
Going Beyond Linear Transformers with Recurrent Fast Weight Programmers

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Transformers with linearised attention (“linear Transformers”) have demonstrated
2 the practical scalability and effectiveness of outer product-based Fast Weight
3 Programmers (FWPs) from the ’90s. However, the original FWP formulation is
4 more general than the one of linear Transformers: a *slow* neural network (NN)
5 continually reprograms the weights of a *fast* NN with *arbitrary* NN architecture. In
6 existing linear Transformers, both NNs are feedforward and consist of a single layer.
7 Here we explore new variations by adding recurrence to the slow and fast nets. We
8 evaluate our novel recurrent FWPs (RFWPs) on two synthetic algorithmic tasks
9 (code execution and sequential ListOps), Wikitext-103 language models, and on the
10 Atari 2600 2D game environment. Our models exhibit properties of Transformers
11 and RNNs. In the reinforcement learning setting, we report large improvements
12 over the LSTM in several Atari games. Our code is publicly available.

13 1 Introduction

14 The Transformer [1] has arguably become the most popular neural network (NN) for processing
15 text data. Its success on neural machine translation quickly transferred to other problems in natural
16 language processing, such as language modelling [2, 3] or question answering [4]. Recently, it has also
17 been applied in other domains such as image processing [5, 6] or mathematical problem solving [7, 8].

18 Conceptually, the Transformer is a deep feedforward NN that processes all elements of a sequence in
19 parallel: unlike in recurrent NNs (RNNs), the computations of a layer for the entire sequence can be
20 packed into one big matrix multiplication. This scales well with the number of parallel processors.

21 Despite the benefits of parallelisation, a major drawback of Transformers is their computational
22 complexity in time and space which is quadratic in sequence length. Furthermore, in its auto-
23 regressive version [1, 2], which is the focus of our work, the state size linearly increases with sequence
24 length. This makes it prohibitive in auto-regressive settings dealing with very long or potentially
25 infinite sequences—forcing practitioners to truncate the temporal context into a fixed-size window
26 and ignoring dependencies going beyond that window. Although certain Transformer modifications
27 try to address this issue [9, 10], reinforcement learning in partially observable environments [11, 12]
28 is still dominated by RNNs such as the Long Short-Term Memory (LSTM; [13]).

29 To scale Transformers to longer sequences, recent works have proposed to linearise the softmax in
30 the attention computation [14], reorganizing the latter in a sequential way. Such models include
31 Katharopoulos et al.’s *Linear Transformer* (LT) [14] and Choromanski et al.’s *Performer* [15]. They
32 enjoy time and space complexities linear in sequence length with states of constant size. While
33 their performance on some tasks does not fully match the one of regular Transformers [16], several
34 improvements have already been proposed [17, 18] (see our review in Sec. 2.2) which makes this
35 Transformer family a promising alternative.

36 Here we go one step further in advancing linear Transformer variants as powerful sequence processing
 37 models, adopting the perspective of “Fast Weight Programmers” (FWPs) [19, 20, 21]. Recent work
 38 emphasised that linearised self-attention is equivalent to an outer product-based FWP from the ’90s
 39 ([18]; reviewed in Sec. 2). Here we explore this connection further and develop more complex
 40 FWPs. The original FWP [19] is a two-NN system: a slow and a fast net, each with arbitrary
 41 architectures. The slow net learns to generate rapid context-dependent weight modifications for the
 42 fast net. In the case of existing linear Transformer variants, the slow and fast nets are simple one
 43 layer feedforward NNs. Here we augment them with recurrent connections to obtain recurrent FWPs
 44 (RFWPs). Recurrence enhances the model’s theoretical power [22] and can help to solve tasks that
 45 naturally require recurrence as a part of the solution.

46 Our experiments on the language modelling dataset Wikitext-103 [23] show that our RFWPs are
 47 competitive compared to regular Transformers. We then study various properties of the proposed
 48 models on two synthetic algorithmic tasks: code execution [24] and sequential ListOps [25]. Finally,
 49 it is straightforward to apply our models to reinforcement learning problems as a drop-in replacement
 50 for LSTMs. Here our RFWPs obtain large improvements over LSTM baselines across many Atari
 51 2600 2D game environments [26]. Although LSTM still works better in a few environments, we
 52 show that our RFWPs generally improve by scaling them up.

53 The main contribution of this work is twofold: (1) from the perspective of FWPs, we study alternative
 54 and potentially more powerful FWPs for sequence processing, demonstrating that NNs can easily learn
 55 to control NNs that are more complex than a single feedforward layer, and (2) from the perspective
 56 of Transformer models, our models augment linear Transformers with recurrence, addressing general
 57 limitations of existing auto-regressive Transformer models.

58 2 Background on Fast Weight Programmers (FWPs)

59 Here we review the general concept of FWPs, as well as two specific instances thereof: the linear
 60 Transformer [14, 15] and the Delta Net [18].

61 2.1 General Formulation

62 We refresh the concept of fast weight controllers or FWPs [19, 20] using modern notation in a
 63 sequence processing scenario. An FWP with trainable parameters θ_{slow} sequentially transforms an
 64 input sequence $\{\mathbf{x}^{(i)}\}_{i=1}^L$ with $\mathbf{x}^{(i)} \in \mathbb{R}^{d_{\text{in}}}$ to an output sequence $\{\mathbf{y}^{(i)}\}_{i=1}^L$ with $\mathbf{y}^{(i)} \in \mathbb{R}^{d_{\text{out}}}$ as

$$\theta_{\text{fast}}^{(i)}, \mathbf{q}^{(i)} = \text{SlowNet}(\{\mathbf{x}^{(j)}\}_{j=1}^i, \{\mathbf{y}^{(j)}\}_{j=0}^{i-1}, \{\theta_{\text{fast}}^{(j)}\}_{j=0}^{i-1}; \theta_{\text{slow}}) \quad (1)$$

$$\mathbf{y}^{(i)} = \text{FastNet}(\{\mathbf{q}^{(j)}\}_{j=1}^i, \{\mathbf{y}^{(j)}\}_{j=0}^{i-1}; \theta_{\text{fast}}^{(i)}) \quad (2)$$

65 where $\mathbf{y}^{(0)}$ and $\theta_{\text{fast}}^{(0)}$ are initial variables. In this two-NN system, the parameters $\theta_{\text{fast}}^{(i)}$ of the NN
 66 FastNet are generated by another NN SlowNet at each time step i . The weights of the fast net are
 67 *fast* in the sense that they may rapidly change at every step of the sequence while the weights of
 68 the slow net θ_{slow} are *slow* because they can only change due to gradient descent during training,
 69 remaining fixed afterwards¹. Eq. 1 expresses a slow NN in its general form. The slow net can generate
 70 certain fast weight modifications in a way that reflects certain architectural choices for the slow and
 71 fast NNs. In addition to the fast weights $\theta_{\text{fast}}^{(i)}$, the slow net also generates or *invents* an input $\mathbf{q}^{(i)}$ to be
 72 fed to the fast net (alternatively $\mathbf{q}^{(i)}$ can simply be $\mathbf{x}^{(i)}$). While the architectures of slow and fast nets
 73 are arbitrary, they are typically chosen to be differentiable such that the entire FWP can be trained in
 74 an end-to-end manner using gradient descent. By interpreting the weights of an NN as a program
 75 [27], the slow net effectively learns to control, or *program*, the fast NN. Thus, the slow net is a neural
 76 programmer of fast weights, and its parameter set θ_{slow} is a compressed form of information used to
 77 produce potentially infinite variations of context-dependent fast weights.

78 In many settings it makes sense to generate the fast weights θ_{fast} incrementally in an iterative fashion,
 79 where the SlowNet is further decomposed into two sub-parts:

$$\mathbf{z}^{(i)}, \mathbf{q}^{(i)} = \text{SlowSubnet}(\{\mathbf{x}^{(j)}\}_{j=1}^i, \{\mathbf{y}^{(j)}\}_{j=0}^{i-1}, \{\theta_{\text{fast}}^{(j)}\}_{j=0}^{i-1}; \theta_{\text{slow}}) \quad (3)$$

$$\theta_{\text{fast}}^{(i)} = \text{UpdateRule}(\theta_{\text{fast}}^{(i-1)}, \mathbf{z}^{(i)}) \quad (4)$$

¹The fast net could also contain some additional slow weights; we ignore this possibility.

80 where UpdateRule takes the fast weights $\theta_{\text{fast}}^{(i-1)}$ from the previous iteration to produce the new fast
 81 weights conditioned on $\mathbf{z}^{(i)}$. In the next section we’ll review recent work from this perspective.

82 2.2 Linear Transformers as Fast Weight Programmers

83 In general, the dimension of the fast weights $\theta_{\text{fast}}^{(i)}$ is too large to be conveniently parameterised by an
 84 NN. Instead, it was proposed [19] to perform a rank-one update via the outer product of two vectors
 85 generated by the slow net. Two recent models directly correspond to such outer product-based FWPs:
 86 linear Transformers [14] and the Delta Net [18].

87 **Linear Transformer.** The “linear Transformer” [14, 15] is a Transformer where the softmax in
 88 the attention is linearised. This is achieved by replacing the softmax with a kernel function ϕ , after
 89 which the self-attention can be rewritten as a basic outer product-based FWP [19, 18]. Previous
 90 works focused on different ϕ maps with properties such as increased capacity [18] or guaranteed
 91 approximation of the softmax in the limit [15, 17]. For our purposes, the particular choice of ϕ is
 92 irrelevant and we simply assume $\phi : \mathbb{R}^{d_{\text{key}}} \rightarrow \mathbb{R}^{d_{\text{key}}}$, simplifying our equations below by writing \mathbf{k}, \mathbf{q}
 93 instead of $\phi(\mathbf{k}), \phi(\mathbf{q})$. Using otherwise the same notation, for each new input $\mathbf{x}^{(i)}$, the output $\mathbf{y}^{(i)}$ is
 94 obtained by:

$$\mathbf{k}^{(i)}, \mathbf{v}^{(i)}, \mathbf{q}^{(i)} = \mathbf{W}_k \mathbf{x}^{(i)}, \mathbf{W}_v \mathbf{x}^{(i)}, \mathbf{W}_q \mathbf{x}^{(i)} \quad (5)$$

$$\mathbf{W}^{(i)} = \mathbf{W}^{(i-1)} + \mathbf{v}^{(i)} \otimes \mathbf{k}^{(i)} \quad (6)$$

$$\mathbf{y}^{(i)} = \mathbf{W}^{(i)} \mathbf{q}^{(i)} \quad (7)$$

95 where the slow weight matrices $\mathbf{W}_k \in \mathbb{R}^{d_{\text{key}} \times d_{\text{in}}}$ and $\mathbf{W}_v \in \mathbb{R}^{d_{\text{value}} \times d_{\text{in}}}$ are used to obtain the *key*
 96 $\mathbf{k}^{(i)} \in \mathbb{R}^{d_{\text{key}}}$ and the *value* $\mathbf{v}^{(i)} \in \mathbb{R}^{d_{\text{value}}}$. The key and value vectors are used to generate new weights
 97 via the outer product $\mathbf{v}^{(i)} \otimes \mathbf{k}^{(i)} \in \mathbb{R}^{d_{\text{value}} \times d_{\text{key}}}$. A further simplification in the equations above is the
 98 omission of attention normalisation which has been experimentally shown to be unnecessary if the ϕ
 99 function produces normalised scores [18].

100 In Eq. 6, the previous fast weight matrix $\mathbf{W}^{(i-1)} \in \mathbb{R}^{d_{\text{value}} \times d_{\text{key}}}$ is updated to yield $\mathbf{W}^{(i)}$ by adding
 101 the update term $\mathbf{v}^{(i)} \otimes \mathbf{k}^{(i)}$. This corresponds to the *sum update rule*. Here the fast NN is a simple
 102 linear transformation as in Eq. 7 which takes as input the (query) vector $\mathbf{q}^{(i)} \in \mathbb{R}^{d_{\text{key}}}$ generated by
 103 the slow weights $\mathbf{W}_q \in \mathbb{R}^{d_{\text{key}} \times d_{\text{in}}}$. Hence, in linear Transformers, the previous Eq. 3 simplifies to:
 104 $\mathbf{z}^{(i)}, \mathbf{q}^{(i)} = \text{SlowSubnet}(\mathbf{x}^{(i)}; \theta_{\text{slow}})$ with $\mathbf{z}^{(i)} = (\mathbf{k}^{(i)}, \mathbf{v}^{(i)})$.

105 **Delta Net.** The Delta Net [18] is obtained by replacing the sum update rule (Eq. 6) in the linear
 106 Transformer with the *delta update rule* [28]:

$$\mathbf{W}^{(i)} = \mathbf{W}^{(i-1)} + \beta^{(i)} (\mathbf{v}^{(i)} - \bar{\mathbf{v}}^{(i)}) \otimes \mathbf{k}^{(i)} \quad (8)$$

107 where $\beta^{(i)} \in \mathbb{R}$ is a fast parameter (learning rate) of the update rule generated by the slow net with
 108 weights $\mathbf{W}_\beta \in \mathbb{R}^{d_{\text{in}} \times 1}$ and the sigmoid function σ :

$$\beta^{(i)} = \sigma(\mathbf{W}_\beta \mathbf{x}^{(i)}) \quad (9)$$

109 and $\bar{\mathbf{v}}^{(i)} \in \mathbb{R}^{d_{\text{value}}}$ is generated as a function of the previous fast weights $\mathbf{W}^{(i-1)}$ and the key $\mathbf{k}^{(i)}$

$$\bar{\mathbf{v}}^{(i)} = \mathbf{W}^{(i-1)} \mathbf{k}^{(i)}. \quad (10)$$

110 Schlag et al. [18] introduce the Delta Net to address a memory capacity problem affecting
 111 linear Transformers with the sum update rule. The corresponding Eq. 3 is: $\mathbf{z}^{(i)}, \mathbf{q}^{(i)} =$
 112 $\text{SlowSubnet}(\mathbf{x}^{(i)}, \mathbf{W}^{(i-1)}; \theta_{\text{slow}})$ with $\mathbf{z}^{(i)} = (\mathbf{k}^{(i)}, \mathbf{v}^{(i)}, \beta^{(i)}, \bar{\mathbf{v}}^{(i)})$. Thus, unlike linear Transform-
 113 ers, the SlowNet in the Delta Net takes the previous fast weights $\mathbf{W}^{(i-1)}$ into account to generate
 114 the new fast weight updates.

115 In summary, these models can be described as FWPs whose slow and fast net each consists of a single
 116 feedforward layer (Eqs. 5 and 7). In the next section we present several extensions.

117 **Other Approaches.** While our focus here is on outer product-based weight generation, which is
 118 an efficient method to handle high dimensional NN weights, there are also other approaches. For
 119 example, instead of generating a new weight matrix, Hypernetworks scale the rows of a slow weight
 120 matrix with a generated vector of appropriate size. In the broad sense of context-dependent weights,
 121 many concepts relate to FWPs, e.g. dynamic convolution [29, 30, 31], lambda networks [32], or
 122 dynamic plasticity [33, 34, 35].

123 3 Fast Weight Programmers with slow or fast RNNs

124 The original formulation of FWPs reviewed in Sec. 2.1 is more general than existing models presented
 125 in Sec. 2.2. In particular, both fast and slow networks in existing linear Transformers are one-layer
 126 feedforward networks (Eq. 7). Here we present FWPs with recurrent fast nets in Sec. 3.1 and FWPs
 127 with recurrent slow nets in Sec. 3.2.

128 3.1 Fast Network Extensions

129 In principle, any NN architecture can be made *fast*. Its fast weight version is obtained by replacing
 130 the networks weights with fast weights parameterised by an additional slow network. For example,
 131 consider a regular RNN layer with two weight matrices \mathbf{W} and \mathbf{R} :

$$\mathbf{h}^{(i)} = \sigma(\mathbf{W}\mathbf{x}^{(i)} + \mathbf{R}\mathbf{h}^{(i-1)}) \quad (11)$$

132 A fast weight version can be obtained by replacing \mathbf{W} and \mathbf{R} with $\mathbf{W}^{(i)}$ and $\mathbf{R}^{(i)}$ which are
 133 controlled as in Eq. 8 with all necessary variables generated by a separate slow net at each time step i .

134 While this view illustrates the generality of FWPs, the angle under which we approach these models
 135 here is slightly different: we introduce recurrence as a way of augmenting existing linear Transform-
 136 ers.

137 **Delta RNN.** We obtain a fast weight RNN called **Delta RNN** by adding an additional recurrent
 138 term to the feedforward fast net of the FWP (Eq. 7):

$$\mathbf{y}^{(i)} = \mathbf{W}^{(i)}\mathbf{q}^{(i)} + \mathbf{R}^{(i)}f(\mathbf{y}^{(i-1)}) \quad (12)$$

139 where $\mathbf{R}^{(i)} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{out}}}$ is an additional fast weight matrix which introduces recurrent connections.
 140 It is also generated by the slow net using the delta update rule, similar to $\mathbf{W}^{(i)}$ in Eq. 8 but with
 141 additional slow weights. We apply an element-wise activation function f to the previous output of
 142 the fast network $\mathbf{y}^{(i-1)}$ to obtain the recurrent query. The choice of activation function is crucial
 143 here because, to achieve stable model behaviour, the elements in key and query vectors should be
 144 positive and sum up to one when the delta update rule is used [18]. We use the softmax function
 145 ($f = \text{softmax}$ in Eq. 12) to satisfy these conditions.

146 Analogous to the Delta RNN, we also construct a **Delta LSTM** with six fast weight matrices. The
 147 exact equations can be found in Appendix A.

148 **Alternative Feedforward Fast Nets.** While the focus of this work is on RNNs, there are also
 149 interesting fast feedforward models to be used in Eq. 7 which might result in stronger feedforward
 150 baselines. For example, we can simply replace the one-layer transformation of Eq. 7 by a K -layer
 151 deep network:

$$\mathbf{h}_k^{(i)} = \mathbf{W}_k^{(i)}f(\mathbf{h}_{k-1}^{(i)}) \quad \text{for } k \in [1..K] \quad \text{with } \mathbf{h}_0^{(i)} = \mathbf{q}^{(i)} \quad (13)$$

$$\mathbf{y}^{(i)} = \mathbf{h}_K^{(i)} \quad (14)$$

152 where the slow network produces all K fast weights $\{\mathbf{W}_k^{(i)}\}_{k=1}^K$ and query $\mathbf{q}^{(i)}$ from a single
 153 input $\mathbf{x}^{(i)}$. In light of the capacity limitation in linear Transformers [18], this might introduce
 154 additional capacity without the need of larger representations, analogous to the trade-off in a multilayer
 155 perceptron (MLP) between narrow & deep versus shallow & wide. We refer to this class of models as
 156 **Delta MLPs**. Again, for stable model behaviour with the delta rule, we apply the softmax activation
 157 f to the vectors to be used as a query.

158 Another interesting approach is to use a Delta Net itself as a fast net, i.e., make the slow weights
 159 in the Delta Net fast (thus obtaining a **Delta Delta Net**). Such a model could in principle learn to
 160 improve the way of generating fast weights depending on the context. While we plan to investigate
 161 this potential in future work, we also include preliminary results of such a model in our language
 162 modelling experiments (Sec. 4.1).

163 We note that there are several previously proposed recurrent fast weight models. For example,
 164 Schmidhuber [21]’s FWP have been revisited by Ba et al. [36]. There, key and value vectors are
 165 not generated within the same time step, unlike in our models or in linear Transformers. Another
 166 related recurrent FWP is the Fast Weight Memory (FWM) [37]. The FWM fits neatly in the RFWP
 167 framework: the slow net is an LSTM and the fast net is a higher-order RNN. However, the FWM is a
 168 single pair of slow and fast nets, and a multi-layer version, as in the linear Transformer family, were
 169 not explored. Similarly, the Metalearned Neural Memory [35] uses an LSTM as its slow net and a
 170 deep MLP as its fast net but again limited to one pair. Others have investigated variants of RNNs
 171 with fast weights for toy synthetic retrieval tasks [38, 39]. In particular, Keller et al. [39] augment the
 172 LSTM with a fast weight matrix in the cell update. In contrast, we make all weights in the LSTM fast
 173 and, importantly, our model specifications build upon the successful Transformer architecture.

174 3.2 Slow Network Extensions

175 In linear Transformers, the slow network is purely feedforward (Eq. 5). It can be made recurrent at
 176 two different levels: within the slow network (i.e. the slow network computes weight updates based
 177 on its own previous outputs: e.g. key, value, query vectors) or via the fast network by taking the fast
 178 net’s previous output as an input. In our preliminary experiments, we found the former to be sub-
 179 optimal (at least in language modelling experiments). So we focus on the latter approach: we make
 180 the slow net in the Delta Net dependent on the previous output of the fast network. We refer to this
 181 model as the **Recurrent Delta Net** (RDN).

182 **Recurrent Delta Net.** We obtain the RDN by modifying the generation of key, value, and query
 183 vectors in Eq. 5 as well as the learning rate in Eq. 9 for the delta update rule. We add additional slow
 184 weights (\mathbf{R}_k , \mathbf{R}_v , \mathbf{R}_q , and \mathbf{R}_β) for recurrent connections which connect the previous output of the
 185 fast net $\mathbf{y}^{(i-1)}$ (Eq. 7) to the new $\mathbf{k}^{(i)}$, $\mathbf{v}^{(i)}$, $\mathbf{q}^{(i)}$, and $\beta^{(i)}$ as follows:

$$\mathbf{k}^{(i)} = \mathbf{W}_k \mathbf{x}^{(i)} + \mathbf{R}_k \tanh(\mathbf{y}^{(i-1)}) \quad (15)$$

$$\mathbf{v}^{(i)} = \mathbf{W}_v \mathbf{x}^{(i)} + \mathbf{R}_v \tanh(\mathbf{y}^{(i-1)}) \quad (16)$$

$$\mathbf{q}^{(i)} = \mathbf{W}_q \mathbf{x}^{(i)} + \mathbf{R}_q \tanh(\mathbf{y}^{(i-1)}) \quad (17)$$

$$\beta^{(i)} = \sigma(\mathbf{W}_\beta \mathbf{x}^{(i)} + \mathbf{R}_\beta \tanh(\mathbf{y}^{(i-1)})) \quad (18)$$

186 While the rest of the model remains as in the Delta Net, with these simple extra recurrent connections
 187 the model becomes a proper RNN. The corresponding dependencies in Eq. 3 are: $\mathbf{z}^{(i)}$, $\mathbf{q}^{(i)} =$
 188 $\text{SlowSubnet}(\mathbf{x}^{(i)}, \mathbf{y}^{(i-1)}, \mathbf{W}^{(i-1)}; \boldsymbol{\theta}_{\text{slow}})$ with $\mathbf{z}^{(i)} = (\mathbf{k}^{(i)}, \mathbf{v}^{(i)}, \beta^{(i)}, \bar{\mathbf{v}}^{(i)})$.

189 In contrast to regular RNNs, all these fast weight variants maintain the standard one-dimensional
 190 states and the two-dimensional fast weight states which both serve as additional memory. In our
 191 experiments (Sec. 4), we demonstrate that slow NNs can learn to control the weights of these rather
 192 complex RNN variants, and that the introduction of recurrence in the fast network introduces a crucial
 193 property that is missing in regular Transformers.

194 4 Experiments

195 We conduct experiments in four different settings. We start by evaluating all models on a language
 196 modelling task (Sec. 4.1) to obtain a performance overview and to discuss computational costs.
 197 Language modelling is an excellent task to evaluate sequence models. However, to highlight their
 198 different capabilities we evaluate our models also on algorithmic tasks. In fact, it is well-known
 199 that the actual capabilities of RNNs differ from one architecture to another [40]. We are interested
 200 in discussing such differences. With that goal in mind, we conduct experiments on two synthetic
 201 algorithmic tasks, code execution (Sec. 4.2) and sequential ListOps (Sec. 4.3), which are designed
 202 to compare elementary sequence processing abilities of models. Finally, we apply our models to
 203 reinforcement learning in 2D game environments (Sec. 4.4) as a replacement for LSTMs.

Table 1: WikiText-103 language model perplexity results with the common *small* setting [17, 18]. For each model, its name, corresponding slow and fast networks, and weight update rule (Update) is specified. All models are trained and evaluated on the span of 256 tokens except for the models in the last two rows (+ full context) which are trained and evaluated without context truncation. Training language models on Wikitext-103 is a resource-demanding experiment which is why we only provide results from a single run. Parameter count is in millions.

Name	Slow net	Update	Fast net	Valid	Test	#Prms
Transformer	-	-	-	33.0	34.1	44.0
Linear Transformer	Feedforward	sum	Linear	37.1	38.3	44.0
Delta Net		delta		34.1	35.2	44.0
Delta MLP	Feedforward	delta	Deep MLP	35.8	36.8	44.3
Delta Delta Net			Delta Net	34.0	35.2	44.6
Delta RNN			RNN	33.8	35.0	44.6
Delta LSTM			LSTM	33.4	34.7	47.3
RDN	Recurrent		Linear	34.1	35.2	44.1
Delta RNN	+ full context			31.8	32.8	44.6
RDN	+ full context			32.5	33.6	44.1

204 4.1 Language Modelling

205 We first evaluate all discussed models on the generic language modelling task. This allows us to
 206 obtain a performance overview and review the computational efficiency of different models. We use
 207 the Wikitext-103 dataset [23] and follow the *small model setting* similar to those used in recent works
 208 by Peng et al. [17] and Schlag et al. [18]. This allows us to train and evaluate different models with a
 209 reasonable amount of compute on this resource-demanding language modelling task.

210 **Perplexity results.** See Table 1 which also serves as a tabular summary recapitulating different
 211 models described in Sec. 2 and 3 with various architectures for slow and fast nets, and the choice
 212 of update rule. The top block of Table 1 shows the performance of the baseline Transformer,
 213 Katharopoulos et al. [14]’s Linear Transformer, and Schlag et al. [18]’s Delta Net. The performance of
 214 models presented in Sec. 3 can be found in the middle block. First of all, Delta MLP performs worse
 215 than the baseline Delta Net despite a slight increase in parameter count (44.3 vs. 44.0 M). This supports
 216 the intuition that it is better to make the slow network aware of the outputs of intermediate layers to
 217 generate fast weights in a deep network, instead of generating fast weights for all layers at a time. In
 218 all other models, the performance never degrades with the proposed architectural augmentation. The
 219 Delta Delta Net yields limited improvements; we plan to study this model in depth in future work.
 220 With the same amount of parameters (44.6 M), Delta RNN yields greater improvements. Among
 221 the models presented here, the Delta LSTM variant exhibits the best performance. This shows that
 222 the slow network successfully controls the rather complex fast LSTM network, although it also
 223 requires more parameters (47.3 M) compared to other models. Finally, the benefits of recurrent
 224 connections added to the baseline Delta Net do not directly translate into practical improvements in
 225 language modelling as demonstrated by the performance of RDN compared to the baseline Delta
 226 Net. Importantly, given a constant memory size w.r.t. sequence length, it is straight-forward to train
 227 and evaluate our RNNs without context truncation (while still limiting the backpropagation span).
 228 Corresponding performances of Delta RNN and RDN are shown in the bottom part of Table 1: they
 229 are better than the regular Transformer with a limited context (256 tokens).

230 While language modelling is useful as a first sanity check (here for example, except for Delta Delta
 231 Net, all models achieve reasonable performance), the task is too generic to identify certain important
 232 aspects of the models, such as real benefits of recurrence. Before we move on to trickier reinforcement
 233 learning applications, Sec. 4.2 and 4.3 will focus on studying such aspects using synthetic algorithmic
 234 tasks.

235 **Computational efficiency.** What are the additional computational costs introduced to linear Trans-
 236 formers by the modifications described in Sec. 3? First of all, none of them affect the core complexity

237 of linear Transformers: they all have a constant space and linear time complexity w.r.t. sequence
 238 length. However, the per-time-step computational costs differ a lot from one model to another, as
 239 quantified here in terms of our implementation’s training speed. All models are implemented using a
 240 custom CUDA kernel except the baseline Transformer for which we use regular PyTorch code. Training
 241 speeds of LT and Delta Net in Table 1 are 66 K and 63 K words per second respectively (vs. 33 K
 242 for the baseline Transformer). The most expensive model is the Delta LSTM. This fast LSTM with
 243 tied input-forget gates has 6 weight matrices, and each of these are manipulated by separate delta rules.
 244 The corresponding speed is 14 K words per second, which is prohibitively slow to scale for more ex-
 245 periments. In contrast, the speeds of Delta RNN and RDN remain reasonable: 41 K and 35 K words
 246 per second respectively. Therefore, the remaining experiments will focus on these two recurrent ar-
 247 chitectures which are promising and practical in terms of both performance and computational costs.

248 4.2 Code Execution Task: Learning to Maintain and Update Variable States

249 In code execution tasks [24], models are trained to sequentially read the input code provided as
 250 word-level text, and to predict the results of the corresponding code execution. We adopt the task
 251 setting from Fan et al. [41] with one conditional and three basic statements. We refer the readers
 252 to Appendix B for a precise description of the task. This code execution task requires models to
 253 maintain the values of multiple variables, which has been shown to be difficult for Transformers with
 254 only feedforward connections [41].

255 The left block of Table 2 shows the results. Following again Fan et al. [41], we control the task
 256 difficulty by modifying the number of variables (3 or 5). The model architectures are fixed: the
 257 LSTM has only one layer with 256 nodes and all Transformer variants have the same architecture
 258 with 4 layers with a hidden size of 256 using 16 heads and an inner feedforward layer size of 1024.

259 We first note that the LSTM is the best performer for both difficulty levels with the smallest perfor-
 260 mance drop through variable increase. In contrast to prior claims [41], the LSTM is clearly capa-
 261 ble of storing the values of multiple variables in a single vector. With three variables, the regular
 262 Transformer already largely underperforms other models with a mutable memory: Delta Net, Delta
 263 RNN, and RDN. Linear Transformers completely fail at this task, likely due to the memory capacity
 264 problem pointed out by Schlag et al. [18]. By increasing number of variables to five, the baseline
 265 Transformers, Delta Net, and RDN become unstable as shown by high standard deviations w.r.t. the
 266 seed. The benefits of recurrent connections introduced in our RDN compared to the baseline Delta
 267 Net become more apparent (76.3 vs. 61.4%). In contrast, Delta RNN remains stable and gives the
 268 best performance (85.1%) among Transformer variants, which shows the benefits of recurrence and
 269 in particular the regular RNN architecture. To match the performance of LSTM on this task, however,
 270 these models need more layers (see Appendix B for additional results).

Table 2: **Test accuracies (%)** with standard deviations on **code execution** (Code Exec) and **sequential ListOps** (Seq ListOps). The difficulty of the task is controlled by the maximum number of possible variables (# variables) for code execution, and the list depth (10 or 15) for ListOps. For code execution with 5 variables, we report means over six seeds. In all other cases, the results are computed with three seeds. For more results, see Appendix B.

	Code Exec (# variables)		Seq ListOps (depth)	
	3	5	10	15
LSTM	99.0 ± 0.1	93.2 ± 6.1	88.5 ± 2.9	24.4 ± 1.1
Transformer	71.8 ± 2.6	35.4 ± 28.2	79.1 ± 0.9	75.3 ± 0.4
Linear Transformer	0.0 ± 0.0	0.0 ± 0.0	64.0 ± 0.3	64.4 ± 0.4
Delta Net	90.7 ± 2.7	61.4 ± 20.0	85.7 ± 1.8	77.6 ± 1.4
Delta RNN	90.8 ± 1.7	85.1 ± 1.9	83.6 ± 1.2	78.0 ± 1.0
RDN	92.6 ± 2.2	76.3 ± 17.6	83.2 ± 0.9	79.2 ± 1.4

271 **4.3 Sequential ListOps: Learning Hierarchical Structure and Computation**

272 The ListOps task [25] requires list operation execution which is a typical test for hierarchical
 273 structure learning. We use a simple variant of ListOps whose detailed descriptions can be found in
 274 Appendix B. For example, the list [MAX 6 1 [FIRST 2 3] 0 [MIN 4 7 1]] is of depth two
 275 and the expected output is 6. While early research comparing self-attention to RNNs [42] has shown
 276 some advantages of recurrence in hierarchical structure learning, more recent work [43] reports
 277 Transformers outperforming LSTMs on ListOps. According to Tay et al. [16], linear Transformer
 278 variants (LT and Performers) underperform other Transformer variants by a large margin on ListOps.

279 The right block of Table 2 shows results for two different depths: 10 and 15. The model architectures
 280 are identical to those used in the code execution task (Sec. 4.2). At depth 10, we find LSTM to perform
 281 best, while mutable memory Transformer variants (Delta Net, Delta RNN, and RDN) outperform the
 282 regular and linear Transformers. At depth 15, the LSTM’s performance drops drastically (to 24.4%),
 283 while the differences between Transformer variants remain almost the same. We note that sequences
 284 are longer for the depth 15 problem (mean length of 185 tokens) than for the depth 10 version (mean
 285 length of 98 tokens). This might add extra difficulty for the small 256-dimensional LSTM. The
 286 performance differences between the baseline Delta Net and the proposed Delta RNN and RDN are
 287 rather small for this task. Importantly, introduction of recurrence does not hurt for this task requiring
 288 hierarchical structure learning, and our novel models outperform both regular and linear Transformers.

289 **4.4 Reinforcement Learning in 2D Game Environments**

290 We finally evaluate the performance of our models as a direct replacement for the LSTM in reinforce-
 291 ment learning settings. In fact, only a limited number of prior works have investigated Transformers
 292 for reinforcement learning. Parisotto et al. [11] evaluate them on the DMLab-30 [44, 45] and on Atari
 293 but in a multi-task setting [46]. Others [41, 12] use toy maze environments. In contrast to Parisotto
 294 et al. [11]’s work, which presents multi-task Atari as a side experiment, we study the Transformer
 295 family of models on the standard Atari 2600 setting [26, 47, 48] by training game specific agents.

296 **Settings.** We train an expert agent on each game separately with the Importance Weighted Actor-
 297 Learner Training Architecture (IMPALA) using the V-trace actor-critic setup [49] and entropy
 298 regularization [50] implemented in Torchbeast [51]. Our model follows the *large* architecture of
 299 Espeholt et al. [49] which consists of a 15-layer residual convolutional NN with one 256-node LSTM
 300 layer which we replace by either the RDN (Sec. 3.2) or the Delta RNN (Sec. 3.1). In line with the
 301 small LSTM used for Atari (only 1 layer with 256 hidden nodes) we also configure a small RDN: 2
 302 layers with a hidden size of 128 using 4 heads, and a feedforward dimension of 512. We find this
 303 small model to perform already surprisingly well. For the rest, we use the same hyperparameters as
 304 Espeholt et al. [49] which can be found in Appendix C.

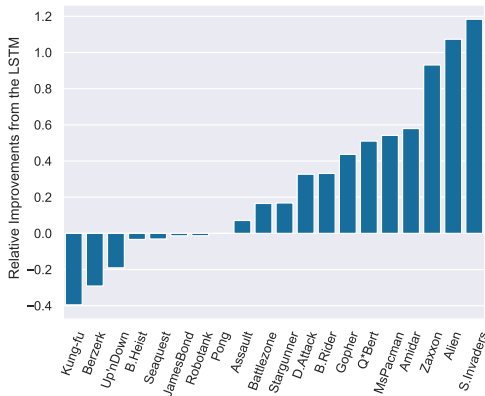


Figure 1: Rel. improvements in test scores obtained by 2-layer RDN compared to LSTM after 50M env. steps.

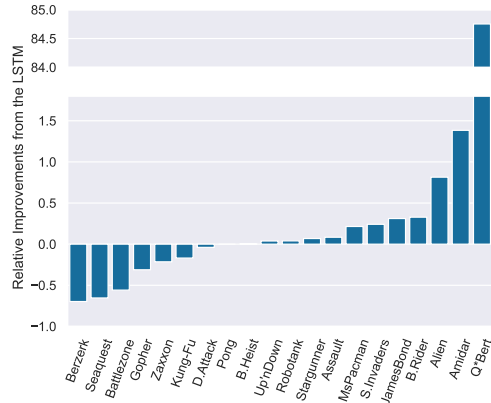


Figure 2: Rel. improvements in test scores obtained by 2-layer RDN compared to LSTM after 200M env. steps.

305 **Main experiments.** We evaluate our models in 20 environments. According to Mott et al. [52], in
 306 about half of them the LSTM outperforms the feedforward baselines — which we confirm in our
 307 setting with 50 M steps (see Appendix C). We report results at 50 M and 200 M environmental steps
 308 of training. As done by Nair et al. [53], we run the trained agent for 30 test episodes. Here we repeat
 309 this evaluation 5 times to report the average score with a standard deviation. The following analysis
 310 focuses on the RDN (Sec. 3.2). A similar study of the Delta RNN can be found in Appendix C.

311 Figure 1 shows relative improvements of RDN over LSTM after 50 M interactions. In 12 games,
 312 the RDN yields improvements over LSTM, whereas on 3 games, the LSTM performs better. In
 313 the remaining 5 games, both reach similar scores. The exact scores can be found in Appendix C.
 314 Interestingly, this trend does not directly extrapolate to the 200 M case, which is presented in Figure
 315 2. With more training, the LSTM surpasses the performance of the RDN on *Battlezone*, *Gopher*,
 316 *Seaquest* and *Zaxxon*, while the RDN catches up on *Up’N Down* and *Kung-Fu Master*. Overall, there
 317 are 6 games on which LSTM clearly outperforms RDN at 200 M steps, whereas in 9 games the result
 318 is the opposite.

319 On a side note, some of the scores achieved by the RDN at 200 M step are excellent: a score of over
 320 170 K and 980 K on *Space Invader* and *Q*Bert* respectively beats the state-of-the-art set by MuZero
 321 [54] and Agent57 [46]. However, a direct comparison is not fair as we train game-specific agents.

322 **Experiments with larger models.** Given the results above, a natural question to ask is whether a
 323 larger model size improves the RDN on games where the LSTM dominates. We focus on four such
 324 environments: *Battlezone*, *Berzerk*, *Gopher*, and *Seaquest* (See Fig. 2). We double the model size to
 325 3.4 M parameters by increasing the number of layers to 4 and the hidden size to 256, with 8 heads. As
 326 shown in Table 3, larger RDN models reduce the gap to the LSTM (except in *Berzerk*). This indicates
 327 that further scaling RDN might be as promising as scaling regular Transformers in other domains.

Table 3: Performance of larger RDN on **games where the LSTM dominates** (200 M steps).

	Battlezone	Berzerk	Gopher	Seaquest
LSTM	24,873 ± 1,240	1,150 ± 92	124,914 ± 22,422	12,643 ± 1,627
RDN	10,980 ± 1,104	348 ± 17	86,008 ± 11,815	4,373 ± 504
RDN larger	28,273 ± 5,333	346 ± 9	118,273 ± 14,872	14,601 ± 712

328 5 Conclusion

329 We explore the connection between linear Transformers and Fast Weight Programmers and propose
 330 various new linear Transformer variants with recurrent connections. Our novel Recurrent Fast Weight
 331 Programmers (RFPWs) outperform previous linear and regular Transformers on a code execution
 332 task and significantly improve over Transformers in a sequential ListOps task. On Wikitext-103
 333 in the “small” model setting, RFPWs perform competitively compared to previous best linear Trans-
 334 former variants for truncated contexts and beat regular Transformers for full contexts. Our RFPWs
 335 can also be used as drop-in replacements for problems where RNNs are still dominant. In particular,
 336 we evaluate them in a reinforcement learning settings on 20 Atari 2600 environments. We obtain
 337 clear improvements over LSTMs across many environments with small models and demonstrate
 338 promising scaling properties for larger models. Our work highlights the usefulness of the FWP
 339 framework from the ’90s and its connection to modern architectures, opening promising avenues
 340 for further research into new classes of recurrent Transformers.

341 **Broader impact and limitations.** Development of Transformers has resulted in many useful real
 342 world applications involving text, audio, and vision. Our work contributes in advancing linear attention
 343 variants thereof which are crucial for auto-regressive problems with long sequences because of their
 344 better time and space complexity. Our methods are evaluated on standard benchmarks which permit a
 345 direct comparison with similar work. However, applying such models to real world applications would
 346 require extra experiments and scaling up these models on larger data can be costly in terms of energy
 347 consumption. Also, training our models on biased data can result in learning undesirable properties.

348 **References**

- 349 [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
350 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proc. Advances in Neural*
351 *Information Processing Systems (NIPS)*, pages 5998–6008, Long Beach, CA, USA, December
352 2017.
- 353 [2] Rami Al-Rfou, Dokook Choe, Noah Constant, Mandy Guo, and Llion Jones. Character-level
354 language modeling with deeper self-attention. In *Proc. Conference on Artificial Intelligence*
355 *(AAAI)*, pages 3159–3166, Honolulu, HI, USA, January 2019.
- 356 [3] Tom B Brown et al. Language models are few-shot learners. In *Proc. Advances in Neural*
357 *Information Processing Systems (NeurIPS)*, Virtual only, December 2020.
- 358 [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of
359 deep bidirectional transformers for language understanding. In *Proc. North American Chapter*
360 *of the Association for Computational Linguistics on Human Language Technologies (NAACL-*
361 *HLT)*, pages 4171–4186, Minneapolis, MN, USA, June 2019.
- 362 [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
363 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
364 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
365 recognition at scale. In *Int. Conf. on Learning Representations (ICLR)*, Virtual only, May 2021.
- 366 [6] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable
367 DETR: Deformable transformers for end-to-end object detection. In *Int. Conf. on Learning*
368 *Representations (ICLR)*, Virtual only, May 2021.
- 369 [7] David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. Analysing mathematical
370 reasoning abilities of neural models. In *Int. Conf. on Learning Representations (ICLR)*, New
371 Orleans, LA, USA, May 2019.
- 372 [8] Francois Charton, Amaury Hayat, and Guillaume Lample. Learning advanced mathematical
373 computations from examples. In *Int. Conf. on Learning Representations (ICLR)*, Virtual only,
374 May 2021.
- 375 [9] Zihang Dai, Zhilin Yang, Yiming Yang, William W Cohen, Jaime Carbonell, Quoc V Le, and
376 Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length
377 context. In *Proc. Association for Computational Linguistics (ACL)*, pages 2978–2988, Florence,
378 Italy, July 2019.
- 379 [10] Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap.
380 Compressive transformers for long-range sequence modelling. In *Int. Conf. on Learning*
381 *Representations (ICLR)*, Virtual only, April 2020.
- 382 [11] Emilio Parisotto, Francis Song, Jack Rae, Razvan Pascanu, Caglar Gulcehre, Siddhant Jayaku-
383 mar, Max Jaderberg, Raphael Lopez Kaufman, Aidan Clark, Seb Noury, et al. Stabilizing Trans-
384 formers for reinforcement learning. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages
385 7487–7498, Virtual only, July 2020.
- 386 [12] Emilio Parisotto and Ruslan Salakhutdinov. Efficient transformers in reinforcement learning
387 using actor-learner distillation. In *Int. Conf. on Learning Representations (ICLR)*, Virtual only,
388 May 2021.
- 389 [13] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):
390 1735–1780, 1997.
- 391 [14] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers
392 are rnns: Fast autoregressive transformers with linear attention. In *Proc. Int. Conf. on Machine*
393 *Learning (ICML)*, Virtual only, July 2020.
- 394 [15] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane,
395 Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking
396 attention with performers. In *Int. Conf. on Learning Representations (ICLR)*, Virtual only, 2021.

- 397 [16] Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao,
398 Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena: A benchmark for efficient
399 transformers. In *Int. Conf. on Learning Representations (ICLR)*, Virtual only, May 2021.
- 400 [17] Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah A Smith, and Lingpeng Kong.
401 Random feature attention. In *Int. Conf. on Learning Representations (ICLR)*, Virtual only, 2021.
- 402 [18] Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast
403 weight memory systems. In *Proc. Int. Conf. on Machine Learning (ICML)*, Virtual only, July
404 2021.
- 405 [19] Jürgen Schmidhuber. Learning to control fast-weight memories: An alternative to recurrent
406 nets. Technical Report FKI-147-91, Institut für Informatik, Technische Universität München,
407 March 1991.
- 408 [20] Jürgen Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic
409 recurrent networks. *Neural Computation*, 4(1):131–139, 1992.
- 410 [21] Jürgen Schmidhuber. Reducing the ratio between learning complexity and number of time
411 varying variables in fully recurrent nets. In *International Conference on Artificial Neural
412 Networks (ICANN)*, pages 460–463, Amsterdam, Netherlands, September 1993.
- 413 [22] Michael Hahn. Theoretical limitations of self-attention in neural sequence models. *Transactions
414 of the Association for Computational Linguistics*, 8:156–171, 2020.
- 415 [23] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
416 models. In *Int. Conf. on Learning Representations (ICLR)*, Toulon, France, April 2017.
- 417 [24] Wojciech Zaremba and Ilya Sutskever. Learning to execute. *Preprint arXiv:1410.4615*, 2014.
- 418 [25] Nikita Nangia and Samuel Bowman. ListOps: A diagnostic dataset for latent tree learning.
419 In *Proc. North American Chapter of the Association for Computational Linguistics (NAACL):
420 Student Research Workshop*, pages 92–99, New Orleans, LA, USA, June 2018.
- 421 [26] Marc G. Bellemare, Georg Ostrovski, Arthur Guez, Philip S. Thomas, and Rémi Munos.
422 Increasing the action gap: New operators for reinforcement learning. In *Proc. AAAI Conf. on
423 Artificial Intelligence*, pages 1476–1483, Phoenix, AZ, USA, February 2016. AAAI Press.
- 424 [27] Jürgen Schmidhuber. Making the world differentiable: On using fully recurrent self-supervised
425 neural networks for dynamic reinforcement learning and planning in non-stationary environ-
426 ments. Technical Report FKI-126-90, Institut für Informatik, Technische Universität München,
427 1990.
- 428 [28] Bernard Widrow and Marcian E Hoff. Adaptive switching circuits. In *Proc. IRE WESCON
429 Convention Record*, pages 96–104, Los Angeles, CA, USA, August 1960.
- 430 [29] Benjamin Klein, Lior Wolf, and Yehuda Afek. A dynamic convolutional layer for short
431 rangeweather prediction. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition
432 (CVPR)*, pages 4840–4848, Boston, MA, USA, June 2015.
- 433 [30] Hyeonwoo Noh, Paul Hongsuck Seo, and Bohyung Han. Image question answering using
434 convolutional neural network with dynamic parameter prediction. In *Proc. IEEE Conf. on
435 Computer Vision and Pattern Recognition (CVPR)*, pages 30–38, Las Vegas, NV, USA, 2016.
- 436 [31] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. In
437 *Proc. Advances in Neural Information Processing Systems (NIPS)*, pages 667–675, Barcelona,
438 Spain, 2016.
- 439 [32] Irwan Bello. Lambdanetworks: Modeling long-range interactions without attention. In *Int.
440 Conf. on Learning Representations (ICLR)*, Virtual only, May 2021.
- 441 [33] Thomas Miconi, Kenneth Stanley, and Jeff Clune. Differentiable plasticity: training plastic
442 neural networks with backpropagation. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages
443 3559–3568, Stockholm, Sweden, July 2018.

- 444 [34] Thomas Miconi, Aditya Rawal, Jeff Clune, and Kenneth O. Stanley. Backpropamine: training
445 self-modifying neural networks with differentiable neuromodulated plasticity. In *Int. Conf. on*
446 *Learning Representations (ICLR)*, New Orleans, LA, USA, May 2019.
- 447 [35] Tsendsuren Munkhdalai, Alessandro Sordoni, Tong Wang, and Adam Trischler. Metalearned
448 neural memory. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, pages
449 13310–13321, Vancouver, Canada, December 2019.
- 450 [36] Jimmy Ba, Geoffrey E Hinton, Volodymyr Mnih, Joel Z Leibo, and Catalin Ionescu. Using
451 fast weights to attend to the recent past. In *Proc. Advances in Neural Information Processing*
452 *Systems (NIPS)*, pages 4331–4339, Barcelona, Spain, December 2016.
- 453 [37] Imanol Schlag, Tsendsuren Munkhdalai, and Jürgen Schmidhuber. Learning associative infer-
454 ence using fast weight memory. In *Int. Conf. on Learning Representations (ICLR)*, Virtual only,
455 May 2021.
- 456 [38] Imanol Schlag and Jürgen Schmidhuber. Gated fast weights for on-the-fly neural program
457 generation. In *NIPS Metalearning Workshop*, Long Beach, CA, USA, December 2017.
- 458 [39] T Anderson Keller, Sharath Nittur Sridhar, and Xin Wang. Fast weight long short-term memory.
459 *Preprint arXiv:1804.06511*, 2018.
- 460 [40] Gail Weiss, Yoav Goldberg, and Eran Yahav. On the practical computational power of finite
461 precision rnns for language recognition. In Iryna Gurevych and Yusuke Miyao, editors, *Proc.*
462 *Association for Computational Linguistics (ACL)*, pages 740–745, Melbourne, Australia, July
463 2018.
- 464 [41] Angela Fan, Thibaut Lavril, Edouard Grave, Armand Joulin, and Sainbayar Sukhbaatar. Address-
465 ing some limitations of Transformers with feedback memory. *Preprint arXiv:2002.09402*, 2020.
- 466 [42] Ke Tran, Arianna Bisazza, and Christof Monz. The importance of being recurrent for modeling
467 hierarchical structure. In *Proc. Conf. on Empirical Methods in Natural Language Processing*
468 *(EMNLP)*, pages 4731–4736, Brussels, Belgium, October 2018.
- 469 [43] Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. Pretrained transformers as universal
470 computation engines. *Preprint arXiv:2103.05247*, 2021.
- 471 [44] Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich
472 Küttler, Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. Deepmind lab.
473 *Preprint arXiv:1612.03801*, 2016.
- 474 [45] Joel Z Leibo, Cyprien de Masson d’Autume, Daniel Zoran, David Amos, Charles Beattie,
475 Keith Anderson, Antonio García Castañeda, Manuel Sanchez, Simon Green, Audrunas Gruslys,
476 et al. Psychlab: a psychology laboratory for deep reinforcement learning agents. *Preprint*
477 *arXiv:1801.08116*, 2018.
- 478 [46] Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskiy,
479 Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the atari human bench-
480 mark. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 507–517, Virtual only, July 2020.
- 481 [47] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
482 Wierstra, and Martin Riedmiller. Playing Atari with deep reinforcement learning. In *NIPS Deep*
483 *Learning Workshop*, Lake Tahoe, NV, USA, December 2013.
- 484 [48] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
485 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
486 Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- 487 [49] Lasse Espeholt, Hubert Soyer, Rémi Munos, Karen Simonyan, Volodymyr Mnih, Tom Ward,
488 Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu.
489 IMPALA: scalable distributed deep-RL with importance weighted actor-learner architectures.
490 In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 1406–1415, Stockholm, Sweden, July
491 2018.

- 492 [50] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap,
493 Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforce-
494 ment learning. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 1928–1937, New York
495 City, NY, USA, June 2016.
- 496 [51] Heinrich Küttler, Nantas Nardelli, Thibaut Lavril, Marco Selvatici, Viswanath Sivakumar, Tim
497 Rocktäschel, and Edward Grefenstette. Torchbeast: A PyTorch platform for distributed RL.
498 *Preprint arXiv:1910.03552*, 2019.
- 499 [52] Alexander Mott, Daniel Zoran, Mike Chrzanowski, Daan Wierstra, and Danilo Jimenez Rezende.
500 Towards interpretable reinforcement learning using attention augmented agents. In *Proc.*
501 *Advances in Neural Information Processing Systems (NeurIPS)*, pages 12329–12338, Vancouver,
502 Canada, December 2019.
- 503 [53] Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro
504 De Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, et al.
505 Massively parallel methods for deep reinforcement learning. In *Deep Learning Workshop,*
506 *International Conference on Machine Learning (ICML)*, Lille, France, July 2015.
- 507 [54] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Si-
508 mon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering
509 Atari, Go, Chess and Shogi by planning with a learned model. *Nature*, 588(7839):604–609,
510 2020.
- 511 [55] Felix A Gers, Jürgen Schmidhuber, and Fred Cummins. Learning to forget: Continual prediction
512 with LSTM. *Neural computation*, 12(10):2451–2471, 2000.

513 Checklist

- 514 1. For all authors...
- 515 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
516 contributions and scope? [\[Yes\]](#)
- 517 (b) Did you describe the limitations of your work? [\[Yes\]](#) See discussion on computational
518 cost in Sec. 4.1, as well as “Broader impact and limitations” paragraph in Conclusion.
- 519 (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See our
520 “Broader impact and limitations” paragraph in Conclusion.
- 521 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
522 them? [\[Yes\]](#)
- 523 2. If you are including theoretical results...
- 524 (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#) Our work is
525 empirical.
- 526 (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#) Our work is empirical.
- 527 3. If you ran experiments...
- 528 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
529 mental results (either in the supplemental material or as a URL)? [\[Yes\]](#) Please check
530 the supplemental material.
- 531 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
532 were chosen)? [\[Yes\]](#) Please check the supplemental material.
- 533 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
534 ments multiple times)? [\[Yes\]](#) We provide mean and std for experiments which are not
535 too expensive to be run multiple times.
- 536 (d) Did you include the total amount of compute and the type of resources used (e.g., type
537 of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) Please check the supplemental
538 material.
- 539 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 540 (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)

- 541 (b) Did you mention the license of the assets? [No] The license of the used assets is
542 discussed in the referenced documents.
- 543 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
544 We provide our code in the supplemental material.
- 545 (d) Did you discuss whether and how consent was obtained from people whose data you're
546 using/curating? [No] We only run on default benchmarks or synthetic data that was not
547 obtained from people.
- 548 (e) Did you discuss whether the data you are using/curating contains personally identifiable
549 information or offensive content? [No] This only applies to the Wikitext-103 dataset
550 which is a standard dataset composed of Wikipedia articles. As such it follows the
551 Wikipedia guidelines which do not contain personally identifiable information or
552 offensive content.
- 553 5. If you used crowdsourcing or conducted research with human subjects...
- 554 (a) Did you include the full text of instructions given to participants and screenshots, if
555 applicable? [N/A]
- 556 (b) Did you describe any potential participant risks, with links to Institutional Review
557 Board (IRB) approvals, if applicable? [N/A]
- 558 (c) Did you include the estimated hourly wage paid to participants and the total amount
559 spent on participant compensation? [N/A]