracket, \quad \mathcal{V}_3 \subseteq \llbracket !\mathbf{int}_W \multimap \mathbf{int}_W \rrbracket.$$

Take  $\mathcal{V}_2$  as in Section 4.1, and

$$\mathcal{V}_3 = \text{span}(\llbracket \underline{a_1} \rrbracket, \dots, \llbracket \underline{a_k} \rrbracket)$$

for some polynomials  $a_1, \dots, a_k \in \mathbb{N}[x]$ . The relevant part of the master algorithm is

$$\begin{array}{c} \underline{\text{newfeed}} \\ \vdots \\ \mathbf{int}_W, \mathbf{!int}_W \multimap \mathbf{int}_W \vdash \mathbf{int}_W \end{array}$$

with linear lambda term

$$(4.18) \quad \underline{\text{newfeed}} = (\lambda p q x. ((q(\text{promote } p))x)).$$

This appears in  $\llbracket \underline{\text{master}}^{\text{out}} \rrbracket_{nl}(g)(h)$  in the form

$$(4.19) \quad \llbracket \underline{\text{newfeed}} \rrbracket(M_2(r_2), M_3(r_3))_{nl}(R)(r_1).$$

To give an example: if  $M_2(r_2) = \sum_i \lambda_i \llbracket i \rrbracket$  and  $M_3(r_3) = \sum_j \mu_j \llbracket a_j \rrbracket$ , then (4.19) is

$$(4.20) \quad \left\{ \sum_j \mu_j a_j \left( \sum_i \lambda_i R^i \right) \right\} (r_1).$$

Note that  $\text{End}_{\mathbb{R}}(\mathscr{W})$  is an algebra, so  $a_j(Y)$  is defined for any linear operator  $Y$ . If the third memory ring is sharply focused at the identity polynomial  $1 \in \mathbb{N}[x]$  then (4.20) is just  $\sum_i \lambda_i R^i(r_1)$ , so we shift the first memory ring according to the pattern on the second memory ring. If however the third memory ring is sharply focused at  $a(x) = x^2$  then for every step of size  $s \in \mathbb{N}$  dictated by the pattern on the second memory ring, we actually execute a step of size  $s^2$ .

### 4.3.1 Dihedral polynomial step NTM

There is a variant of the polynomial step NTM where we use binary integers instead of integer types, following Example 3.6. We keep the same alphabet, but take

$$\mathscr{V}_1 = \bigoplus_{i \in \mathbf{a}} \mathbb{R}i, \quad \mathscr{V}_2 \subseteq \llbracket \mathbf{bint}_W \rrbracket, \quad \mathscr{V}_3 \subseteq \llbracket !\mathbf{bint}_W \multimap \mathbf{bint}_W \rrbracket.$$

Take  $\mathscr{V}_2$  as in Section 4.1, and

$$\mathscr{V}_3 = \text{span}(\llbracket \underline{\text{id}} \rrbracket, \llbracket \underline{\text{repeat}} \rrbracket).$$

Where  $\underline{\text{id}}$  is the identity function, and  $\underline{\text{repeat}}$  repeats a binary integer [34, §3.2]. Recall that a binary integer, say 001, is interpreted as a word  $R'RR$  in the dihedral group and acts on the first memory ring according to the action of this group on  $\mathscr{W}$ . A sequence of binary integers stored on the second memory ring may therefore be thought of as a sequence of group elements  $g_0, \dots, g_4$ . If the third memory ring is sharply focused at the identity function, this sequence will be applied as it stands. If, however, the third memory ring is sharply focused at the repeat function, the sequence that gets applied will be  $g_0^2, \dots, g_4^2$  as the underlying binary sequence being repeated

$$\llbracket \underline{\text{repeat}} \rrbracket_{nl} \llbracket 001 \rrbracket = \llbracket 001001 \rrbracket$$

corresponds to the group elements being squared

$$R'RR \mapsto R'RRR'RR.$$

## 5 Unpolished remarks

**Remark 5.1.** A more powerful system would substitute a neural theorem prover along the lines of [30, 32, 44] in place of the libraries of functional programs  $\mathcal{P}_i$ . 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 [31].

**Remark 5.2.** In functional programming languages like differential  $\lambda$ -calculus [38] and differential linear logic [33] 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 a generalisation of automatic differentiation [37] which is widely used to compute derivatives of complicated real-valued functions, for example in the backpropagation algorithms of machine learning packages like TensorFlow [35, §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>2</sup> The idea of differential  $\lambda$ -calculus is similar, but more complex [39]. 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.

**Remark 5.3.** One method of adding algorithmic elements to deep neural networks that has been explored extensively is the idea of *attention networks* which implement a form of “distributed decision making” [2, p.1].

According to [2, §III.A] one of the problems with the original encoder-decoder network approach to machine translation [17] is that the model compresses the information in the input sequence into a vector without any presupposed structure; intuitively this makes the representations learned by the model harder to understand, and the output sequence harder to construct. From this perspective the LLRNN versions of the NTM models in previous sections can be thought of as structured encoder-decoder models with attention, where the representation learned by the encoder is the vector state of the memory rings after all input symbols have been read, as well as the matrix  $V$  mapping symbols to hidden state vectors. In our examples these vectors are highly structured, in the precise sense that the coefficients are highly correlated under the dynamics of the RNN.

These correlations ultimately arise from the fact that the coefficient spaces of the memory rings are spanned by linear logic denotations which imbue the coefficients with a dynamic meaning. Thus linear logic provides a structured algorithmic “scaffolding” which is populated by the encoder. For example, in Section 4.1 the meaning of the symbol  $S$  (which is implicit in the training data) is made explicit by learning the contents of the fourth memory ring, and the association between  $S$  and a manipulation of this ring, which is encoded by the coefficients of  $V(S)$ .

---

<sup>2</sup>See the “gradients” method of tensorflow/python/ops/gradients.py in TensorFlow 0.10

**Remark 5.4.** What is the link between attention mechanisms and the vanishing gradient problem? Clearly one of the purposes of attention mechanisms and external memories is to avoid the problem of “forgetting” large scale structure in the input to a recurrent network. The explicit way to do this is to restrict one part of the recurrent matrix  $H$  to be the identity [45]. More generally, we could restrict a part of the current matrix to be an element of the orthogonal group, i.e. a rotation. But then if look at that matrix, one part of it is precisely predicting a “rotation” to apply to the other part, and we are basically in the situation of a NTM.

We write  $\Phi = \llbracket \text{less} \rrbracket$  for the proof of Definition 5.7.

**Lemma 5.5.** *There is a proof more such that denotation gives a linear map*

$$\Psi = \llbracket \text{more} \rrbracket : \llbracket \mathbf{bint}_W \rrbracket \longrightarrow \llbracket \mathbf{bint}_{W \& W} \rrbracket$$

with the following properties:

(i)  $\Phi \circ \Psi = 1.$

(ii) *On the subspace of block matrices*

$$B = \left\{ \begin{pmatrix} \gamma_1 & 0 \\ 0 & \gamma_2 \end{pmatrix} \mid \gamma_i \in \text{End}_{\mathbb{R}}(\mathcal{W}) \right\} \subseteq \text{End}_{\mathbb{R}}(\mathcal{W} \oplus \mathcal{W})$$

we have

$$(5.1) \quad \Psi(\llbracket \underline{S}_W \rrbracket)|_{B \times B} = \llbracket \underline{S}_{W \& W} \rrbracket|_{B \times B}$$

for all  $S \in \{0, 1\}^*$ .

*Proof.* Let  $\pi_i, \iota_i$  denote the projections from and injections into  $\llbracket W \& W \rrbracket = \mathcal{W} \oplus \mathcal{W}$ . We define for  $p \in \llbracket \mathbf{bint}_W \rrbracket$

$$\Psi(p)(\alpha, \beta) = (p(\pi_1 \alpha \iota_1, \pi_1 \beta \iota_1), p(\pi_2 \alpha \iota_2, \pi_2 \beta \iota_2)).$$

It is clear that this the denotations of a proof. We observe that

$$\begin{aligned} (\Phi \Psi(p))(\alpha, \beta) &= \pi_1 \Psi(p)(\alpha \oplus 1_{\mathcal{W}}, \beta \oplus 1_{\mathcal{W}}) \\ &= \pi_1 \left( p(\pi_1(\alpha \oplus 1_{\mathcal{W}}) \iota_1, \pi_1(\beta \oplus 1_{\mathcal{W}}) \iota_1), p(\pi_2(\alpha \oplus 1_{\mathcal{W}}) \iota_2, \pi_2(\beta \oplus 1_{\mathcal{W}}) \iota_2) \right) \\ &= \pi_1(p(\alpha, \beta), p(1_{\mathcal{W}}, 1_{\mathcal{W}})) \\ &= p(\alpha, \beta) \end{aligned}$$

which proves (i). It is clear that (ii) holds. □



**Remark 5.6.** The purpose of linear logic in this context is to give a coherent and general way of writing down such polynomials, which realise a given computational idea when inserted into the evolution equation of an RNN.

In this case the term in (4.13) is

$$(5.2) \quad \left( \sum_i \lambda_i R^i \right)^a \left( \sum_j \mu_j R^j \right)^b (r_1)$$

which can be expanded as

$$(5.3) \quad \sum_{u_0+\dots+u_L=a} \sum_{v_0+\dots+v_L=b} \frac{a!}{u_0! \dots u_L!} \frac{b!}{v_0! \dots v_L!} \lambda_0^{u_0} \dots \lambda_L^{u_L} \mu_0^{v_0} \dots \mu_L^{v_L} R^{\sum_t t(u_t+v_t)}(r_1).$$

This may be written as  $\sum_p f_p(\lambda, \mu) R^p(r_1)$  for some polynomials  $f_p$  in variables  $\{\lambda_i, \mu_j\}_{i,j}$ .

## 5.1 Decaying pattern NTM

We take the same alphabet  $\mathbf{a}$  as in the multiple pattern NTM, and the same coefficient space  $\mathcal{V}_1$  for the first memory ring. There are two additional memory rings, with

$$\mathcal{V}_2 \subseteq \llbracket \mathbf{bint}_{W \& W} \rrbracket, \quad \mathcal{V}_3 \subseteq \llbracket \mathbf{bint}_{W \& W} \multimap \mathbf{bint}_W \& \mathbf{bint}_W \rrbracket.$$

An example of a proof in  $\mathcal{V}_3$  is the algorithm tail which takes a binary sequence and returns a pair by separating the tail, for example

$$\llbracket \text{tail} \rrbracket (1101_{W \& W}) = (110_W, 1_W).$$

The idea is that if the third memory ring is sharply focused at tail then we update the second memory ring by replacing a vector  $\llbracket 1101 \rrbracket$  at a vertex of the ring by its truncation  $\llbracket 110 \rrbracket$ , and we moreover use the tail  $\llbracket 1 \rrbracket$  to update the first memory ring by reading it as a power of the rotation operator. In this way the contents of the second memory ring *decay* over time. To realise this we need the following auxiliary maps:

**Definition 5.7.** Given a type  $W$  with  $\mathcal{W} = \llbracket W \rrbracket$  let less be the proof of

$$\mathbf{bint}_{W \& W} \vdash \mathbf{bint}_W$$

with the property that for  $\alpha, \beta \in \text{End}_{\mathbb{R}}(\llbracket W \rrbracket)$  and  $q \in \llbracket \mathbf{bint}_{W \& W} \rrbracket$

$$\llbracket \text{less} \rrbracket (q)(\alpha, \beta) = \pi_1 q(\alpha \oplus 1_{\mathcal{W}}, \beta \oplus 1_{\mathcal{W}})$$

where  $\pi_i$  denote the projections from  $\mathcal{W} \oplus \mathcal{W}$ . Then

$$\Phi \llbracket S_{W \& W} \rrbracket = \llbracket S_W \rrbracket$$

for all  $S \in \{0, 1\}^*$ .

There does not appear to be a proof more in linear logic which gives a section of less in the sense that its denotation maps  $\llbracket \underline{S}_W \rrbracket$  to  $\llbracket \underline{S}_{W \& W} \rrbracket$  (see however Lemma 5.5). But we can just fix one ad hoc (this is simulating second order linear logic in an ugly way).

**Definition 5.8.** We fix a linear map  $\psi : \llbracket \mathbf{bint}_W \rrbracket \longrightarrow \llbracket \mathbf{bint}_{W \& W} \rrbracket$  with the property that

$$\psi \llbracket \underline{S}_W \rrbracket = \llbracket \underline{S}_{W \& W} \rrbracket$$

for a finite collection of sequences  $S$  that are relevant to the model.

The master algorithm at a given time step applies  $M_3(r_3) \in \mathcal{V}_3$  to  $M_2(r_2) \in \mathcal{V}_2$  to get

$$\llbracket \underline{\text{feed}} \rrbracket(M_2(r_2), M_3(r_3)) \in \llbracket \mathbf{bint}_W \rrbracket \oplus \llbracket \mathbf{bint}_W \rrbracket.$$

The second component is then applied to the contents of the first memory ring, as in

$$r_1^{(t+1)} = \left( \pi_2 \llbracket \underline{\text{feed}} \rrbracket(M_2(r_2), M_3(r_3)) \right) (1_{\mathcal{W}}, R)(r_1^{(t)}).$$

Simultaneously, we update the memory state of the second ring according to the following rule: it becomes the linear map  $\mathcal{W} \longrightarrow \mathcal{V}_2$  defined by

$$w \longmapsto \psi \pi_1 \llbracket \underline{\text{feed}} \rrbracket(M_2(w), M_3(r_3)).$$

This kind of structure could in principle be used by the controller to predict that an event will happen a fixed amount of time in the future, by pushing “countdowns”  $\llbracket \underline{10000} \rrbracket$  onto the memory ring, and have the controller react to a symbol coming off the first memory ring after it is shifted by the final  $\llbracket \underline{1} \rrbracket$ .

**Remark 5.9.** Actually we seem to have been confused about the kinds of polynomials we get. It seems that unless feed uses non-linear inputs we never get powers of variables. So we should find an example where this is the case.

## References

- [1] D. Bahdanau, K. Cho and Y. Bengio, *Neural machine translation by jointly learning to align and translate*, arXiv preprint arXiv:1409.0473.
- [2] K. Cho, A. Courville and Y. Bengio, *Describing multimedia content using attention-based encoder-decoder networks*, IEEE Transactions on Multimedia, 17(11), 1875-1886, (2015).
- [3] A. Neelakantan, Q. V. Le and I. Sutskever, *Neural Programmer: Inducing Latent Programs with Gradient Descent*, (2015).
- [4] S. Reed and N.de Freitas, *Neural Programmer-Interpreters*, (2016).

- [5] Graves, Alex, et al. "Hybrid computing using a neural network with dynamic external memory." *Nature* 538.7626 (2016): 471-476.
- [6] G. Frege, *Begriffsschrift, a formula language, modeled upon that of arithmetic, for pure thought*, (1879). An english translation appears in *From Frege to Gödel. A source book in mathematical logic, 1879–1931*, Edited by J. van Heijenoort, Harvard University Press, 1967.t
- [7] R. Blute, T. Ehrhard and C. Tasson, *A convenient differential category*, arXiv preprint [[arXiv:1006.3140](https://arxiv.org/abs/1006.3140)], 2010.
- [8] M. H. Sørensen and P. Urzyczyn, *Lectures on the Curry-Howard isomorphism* (Vol. 149), Elsevier, (2006).
- [9] Y. LeCun, Y. Bengio and G. Hinton, *Deep learning*, *Nature*, 521(7553), pp.436–444 (2015).
- [10] J. Elman, *Finding structure in time*, *Cognitive science*, 14(2):179-211, 1990.
- [11] A. Graves, G. Wayne and I. Danihelka, *Neural turing machines*, arXiv preprint [arXiv:1410.5401](https://arxiv.org/abs/1410.5401) (2014).
- [12] A. Joulin and T. Mikolov, *Inferring algorithmic patterns with stack-augmented recurrent nets*, *Advances in Neural Information Processing Systems*, 2015.
- [13] A. Graves, *Hybrid computing using a neural network with dynamic external memory*, *Nature* 538.7626 (2016): 471–476.
- [14] 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.
- [15] O. Irsoy and C. Cardie, *Modeling compositionality with multiplicative recurrent neural networks*, arXiv preprint [arXiv:1412.6577](https://arxiv.org/abs/1412.6577) (2014).
- [16] 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.
- [17] I. Sutskever, O. Vinyals and Q. V. Le, *Sequence to sequence learning with neural networks*, *Advances in neural information processing systems*, 2014.
- [18] 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.

- [19] 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.
- [20] 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.
- [21] J. B. Pollack, *The induction of dynamical recognizers*, Machine Learning, 7(2-3), 227-252 (1991).
- [22] D. Murfet, *Logic and linear algebra: an introduction*, preprint (2014) [arXiv:1407.2650].
- [23] D. Murfet, *On Sweedler’s cofree cocommutative coalgebra*, J. Pure and Applied Algebra **219** (2015) 5289–5304.
- [24] J.-Y. Girard, *Linear Logic*, Theoretical Computer Science **50** (1987), 1–102.
- [25] A. Joulin and T. Mikolov, *Inferring algorithmic patterns with stack-augmented recurrent nets*, Advances in Neural Information Processing Systems, 2015.
- [26] E. Grefenstette, et al, *Learning to transduce with unbounded memory*, Advances in Neural Information Processing Systems, 2015.
- [27] J. Weston, C. Sumit and B. Antoine, *Memory networks*, preprint (2014) [arXiv:1410.3916].
- [28] W. Zaremba, et al., *Learning Simple Algorithms from Examples*, preprint (2015) [arXiv:1511.07275].
- [29] P-A. Melliès, *Categorical semantics of linear logic*, in : Interactive models of computation and program behaviour, Panoramas et Synthèses 27, Société Mathématique de France, 2009.
- [30] A. A. Alemi, F. Chollet, G. Irving, C. Szegedy and J. Urban, *DeepMath-Deep Sequence Models for Premise Selection*, arXiv preprint arXiv:1606.04442.
- [31] M. Balog, A. L. Gaunt, M. Brockschmidt, S. Nowozin and D. Tarlow, *Deep-Coder: Learning to Write Programs*, arXiv preprint arXiv:1611.01989.
- [32] T. Rocktäschel and S. Riedel, *Learning Knowledge Base Inference with Neural Theorem Provers*, In NAACL Workshop on Automated Knowledge Base Construction (AKBC) 2016.
- [33] T. Ehrhard, *An introduction to Differential Linear Logic: proof-nets, models and antiderivatives*, [arXiv:1606.01642] (2016).

- [34] J. Clift and D. Mufet, *Cofree coalgebras and differential linear logic*, preprint.
- [35] Abadi, Martn, et al. *TensorFlow: A system for large-scale machine learning* arXiv preprint arXiv:1605.08695 (2016).
- [36] Abadi, Martn, et al. *Tensorflow: Large-scale machine learning on heterogeneous distributed systems*, arXiv preprint arXiv:1603.04467 (2016).
- [37] A. Griewank and A. Walther, *Evaluating derivatives: principles and techniques of algorithmic differentiation*, Siam (2008).
- [38] T. Ehrhard and L. Regnier, *The differential  $\lambda$ -calculus*, Theoretical Computer Science 309, pp. 1–41, (2003).
- [39] 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. [\[URL\]](#).
- [40] M. Minsky, *Logical versus analogical or symbolic versus connectionist or neat versus scruffy*, AI magazine, 12(2), 34 (1991).
- [41] J.-Y. Girard, Y. Lafont and P. Taylor, *Proofs and Types*, Cambridge Tracts in Theoretical Computer Science 7, Cambridge University Press, 1989.
- [42] N. Benton, G. Bierman, V. de Paiva and M. Hyland, *Term assignment for intuitionistic linear logic*, Technical report 262, Computer Laboratory, University of Cambridge, 1992.
- [43] I. Goodfellow, Y. Bengio and A. Courville, *Deep learning*, MIT Press.
- [44] Loos, S., Irving, G., Szegedy, C., Kaliszyk, C. (2017). Deep Network Guided Proof Search. arXiv preprint arXiv:1701.06972.
- [45] T. Mikolov, A. Joulin, S. Chopra, M. Mathieu and M. A. Ranzato, *Learning longer memory in recurrent neural networks*, arXiv preprint arXiv:1412.7753 (2014).