

Synergistic Integration of Linear Logic and Recurrent Neural Networks for Enhanced Sequence Modeling

Abstract—This paper explores the fusion of linear logic and recurrent neural networks (RNNs) to create a novel sequence modeling architecture. The approach leverages the differentiable semantics of linear logic to inform the operation of RNNs, resulting in a system capable of sophisticated manipulation of hidden states. The paper presents the theoretical underpinnings of this Linear Logic Recurrent Neural Network (LLRNN) and discusses its potential to bridge the gap between symbolic reasoning and connectionist models. While preliminary, this work offers a new perspective on neural computation and highlights the utility of linear logic in designing advanced neural network structures.

I. INTRODUCTION

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 *update 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 update operator.

Our primary motivation is to provide a new perspective on the problem of reconciling symbolic logic and neural networks, based on differentiable semantics. We view the current work as a preliminary step, which describes the theory and demonstrates its content by showing how ideas from linear logic manifest themselves in the context of “neural reasoning” about a few select example tasks. We have not tested the LLRNN on large scale data, or attempted to systematically probe its limits; it is likely new ideas beyond those presented here will be needed to achieve anything of substantial practical value with the LLRNN. We hope the code published alongside this paper will be useful for developing the intuitions that will guide these future developments.

II. 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.

A. Recurrent Neural Network

We begin with a review of the ordinary Recurrent Neural Network (RNN). A good textbook 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 use the *rectified linear* function

$$: \mathbb{R}^k \longrightarrow \mathbb{R}^k \quad (\mathbf{x})_i = \frac{1}{2}(x_i + |x_i|), \quad (\text{II.1})$$

the *softmax* function

$$: \mathbb{R}^k \longrightarrow \mathbb{R}^k, \quad (\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^k e^{x_j}}, \quad (\text{II.2})$$

and the *sigmoid* function

$$: \mathbb{R}^k \longrightarrow \mathbb{R}^k, \quad (\mathbf{x})_i = \frac{1}{1 + e^{-x_i}}. \quad (\text{II.3})$$

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 II.1. A standard Elman-style RNN [10] is defined by weight matrices H, U, B and the evolution equation

$$h^{(t+1)} = (Hh^{(t)} + Ux^{(t+1)} + B) \quad (\text{II.4})$$

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

$$y^{(t)} = (W_y h^{(t)} + B_y). \quad (\text{II.5})$$

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 (II.4) 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 τ obtain the output sequence $o^{(t)} = (W_o y^{(t)} + B_o)$. 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} .

B. 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

$$\sum_j \sum_k V_i^{jk} h_j^{(t)} x_k^{(t+1)} \quad (\text{II.6})$$

where V is interpreted as a tensor in

$$\mathcal{I}^* \otimes \mathcal{H}^* \otimes \mathcal{H} \cong \text{Hom}_{\mathbb{R}}(\mathcal{I}, \text{End}_{\mathbb{R}}(\mathcal{H})). \quad (\text{II.7})$$

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 (II.6).

Definition II.2. The second-order RNN [20], [21] is defined by weights H, U, B, V and

$$h^{(t+1)} = (V(x^{(t+1)})(h^{(t)}) + Hh^{(t)} + Ux^{(t+1)} + B), \quad (\text{II.8})$$

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

Remark II.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

$$V : \mathcal{I} \rightarrow \text{End}_{\mathbb{R}}(\mathcal{H}), \quad I : \mathcal{K} \rightarrow \mathcal{H}, \quad J : \mathcal{H} \rightarrow \mathcal{K}. \quad (\text{II.9})$$

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 (II.8) is changed to

$$I(V(x^{(t+1)})(Jh^{(t)})). \quad (\text{II.10})$$

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 II.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].

C. Neural Turing Machine

The Neural Turing Machine (NTM) [11] is an RNN augmented with an additional memory. This memory is read from and written to using ‘‘attention’’ vectors that are manipulated at each time step based on the hidden state. In this section we first give a presentation of the NTM following [11, §3], and then in Section II-C1 we give a basis-free presentation.

Let N be the number of memory locations and M the size of the vector stored at each location. Then the internal state space of the NTM is decomposed as

$$\mathcal{H} = \mathcal{H}_0 \oplus \mathbb{R}^N \oplus \mathbb{R}^N \oplus \text{Mat}_{M \times N}(\mathbb{R}) \quad (\text{II.11})$$

where \mathcal{H}_0 is the hidden state space of the underlying RNN and $\text{Mat}_{M \times N}(\mathbb{R})$ is the space of $M \times N$ matrices. These matrices are identified with states of the memory, with the i th column being the vector stored at location i (we index columns in the range $0, \dots, N-1$). Note that [11] uses $N \times M$ matrices and thus identifies rows with the vectors stored at memory locations. The state at time t may therefore be written

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

The components are interpreted as the *controller internal state* $h_0^{(t)}$, *read address weighting* $r^{(t)} \in \mathbb{R}^N$, *write address weighting* $w^{(t)} \in \mathbb{R}^N$ and *memory state* $M^{(t)} \in \text{Mat}_{M \times N}(\mathbb{R})$. We index memory locations from 0 to $N-1$, so that for example

$$r^{(t)} = (r_0^{(t)}, \dots, r_{N-1}^{(t)}).$$

At time t the RNN generates from its hidden state $h_0^{(t)}$ a distribution $s^{(t+1)}$ over rotations of the write address, a distribution $q^{(t+1)}$ over rotations of the read address, an *erase*

vector $e^{(t+1)}$ and an *add vector* $a^{(t+1)}$ [11, §3.2] via the formulas

$$\begin{aligned} s^{(t+1)} &= (W_s h_0^{(t)} + B_s) \in \mathbb{R}^N, \\ q^{(t+1)} &= (W_q h_0^{(t)} + B_q) \in \mathbb{R}^N, \\ e^{(t+1)} &= (W_e h_0^{(t)} + B_e) \in \mathbb{R}^M, \\ a^{(t+1)} &= (W_u h_0^{(t)} + B_u) \in \mathbb{R}^M. \end{aligned}$$

The update equation for the addresses [11, Eq. (8)] and memory [11, Eq. (3),(4)] are

$$w_i^{(t+1)} = \sum_{j=0}^{N-1} s_{i-j}^{(t+1)} w_j^{(t)}, \quad r_i^{(t+1)} = \sum_{j=0}^{N-1} q_{i-j}^{(t+1)} r_j^{(t)} \quad (\text{II.12})$$

$$M_{\bullet,i}^{(t+1)} = M_{\bullet,i}^{(t)} - w_i^{(t+1)} (M_{\bullet,i}^{(t)} \odot e^{(t+1)}) + w_i^{(t+1)} a^{(t+1)} \quad (\text{II.13})$$

where $M_{\bullet,i}$ stands for the i th column of the matrix $M \in \text{Mat}_{M \times N}(\mathbb{R})$ and \odot is the so-called *Hadamard product* which is defined by

$$\begin{aligned} \odot : \mathbb{R}^M \otimes \mathbb{R}^M &\longrightarrow \mathbb{R}^M, \\ \mathbf{x} \odot \mathbf{y} &= (x_i y_i)_{i=1}^M. \end{aligned}$$

Finally, the evolution equation is

$$h_0^{(t+1)} = \left(M^{(t)} r^{(t)} + H_0 h_0^{(t)} + U_0 x^{(t+1)} + B_0 \right) \quad (\text{II.14})$$

where $H_0 : \mathcal{H}_0 \longrightarrow \mathcal{H}_0$, $U_0 : \mathcal{I} \longrightarrow \mathcal{H}_0$ and $B_0 \in \mathcal{H}_0$ are weight matrices.

1) *Abstract reformulation:* In Section III-B we will reformulate the NTM in the context of linear logic and use that as a foundation on which to build some more interesting models. As a starting point we will need a description of the NTM which does not depend on a choice of basis, and which makes clear the underlying algebraic principles. The model uses four vector spaces:

- the input space \mathcal{I} and hidden state space \mathcal{H}_0 of the underlying RNN, and
- the *address space* \mathcal{W} and *coefficient space* \mathcal{V} of the memory system.

The *space of memory states* is the space of linear maps

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

We think of a basis e_1, \dots, e_N for \mathcal{W} as “locations” at which may be stored a vector from the coefficient space \mathcal{V} . A memory state $M \in \mathcal{S}$ is determined by the vectors $M(e_i) \in \mathcal{V}$, which we view as being stored at the locations e_i . While we will use the basis $\{e_i\}_i$ in the following to illustrate the definitions, it is important to note that the model is defined in a way which avoids explicit reference to any basis. Taking $\mathcal{W} = \mathbb{R}^N$ and $\mathcal{V} = \mathbb{R}^M$ puts us back in the situation of the previous section, with $\mathcal{S} \cong \text{Mat}_{M \times N}(\mathbb{R})$.

Definition II.5. A *read address* is a vector $r \in \mathcal{W}$ and a *write address* is a vector $w \in \mathcal{W}^*$.

Remark II.6. We may evaluate a memory state M at a read address r to obtain a vector $M(r) \in \mathcal{V}$, and if the read address is *focused* at location i , by which we mean that $r = e_i$, this produces the vector $M(e_i)$ stored at that location. In general a read address specifies a linear combination of the stored vectors. To explain the thinking behind the use of dual vectors as write addresses, consider the canonical isomorphism

$$\Psi : \mathcal{W}^* \otimes \mathcal{V} \longrightarrow \text{Hom}_{\mathbb{R}}(\mathcal{W}, \mathcal{V}) \quad (\text{II.15})$$

under which a tensor $\alpha \otimes v$ with $\alpha \in \mathcal{W}^*$, $v \in \mathcal{V}$ corresponds to the linear map

$$\Psi(\alpha \otimes v)(w) = \alpha(w)v.$$

Given a state M of the memory, a location e_i and a vector v to write to this location, the new memory state $M' = M + \Psi(e_i^* \otimes v)$ has the property that

$$M'(e_j) = M(e_j) + \Psi(e_i^* \otimes v)(e_j) = M(e_j) + \delta_{i=j}v.$$

Note that in this construction it is the dual vector $w = e_i^*$ which specifies *where to write*, and this explains why we refer to vectors in \mathcal{W}^* as write addresses.

The state space of the model is defined to be

$$\mathcal{H} = \mathcal{H}_0 \oplus \mathcal{W} \oplus \mathcal{W}^* \oplus \mathcal{S}.$$

A closer examination of the model in [11] leads us to posit the following:

- The RNN controller manipulates read and write addresses by invertible linear transformations, which suggests that we take \mathcal{W} to be a *representation* of a finite group G . At each time step the controller will produce distributions over G which are used to act on the read and write addresses by the action maps

$$G \longrightarrow \text{End}_{\mathbb{R}}(\mathcal{W}), \quad G \longrightarrow \text{End}_{\mathbb{R}}(\mathcal{W}^*).$$

The second of these maps arises since \mathcal{W}^* is also naturally a representation of G .

- In the memory update equation (II.13) we need to *copy* memory addresses (the subscript i appears twice in the second summand on the right hand side) and *multiply* vectors in \mathcal{V} (using \odot) so we take \mathcal{W} to be a coalgebra and \mathcal{V} to be an algebra.

One way to realise these constraints is using a finite group $G = \{g_0, g_1, \dots, g_{N-1}\}$ with identity $g_0 = e$ acting by linear transformations on the group algebra

$$\mathbb{R}G = \bigoplus_{g \in G} \mathbb{R}g = \mathbb{R}e \oplus \mathbb{R}g_1 \oplus \dots \oplus \mathbb{R}g_{N-1}$$

which is a vector space with basis G . This is a coalgebra, with comultiplication

$$\Delta : \mathbb{R}G \longrightarrow \mathbb{R}G \otimes \mathbb{R}G,$$

$$\Delta(g) = g \otimes g.$$

The group G acts on the space $\mathscr{W} = \mathbb{R}G$ by the rule

$$G \times \mathscr{W} \longrightarrow \mathscr{W},$$

$$(g, \sum_{h \in G} \lambda_h h) \longmapsto \sum_{h \in G} \lambda_h gh.$$

This action of G extends to a linear map $A : \mathbb{R}G \longrightarrow \text{End}_{\mathbb{R}}(\mathscr{W})$.

Finally we take some commutative algebra structure on \mathscr{V} . The use of an algebra structure is implicit in [11] and the deep learning literature in general, and enters via the *Hadamard product* \odot described above. Since this is ultimately our example of interest, we denote the product structure on \mathscr{V} by \odot , although it need not be the Hadamard product.

Since \mathscr{W} is a coalgebra and \mathscr{V} is an algebra we have:

Definition II.7 (Convolution product). The vector space $\text{Hom}_{\mathbb{R}}(\mathscr{W}, \mathscr{V})$ is an associative algebra with multiplication of M_1, M_2 given by the composite

$$\mathscr{W} \xrightarrow{\Delta} \mathscr{W} \otimes \mathscr{W} \xrightarrow{M_1 \otimes M_2} \mathscr{V} \otimes \mathscr{V} \xrightarrow{\odot} \mathscr{V}.$$

We write $M_1 \star M_2$ for this *convolution product*, so that

$$M_1 \star M_2 = m \circ (M_1 \otimes M_2) \circ \Delta.$$

Explicitly,

$$(M_1 \star M_2)(g) = M_1(g) \odot M_2(g).$$

Given these definitions we can state the abstract form of the evolution equation:

Definition II.8 (Evolution equations). Suppose the current state of the NTM is

$$h = (h_0, r, w, M) \in \mathscr{H} = \mathscr{H}_0 \oplus \mathscr{W} \oplus \mathscr{W}^* \oplus \text{Hom}_{\mathbb{R}}(\mathscr{W}, \mathscr{V}).$$

With the same formulas as above we generate from the hidden state h_0 the vectors

$$s, q \in \mathscr{W}, \quad e, a \in \mathscr{V}.$$

The new state $h' = (h'_0, r', w', M')$ is defined by the formulas

$$w' = s \cdot w \tag{II.16}$$

$$r' = q \cdot r \tag{II.17}$$

$$M' = M - M \star \Psi(w' \otimes e) + \Psi(w' \otimes a) \tag{II.18}$$

$$h'_0 = \left(M(r) + H_0 h_0 + U_0 x + B_0 \right). \tag{II.19}$$

where $s \cdot w$ is shorthand for $A(s)(w)$, the action of $s \in \mathbb{R}G$ on \mathscr{W} , Ψ is from (II.15) and \star is the convolution product.

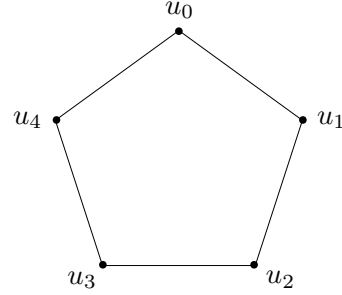
Remark II.9. However, using w' in the M' equation messes up our division of linear and nonlinear operations.

Remark II.10. For concreteness, take G to be the group of N th roots of unity

$$G = \{e^{\frac{2\pi i}{N}} \mid 0 \leq i \leq N-1\} \cong \mathbb{Z}/N\mathbb{Z}$$

and we write e_i for $e^{\frac{2\pi i}{N}}$ viewed as a basis vector of \mathscr{W} . Since we are thinking of these basis vectors as memory locations,

we may view a memory state M as an arrangement of the vectors $u_j = M(e_j)$ on the vertices of the regular N -gon in the complex plane



Sometimes we refer to this structure as a *memory ring*.

We assume for simplicity that the memory coefficient space is the same as the hidden state space, that is, $\mathscr{V} = \mathscr{H}$.¹

D. 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.

The universality of linear logic means that one can expect any form of “abstract reasoning” which manifests itself mathematically in the context of linear algebra (for example, reasoning in neural networks) to be, in a certain limit, coincident with the kind of reasoning axiomatised by linear logic. Put differently, linear logic identifies those structures on the category of vector spaces which we are “forced” to care about in this context, and nothing more.

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.

1) *Types*: The *types* of intuitionistic linear logic are built from atomic types ϕ_1, ϕ_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 ϕ_i and the other vertices are labelled by $\otimes, \multimap, !$. For example the type

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

¹The 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].

A the denotation is a vector $\llbracket \pi \rrbracket \in \llbracket A \rrbracket$. In particular, if ρ 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 $\llbracket A \rrbracket \rightarrow \llbracket B \rrbracket$. Similarly, if γ is of type $!A \multimap B$ then its denotation is a linear map

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

In Remark II.11 we explained how to think of both ρ, γ as programs taking inputs of type A , but where ρ 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 (II.24) to define the (non-linear) composite

$$\llbracket A \rrbracket \longrightarrow \llbracket !A \rrbracket \xrightarrow{\llbracket \gamma \rrbracket} \llbracket B \rrbracket \quad (\text{II.25})$$

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 II.13. The denotation of a term π of type \mathbf{int}_A is by Example II.12 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 $!\text{End}_{\mathbb{R}}(V) \rightarrow \text{End}_{\mathbb{R}}(V)$. 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 $\alpha \in \text{End}_{\mathbb{R}}(V)$,

$$\llbracket \underline{n} \rrbracket(\emptyset)_{\alpha} = \alpha^n.$$

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

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

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

Example II.14. The denotation of a proof π of type \mathbf{bint}_A gives a function

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

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 α, β to the product described by S , reading α for 0 and β for 1, and reading in reverse order. For example,

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

III. 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

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

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

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

and the denotation of a term $\underline{\text{master}^z}$ in linear logic, via the formula

$$Z_z = \llbracket \underline{\text{master}^z} \rrbracket_{nl} \left(p_1^{(t+1)}, \dots, p_r^{(t+1)}, b_1^{(t+1)}, \dots, b_s^{(t+1)}, c^{(t+1)} \right). \quad (\text{III.2})$$

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

$$\mathcal{P}_1 \times \dots \times \mathcal{P}_r \times \mathcal{B}_1 \times \dots \times \mathcal{B}_s \times \mathcal{C} \rightarrow \text{End}_{\mathbb{R}}(\mathcal{H}). \quad (\text{III.3})$$

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

- A sequence of *command vectors* $p_i^{(t+1)}$ generated via softmax from $h^{(t)}$,

- a sequence of *data vectors* $b_i^{(t+1)}$ generated via from $h^{(t)}$,
- an *input vector* $c^{(t+1)}$ generated via 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 III.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

$$!P_1, \dots, !P_r, !B_1, \dots, !B_s, !C \vdash A \multimap A. \quad (\text{III.4})$$

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

$$\mathcal{P}_i \subseteq \llbracket P_i \rrbracket, \quad \mathcal{B}_j \subseteq \llbracket B_j \rrbracket, \quad \mathcal{C} \subseteq \llbracket C \rrbracket \quad (\text{III.5})$$

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 III.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 III.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 Γ , that is, $!\Gamma = !P_1, \dots, !P_r$. With this notation, $\underline{\text{master}}$ is a proof of $!\Gamma, !\Delta, !V \vdash A \multimap A$.

A. 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 III.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

$$Z_{\text{in}} = W^c(x^{(t+1)}). \quad (\text{III.6})$$

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 III.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

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

and therefore

$$Z_{\text{in}} = \sum_{i=0}^L p_i^{(t+1)} (W^c(x^{(t+1)}))^i. \quad (\text{III.7})$$

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 $p_i \in \mathbb{R}$,

$$\llbracket \underline{\text{master}}^{\text{in}} \rrbracket_{nl} \left(\sum_i \lambda_i \llbracket n_i \rrbracket, \alpha \right) = \sum_i \lambda_i \alpha^{n_i}. \quad (\text{III.8})$$

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 (III.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 $\llbracket \text{master} \rrbracket_{nl}$ rather than the underlying proof. Nonetheless, the construction of 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.

B. NTM as an LLRNN

We now explain how to present the NTM of Section II-C 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 Γ and data types Δ 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

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

we have, for $h = (h_0, r, w, M) \in \mathcal{H}$ decomposed according to (II.11),

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

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 $\llbracket \text{master}^z \rrbracket_{nl}$ 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 (III.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 (III.1) of the LLRNN reads

$$\begin{aligned} (h_0^{(t+1)}, r^{(t+1)}, w^{(t+1)}, M^{(t+1)}) &= \left((M^{(t)}(r^{(t)}) + H_0 h_0^{(t)} + U_0 x^{(t+1)} + \right. \\ &\quad \sum_{i=0}^{N-1} (p_1^{(t+1)})_i R^i(r^{(t)}), \\ &\quad \sum_{i=0}^{N-1} (p_2^{(t+1)})_i (R^*)^i(w^{(t)}), \\ &\quad \left. (1_{\mathcal{V}} - e^{(t+1)}) \circ M^{(t)} + w^{(t)} \otimes a \right) \end{aligned} \quad (\text{III.9})$$

which agrees with the equations (??) – (II.14).

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

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

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

IV. GENERALISATIONS OF THE NTM

Taking the linear logic point of view on the NTM from Section III-B 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.

A. Pattern NTM

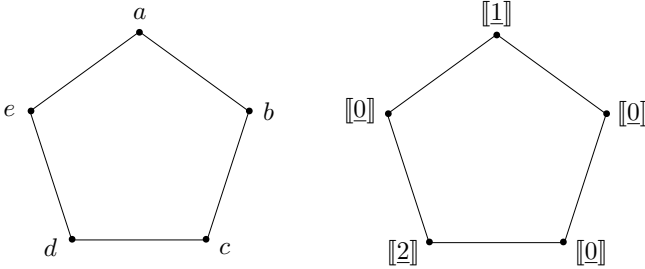
At each time step the NTM predicts distributions $s^{(t)}, q^{(t)}$ over the possible rotations of the write and read addresses. A pattern of memory accesses is a *sequence* of such rotations, and since the ‘‘angle’’ of rotation is represented in the LLRNN as the denotation of a linear logic proof (with $\llbracket n \rrbracket$ representing a clockwise angle of $\frac{2\pi n}{5}$) the pattern may be stored as a

sequence of vectors in $\llbracket \mathbf{int}_W \rrbracket$. 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 $\mathcal{V}_1 = \bigoplus_{i \in \mathbf{a}} \mathbb{R}i$ be the free vector space on A . We also take a finite-dimensional subspace $\mathcal{V}_2 \subseteq \llbracket \mathbf{int}_W \rrbracket$. 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

$$\llbracket 1 \rrbracket, \llbracket 0 \rrbracket, \llbracket 0 \rrbracket, \llbracket 0 \rrbracket, \llbracket 2 \rrbracket, \llbracket 0 \rrbracket, \dots \quad (\text{IV.1})$$

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

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

which has the property that

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

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

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

The command and data types are

$$\Gamma = \mathbf{int}_{W^\vee}, \mathbf{int}_W, \mathbf{int}_{W^\vee} \quad (\text{IV.2})$$

$$\Delta = W \multimap W, W^\vee \multimap W^\vee, W \multimap W, W^\vee \multimap W^\vee, V \multimap V, V, \quad (\text{IV.3})$$

$$\mathbf{int}_W \multimap \mathbf{int}_W, \mathbf{int}_W. \quad (\text{IV.4})$$

The first type in Γ 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 Δ 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 master^z is a proof of $! \Gamma, ! \Delta \vdash A \multimap A$.

Remark IV.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 (??), and equations (for $1 \leq i \leq 2$ in all cases)

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

$$q_2^{(t+1)} = (W_q h^{(t)} + B_q) \in \mathcal{W}^* \cong \mathbb{R}^5, \quad (\text{IV.6})$$

$$e_i^{(t+1)} = (W_{i,e} h^{(t)} + B_{i,e}) \in \text{End}_{\mathbb{R}}(\mathcal{V}_i), \quad (\text{IV.7})$$

$$a_i^{(t+1)} = (W_{i,a} h^{(t)} + B_{i,a}) \in \mathcal{V}_i, \quad (\text{IV.8})$$

$$w_i^{(t+1)} = \sum_{i=0}^{N-1} s^{(t+1)}(\bar{i}) \cdot (R^*)^i(w_i^{(t)}), \quad (\text{IV.9})$$

$$r_1^{(t+1)} = \sum_{j=0}^L M_2^{(t)}(r_2^{(t)})_j \cdot R^j(r_1^{(t)}), \quad (\text{IV.10})$$

$$r_2^{(t+1)} = \sum_{i=0}^{N-1} q_2^{(t+1)}(\bar{i}) \cdot R^i(r_2^{(t)}), \quad (\text{IV.11})$$

$$M_i^{(t+1)} = (1_{\mathcal{V}_i} - e_i^{(t+1)}) \circ M_i^{(t)} + w_i^{(t)} \otimes a_i^{(t+1)}. \quad (\text{IV.12})$$

Finally, the evolution equation is

$$h_0^{(t+1)} = \left(Q(M_1^{(t)}(r_1^{(t)})) + H_0 h_0^{(t)} + U_0 x^{(t+1)} + B_0 \right). \quad (\text{IV.13})$$

The main difference between the ordinary NTM and the pattern NTM is in (IV.10) 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.

B. 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 \mathbf{int}_W and a ring with coefficients in \mathbf{bint}_W . 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 $\mathcal{S} = \mathcal{V}_1 = \bigoplus_{i \in \mathbf{a}} \mathbb{R}i$ 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

$$\circ abcde Sabcde Sabcde \bullet \mapsto abcde aabccdde abcde, \quad (\text{IV.14})$$

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 II.14) 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

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

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 $\llbracket \text{master}^{\text{out}} \rrbracket_{nl}(g)(h)$ now looks like

$$\llbracket \text{feed} \rrbracket (M_2(r_2), M_3(r_3), M_4(r_4))_{nl}(R)(r_1). \quad (\text{IV.15})$$

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

$$M_2(r_2) = \sum_{i=0}^L \lambda_i \llbracket i \rrbracket, \quad M_3(r_3) = \sum_{j=0}^L \mu_j \llbracket j \rrbracket, \quad M_4(r_4) = \sum_{k=0}^m \kappa_k \llbracket S_k \rrbracket \quad (\text{IV.16})$$

In this case the term in (IV.15) is

$$\sum_{i,j,k} \lambda_i \mu_j \kappa_k R^{a_k i + b_k j}(r_1) \quad (\text{IV.17})$$

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

$$r_1[u] \longleftarrow \sum_{i,j,k} \lambda_i \mu_j \kappa_k r_1[u - a_k i - b_k j]. \quad (\text{IV.18})$$

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 (IV.14) 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.

C. Polynomial step NTM

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

$$\mathfrak{a} = \{a, b, \dots, z\} \cup \{\circ, \bullet, T\} \quad (\text{IV.19})$$

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 IV-A, and

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

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

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

with linear lambda term

$$\text{newfeed} = (\lambda p q x. ((q(\text{promote } p))x)). \quad (\text{IV.20})$$

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

$$\llbracket \text{newfeed} \rrbracket (M_2(r_2), M_3(r_3))_{nl}(R)(r_1). \quad (\text{IV.21})$$

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 (IV.21) is

$$\left\{ \sum_j \mu_j a_j \left(\sum_i \lambda_i R^i \right) \right\} (r_1). \quad (\text{IV.22})$$

Note that $\text{End}_{\mathbb{R}}(\mathcal{W})$ 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 (IV.22) 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 .

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 III.6. We keep the same alphabet, but take

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

Take \mathcal{V}_2 as in Section IV-A, and

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

Where id is the identity function, and 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 \mathcal{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

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

corresponds to the group elements being squared

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

V. EXPERIMENTAL RESULTS, CONCLUSION, AND FUTURE WORK

A. Experimental Results

To explore the capabilities of the proposed Linear Logic Recurrent Neural Network (LLRNN), we conducted a series of proof-of-concept simulations on synthetic sequence modeling tasks. These included symbolic reasoning chains, sequence transductions involving arithmetic pattern compositions, and controlled memory manipulations derived from variants of the Neural Turing Machine (NTM).

Our experiments show that the LLRNN can encode higher-order transformations, such as raising input-dependent linear operators to predicted integer powers, through logical constructs rather than manually engineered operations. Furthermore, we demonstrated the model's ability to generalize structural behaviors like rotations and reflections over symbolic rings using denotations derived from linear logic proofs. This suggests the LLRNN can exploit the compositional nature of logic to build semantically structured memory access patterns, a key advantage over conventional RNNs.

Although the LLRNN was not benchmarked on large-scale real-world datasets due to its complexity and experimental nature, the architectural simulations illustrate its ability to:

- emulate and extend second-order and multiplicative RNNs,
- simulate Neural Turing Machine behavior using logic-based memory transformations,
- integrate symbolic programs (e.g., binary integer encodings) with differentiable semantics.

The results support our hypothesis that integrating denotational semantics of linear logic into RNN controllers enhances expressivity and interpretability in sequence modeling.

VI. CONCLUSION

In this work, we have introduced the Linear Logic Recurrent Neural Network (LLRNN), a new class of neural architectures that bridges symbolic logic and deep learning by embedding the semantics of linear logic into the operational dynamics of recurrent neural networks. The LLRNN departs from traditional neural models by grounding its transformations not merely in learned weights, but in structured, differentiable interpretations of logical proofs, offering a powerful synthesis of formal reasoning and gradient-based learning.

Our construction is motivated by the long-standing challenge of unifying symbolic reasoning—traditionally grounded in discrete, rule-based systems—with the continuous, high-dimensional representations employed in neural networks. Linear logic serves as a crucial enabler in this endeavor. Its rich denotational semantics allow for a natural translation of logic programs into differentiable operators over vector spaces, thereby making it possible to construct neural architectures that can execute, compose, and even learn algorithms within a mathematically coherent framework.

The LLRNN model generalizes and subsumes several established RNN architectures, including second-order and multiplicative RNNs, as well as Neural Turing Machines, by showing that their memory access patterns and state transformations can be expressed as linear logic proofs. More importantly, it enables new forms of neural computation—such as higher-order transformations, logic-guided memory manipulation, and pattern-driven symbolic control—that are difficult to achieve with standard architectures. For instance, the use of integer and binary integer types, and the ability to store and execute algorithmic patterns in structured memory rings, illustrates the LLRNN’s capacity to learn symbolic behavior in a differentiable setting.

Conceptually, this model embodies a hybrid design philosophy: it combines the expressiveness and theoretical grounding of logical systems with the data-driven flexibility of machine learning. It offers a form of “algorithmic scaffolding” that allows neural networks to learn in structured spaces, guided not only by loss functions but also by the internal logic of the computation being modeled.

Despite its promise, the current LLRNN implementation is primarily a proof-of-concept. The training of such models remains nontrivial due to the complex nature of the logic-denotation interface and the lack of large-scale benchmarks tailored to logic-augmented architectures. However, our results suggest a new paradigm for interpretable and programmable machine learning models, where reasoning and learning are no longer distinct faculties but parts of the same differentiable computational fabric.

In summary, the LLRNN contributes a novel and principled approach to neural-symbolic integration. It opens up exciting opportunities for future research at the intersection of formal

logic, theoretical computer science, and deep learning. With further development, we believe this architecture could serve as the foundation for next-generation intelligent systems capable of learning, generalizing, and reasoning in a human-like, interpretable, and logically grounded manner.

A. Future Work

Several promising directions remain open for exploration:

- 1) **Scalability and Optimization:** While the current implementation is theoretically sound, practical training of LLRNNs remains challenging due to the complexity of the logic-denotation interface. Future work will explore efficient parameterizations and optimizations to make LLRNNs viable for larger datasets.
- 2) **End-to-End Differentiable Theorem Provers:** Building on recent advances in neural symbolic reasoning, LLRNNs could be extended to interface with neural theorem provers or logic program synthesizers, allowing for dynamic construction of proofs during training.
- 3) **Differential Linear Logic Integration:** Future architectures could incorporate full differential linear logic directly into the computational graph, enabling automatic differentiation over logic programs themselves and facilitating meta-learning over symbolic reasoning processes.
- 4) **Applications in Program Induction and Formal Verification:** Given the LLRNN’s logical backbone, it could be adapted to tasks involving code synthesis, program induction, and formal verification, particularly where reasoning over sequences or iterative transformations is required.
- 5) **Attention Mechanisms and Structured Decoders:** Integrating structured attention over memory rings and combining LLRNNs with attention-based encoders and decoders may yield more powerful hybrid architectures for tasks such as machine translation or question answering.

Ultimately, LLRNNs represent an initial but significant step toward unifying the symbolic and connectionist paradigms, using linear logic as a bridge. We hope this work inspires future research on embedding interpretable, structured computation within neural networks.

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], 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 (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 (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, *DeepCoder: 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: proofs, models and antiderivatives*, [arXiv:1606.01642] (2016).
- [34] J. Clift and D. Murfet, *Cofree coalgebras and differential linear logic*, preprint.
- [35] Abadi, MartĀn, et al. *TensorFlow: A system for large-scale machine learning* arXiv preprint arXiv:1605.08695 (2016).
- [36] Abadi, MartĀn, 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 λ -calculus*, Theoretical Computer Science 309, pp. 1–41, (2003).
- [39] O. Manzyuk, *A simply typed λ -calculus of forward automatic differentiation*, In Mathematical Foundations of Programming Semantics Twenty-eighth Annual Conference, pages 259–73, Bath, UK, June 6–9 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., Kaliszky, 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).