
Neurocoder: Learning General-Purpose Computation Using Stored Neural Programs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Artificial Neural Networks are functionally equivalent to special-purpose comput-
2 ers. Their inter-neuronal connection weights represent the learnt Neural Program
3 that instructs the networks on how to compute the data. However, without stor-
4 ing Neural Programs, they are restricted to only one, overwriting learnt programs
5 when trained on new data. Here we design Neurocoder, a new class of general-
6 purpose neural networks in which the neural network “codes” itself in a data-
7 responsive way by composing relevant programs from a set of shareable, modular
8 programs stored in external memory. For the first time, a Neural Program is ef-
9 ficiently treated as a datum in memory. Integrating Neurocoder into current neu-
10 ral architectures, we demonstrate new capacity to learn modular programs, reuse
11 simple programs to build complex ones, handle pattern shifts and remember old
12 programs as new ones are learnt, and show substantial performance improvement
13 in solving object recognition, playing video games and continual learning tasks.

14 1 Introduction

15 From its inception in 1943 until recently, the fundamental architectures of Artificial Neural Net-
16 works remained largely unchanged - a program is executed by passing data through a network of
17 artificial neurons whose inter-neuronal connection weights are learnt through training with data.
18 These inter-neuronal connection weights, or Neural Programs, correspond to a program in modern
19 computers [32]. Memory Augmented Neural Networks (MANN) are an innovative solution allow-
20 ing networks to access external memory for manipulating data [11, 12]. But they were still unable
21 to store Neural Programs in such external memory, and this severely limits machine learning. Stor-
22 ing inter-neuronal connection weights only in their network does not permit modular separation of
23 Neural programs and is analogous to a computer with one fixed program. Recent works introduce
24 *conditional computation* via adjusting or activating parts of a network in an input-dependent manner
25 [39, 33, 4, 13, 28], but networks remain monolithic. Current networks forget when retrained, old
26 inter-neuronal connection weights are merged with new ones or erased.

27 The brain is modular, not a monolithic system [8, 6]. Neuroscience research indicates that the brain is
28 divided into functional modules [19, 7, 9]. If the neural program for each module is kept in separate
29 networks, networks proliferate. Modular neural networks, another form of conditional computation,
30 combine the output of multiple expert networks, but as the experts grow, the networks grow drasti-
31 cally [20, 14, 35, 29]. This requires huge computational storage and introduces redundancy as these
32 experts do not share common basic programs.

33 *A pathway out of this bind is to keep such basic programs in memory and combine them as required.*
34 This brings neural networks towards modern general-purpose computers that use the stored-program
35 principle [37, 40] to efficiently access reusable programs in external memory. Here we show how
36 Neurocoder, a new neural framework, introduces a new class of general-purpose conditional compu-

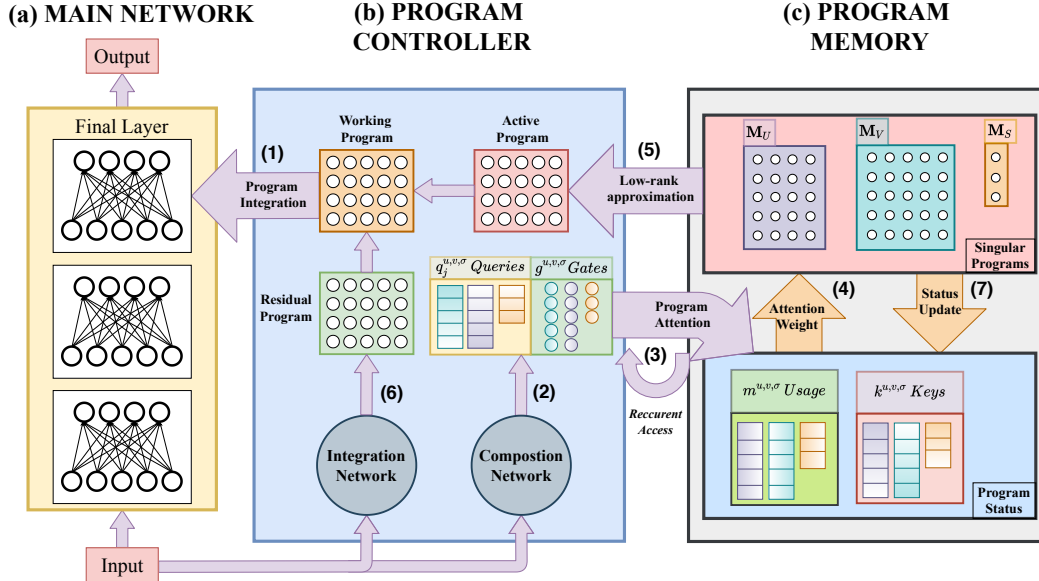


Figure 1: Neurocoder (a) The *Main Network* uses a *working program* to compute the output for the input. Here only the final layer of the Main Network is adaptively loaded with the working program (I). Other layers use traditional Neural Programs as connection weights (fixed-after-training). (b) The Program Controller’s *composition network* controls access to the Program Memory, emitting queries and interpolating gate control signals in response to the input (2). It then performs recurrent multi-head program attention to the Program Status (3), triggering attention weights to the Singular Programs (4). The attended Singular Programs form an *active program* using low-rank approximation (5). *Residual program* produced by the Program Controller’s *integration network* (6) plus the active program derives the *working program*. (c) The Program Memory stores the representations (singular programs) required to reconstruct the active program to be used by the Program Controller. Access is controlled through the Program Status including keys (k), and slot usage (m) that are updated during the training and computation (7).

37 tation machines in which a neural network can be “coded” in an input-dependent manner. Efficient
 38 decomposition of Neural Programs creates shareable modular components that can reconstruct the
 39 whole program space. These components change their “shapes” based on training and are stored
 40 in an external Program Memory. Then, in a data-responsive way, a Program Controller retrieves
 41 relevant components to build the Neural Program. This is analogous to shape-shifting Lego bricks
 42 that can be reused to build unlimited shapes and structures (See Appendix Fig. 4).

43 Using adaptive modular components vastly increases the learning capacity of the neural network
 44 by allowing re-utilisation of parameters, effectively curbing network growth as programs increase.
 45 More importantly, unlike pre-defined sub-networks or modules [20, 1] that combine at activation
 46 level, the construction of our modular components is dynamic and performed on the weight space.
 47 The Neural Program construction is learnt through training via traditional backpropagation [30] as
 48 the architecture is end-to-end differentiable.

49 2 Methods

50 2.1 System overview

51 A Neurocoder is a neural network (Main Network) coupled to an external Program Memory through
 52 a Program Controller. The *working program* of the Main Network processes the input data to pro-
 53 duce the output. This working program is “coded” by the Program Controller by creating an input-
 54 dependent *active program* from the Program Memory (Fig. 1). The following gives a high-level
 55 description of the Neurocoder framework and then the details.

56 **Neurocoder stores Singular Value Decomposition of Neural Programs in Program Memory**

57 The Neural Program needs to be stored efficiently in Program Memory. This is challenging as there
 58 may be millions of inter-neuronal connection weights, thus storing them directly ([22]) is grossly
 59 inefficient. Instead, the Neurocoder forms the basis of a subspace spanned by Neural Programs and
 60 stores the singular values and vectors of this subspace in memory slots of the Program Memory
 61 (hereafter referred to as *singular programs*). Based on the input, relevant singular programs are
 62 retrieved, a new program is reconstructed and then loaded in the Main Network to process the input.
 63 This representational choice significantly reduces the number of stored elements and allows each
 64 singular program to effectively represent a unitary function of the active program.

65 The *active program* matrix \mathbf{P} can be composed by standard low-rank approximation as

$$\mathbf{P} = \mathbf{U}\mathbf{S}\mathbf{V}^T = \sum_n^{r_m} \sigma_n u_n v_n^T \quad (1)$$

66 where \mathbf{U} and \mathbf{V} are matrices of the left and right singular vectors, and \mathbf{S} the matrix of singular values.
 67 r_m is the total number of components we want to retrieve. $\{\sigma_n\}_{n=1}^{r_m}$ is the attended singular values,
 68 $\{u_n\}_{n=1}^{r_m}$ and $\{v_n\}_{n=1}^{r_m}$ the attended singular vectors of \mathbf{S} , \mathbf{U} , and \mathbf{V} , respectively. The Program
 69 Memory is crafted as three *singular program memories* $\{\mathbf{M}_U, \mathbf{M}_V, \mathbf{M}_S\}$ —each of their memory
 70 slot stores a singular component or singular program. *The process “codes” the active program*
 71 *using singular programs from the program memories.* The coding is conditioned on input x_t , yet
 72 we drop index t for notation simplification and leave the details on the computation of σ_n, u_n, v_n in
 73 Sec. 2.2.

74 The Program Memory also maintains the status for each singular program in terms of access and
 75 usage. To access a singular program, *program keys* (k) are used. These keys are low-dimensional
 76 vectors that represent the singular program function and computed by a neural network that ef-
 77 fectively compresses the singular program. The *program usage* (m) measures memory utilisation,
 78 recording how much a memory slot is used in constructing a program. The components of the
 79 Program Memory are summarised in Fig. 1 (c).

80 **Recurrent multi-head program attention mechanisms for program storage and retrieval**

81 Neural networks use the concept of *differentiable attention* to access memory [11, 2]. This de-
 82 fines a weighting distribution over the memory slots essentially weighting the degree to which each
 83 memory slot participates in a read or write operation. This is unlike conventional computers that use
 84 a unique address to access a single memory slot.

85 Here we use two kinds of attention. First is *content-based attention* [11, 12] to ensure that the singu-
 86 lar program is selected based on its functionality and the data input. This is achieved by producing
 87 a query vector based on the input and comparing it to the program keys (k) using cosine similarity.
 88 Higher cosine similarity scores indicate higher attention weights to the singular programs associated
 89 with those program keys. Second, to encourage better memory utilisation, higher attention weights
 90 are assigned to slots with lower program usage (m) through *usage-based attention* [12, 31]. The
 91 attention weights from the two schemas are then combined using interpolating gates to compose the
 92 final attention weights to the Program Memory.

93 We adapt multi-head attention [11, 38] that applies multiple attentions in parallel to retrieve H singu-
 94 lar components. Besides, we introduce a recurrent attention mechanism, in which multi-head access
 95 is performed recurrently in J steps. The j -th set of H retrieved components is conditioned on the
 96 previous ones. This recurrent, multi-head attention allows the composition network to incrementally
 97 search for optimal components for building relevant active programs.

98 **Neurocoder learns to “code” a relevant working program via training**

99 The structure of the Program Memory and the role of the Program Controller facilitates the au-
 100 tomatic construction of working programs via training. The Program Controller controls memory
 101 access through its *composition network* that creates the *attention weight* defining how to weight the
 102 singular programs in the memories. A weighted summation of the singular programs results in the
 103 attended singular program. Applying the recurrent multi-head attention described earlier, multiple
 104 attended singular programs are retrieved to construct an active program (Eq. 1). Then the Program
 105 Controller generates a *residual program* using its *integration network*, adding to the active program

106 to produce the working program of the Main Network. This addition enables creation of flexible
 107 higher-rank working programs, which compensates for the low-rank coding process. The structure
 108 of the Program Controller is illustrated in Fig. 1 (b).

109 The singular programs are trained to represent unitary functions necessary for any computation
 110 whilst the composition and integration networks are trained to compose the relevant programs for
 111 the considering task. As such, beside minimising the task loss, we enforce orthogonality of stored
 112 singular vectors by minimising $\mathcal{L}_o = \mathbf{M}_U \mathbf{M}_U^\top - \mathbf{I} + \mathbf{M}_V \mathbf{M}_V^\top - \mathbf{I}$. The parameters of the networks,
 113 and the stored singular programs are adjusted using gradient training via minimising the total loss

$$\mathcal{L} = \mathcal{L}_{task} + a\mathcal{L}_o \quad (2)$$

114 where \mathcal{L}_{task} represents the supervised task loss and \mathcal{L}_o represents the orthogonal loss weighted by
 115 a hyper-parameter a to enforce orthogonality of the singular vectors.

116 2.2 Attention mechanisms for Program Memory

117 Here we describe program attention mechanisms used in this paper. Given $w_{in}^u, w_{in}^v, w_{in}^\sigma$ (jointly
 118 denoted as $w_{in}^{u,v,\sigma}$)—the attention weight to the i -th slot of the singular program memories $\mathbf{M}_U, \mathbf{M}_V$
 119 and \mathbf{M}_S , we retrieve the n -th singular vector as follows,

$$u_n = \sum_{i=1}^{P_u} w_{in}^u \mathbf{M}_U(i) \quad (3)$$

$$v_n = \sum_{i=1}^{P_v} w_{in}^v \mathbf{M}_V(i) \quad (4)$$

120 For the singular values, we need to enforce $\sigma_1 > \sigma_2 > \dots > \sigma_{r_m} > 0$, thus we retrieve using

$$\sigma_n = \begin{cases} \text{softplus} \left(\sum_{i=1}^{P_s} w_{in}^\sigma \mathbf{M}_S(i) \right) & n = r_m \\ \sigma_{n+1} + \text{softplus} \left(\sum_{i=1}^{P_s} w_{in}^\sigma \mathbf{M}_S(i) \right) & n < r_m \end{cases} \quad (5)$$

121 Here, P_u, P_v and P_s are the number of memory slots of $\mathbf{M}_U, \mathbf{M}_V$ and \mathbf{M}_S , respectively. In this
 122 paper, we set $P = P_u = P_v = P_s$ as the number of memory slots of the Program Memory. We note
 123 that these notations are specified for some data input x_t and the index n later maps to an attention
 124 head h , and an attention step j , hence the full notation should be $w_{tijh}^{u,v,\sigma}$. To simplify notations, we
 125 will drop u, v, σ from now and describe the computation of a representative w_{tijh} for any of the
 126 three program memories in the following parts.

127 Recurrent Access to the Program Memory via the composition network

128 To perform program attention, the Program Controller employs a composition network (denoted
 129 as f_θ), which takes the current input x_t and produce *program composition control signals* (ξ_t^p).
 130 If f_θ performs all attentions concurrently via multi-head attention as in [11, 38], it may lead to
 131 program collapse [22]. To have a better control of the component formation and alleviate program
 132 collapse, we propose to recurrently attend to the program memory. To this end, we implement f_θ as
 133 a recurrent neural network (LSTM [16]) and let it access the program memory J times, resulting in
 134 $\xi_t^p = \{\xi_{tj}^p\}_{j=1}^J$. At access step j , the recurrent network updates its hidden states and generates ξ_{tj}^p
 135 using recurrent dynamics as

$$\xi_{tj}^p, h_j = f_\theta(x_t, h_{j-1}) \quad (6)$$

136 where h_0 is initialized as zeros and ξ_{tj}^p is the program composition control signal at step j that
 137 depends on both on the input data x_t and the the previous state h_{j-1} . Particularly, the control signal
 138 contains the queries and the interpolation gates for each head to compute the program attention
 139 weight: $\xi_{tj}^p = \{q_{tijh}, g_{tijh}\}_{h=1}^H$. Here, at each attention step, we perform multi-head attention with
 140 H as the number of attention heads and thus, each ξ_{tj}^p consists of H pairs of queries and gates.
 141 Hence, the total number of retrieved components $r_m = J \times H$ and the index $n = j \times H + h$.

142 **Attending to Programs by “Name”**

143 Inspired by the content-based attention mechanism for data memory [11], we use the query to look
 144 for the singular programs. In computer programming, to find the appropriate program for some
 145 computation, we often refer to the program description or at least the name of the program. Here, we
 146 create the “name” for our neural programs by compressing the program content to a low-dimensional
 147 key vector. As such, we employ a neural network (f_φ) to compute the program memory keys as

$$k_i = f_\varphi(\mathbf{M}(i)) \quad (7)$$

148 where $k_i \in \mathbb{R}^K$ and i is the row index of the program memory. Here, f_φ learns to compress each
 149 memory slot into a K -dimensional vector. As the singular programs evolve, their keys get updated.
 150 In this paper, we update the program keys after each learning iteration during training.

151 Finally the content-based program memory attention c_{tjih} is computed using cosine distance be-
 152 tween the program keys k_i and the queries q_{tjh} as

$$c_{tjih} = \text{softmax}^{(i)} \left(\frac{q_{tjh} \cdot k_i}{\|q_{tjh}\| \cdot \|k_i\|} \right) \quad (8)$$

153 **Making Every Program Count**

154 Similarly to [12, 31], in addition to the content-based attention, we employ a least-used reading
 155 strategy to encourage the Program Controller to assign different singular programs to different com-
 156 ponents. In particular, we calculate the memory usage for each program slot across attentions as
 157

$$m_{tjih} = \max_{\tilde{j} \leq j} (w_{ti\tilde{j}h}) \quad (9)$$

158 Since we want to consider only l_I amongst P memory slots that have smallest usages, let $\hat{m}_{tjih}^{l_I}$
 159 denote the value of the l_I -th smallest usage, then the least-used attention is computed as

$$l_{tjih} = \begin{cases} \max_i (m_{tjih}) - m_{tjih} & ; m_{tjih} \leq \hat{m}_{tjih}^{l_I} \\ 0 & ; m_{tjih} > \hat{m}_{tjih}^{l_I} \end{cases} \quad (10)$$

160 The final program memory attention is computed as

$$w_{tjih} = \text{sigmoid}(g_{tjih}) c_{tjih} + (1 - \text{sigmoid}(g_{tjih})) l_{tjih} \quad (11)$$

161 Since the usage record are computed along the memory accesses, the multi-step Neurocoder utilises
 162 this attention mechanism better than the single-step Neurocoder, creating different attention styles
 163 (see Sec. 3.2). The composition the active program \mathbf{P}_t is illustrated in Appendix’s Fig. 5.

164 **2.3 Program Integration via the integration network**

165 Since the working program \mathbf{P}_t only contains top r_m principal components, it is low-rank and may
 166 be not flexible enough for sophisticated computation. We propose to enhance \mathbf{P}_t with a residual
 167 program \mathbf{R} — a traditional connection weight trained as the integration network’s parameters, which
 168 is constant after training w.r.t t . The residual program represents the sum of the remaining less
 169 important components. To this end, we suppress \mathbf{R} with a multiplier that is smaller than σ_{tr_m} — the
 170 smallest singular value of the main components - resulting in the integration formula

$$W_t = \mathbf{P}_t + w_t^r \sigma_{tr_m} \mathbf{R} \quad (12)$$

171 where $w_t^r = \text{sigmoid}(f_\phi(x_t))$ is an adaptive gating value that controls the contribution of the
 172 residual program. f_ϕ is the integration network in the Program Controller and hence, in our imple-
 173 mentation, the integration control signal sent by the Program Controller is $\lambda_t^p = \{w_t^r, \sigma_{tr_m}\}$. We
 174 note that in our experiments, the program integration can be disabled (W_t is directly set to \mathbf{P}_t) to
 175 prove the contribution of \mathbf{P}_t or reduce the number of parameters. The working program W_t is then
 176 used by the Main Network to execute the input data x_t (see (Fig. 1 (a))). For example, with linear

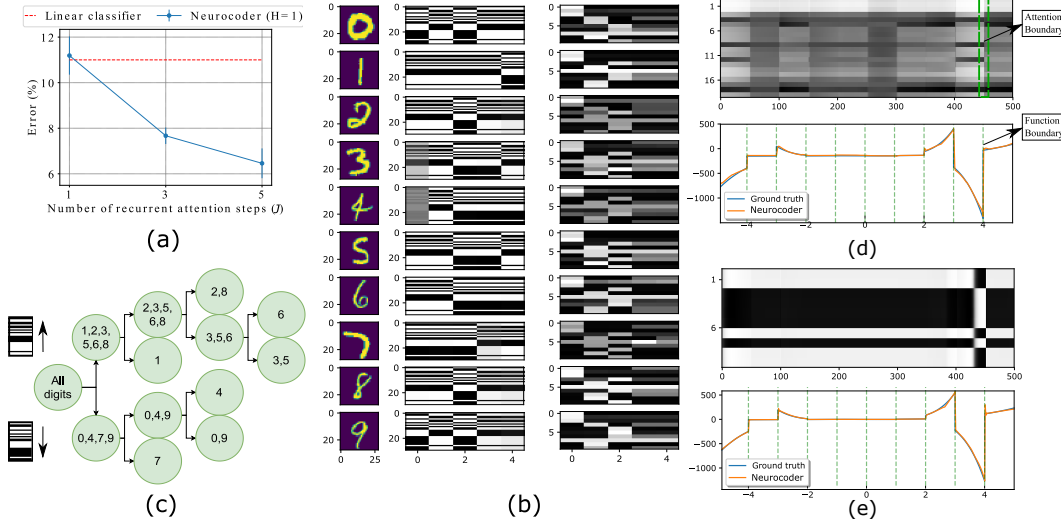


Figure 2: **(a)** MNIST test set classification error vs the number of steps (J) in Neurocoder (blue), compared with a linear classifier (red). **(b)** *1st column*: Digit images; *Middle column*: Single-step attention weights for 30 slots in \mathbf{M}_U (vertical axis) for first 3 singular vectors (horizontal axis) for each digit; *Last column*: Multi-step attention weights for 10 slots in \mathbf{M}_U (vertical axis) for first 3 singular vectors (horizontal axis). Multi-step attention is able to produce far more diverse patterns with fewer slots - 10 slots compared to single-step 30 slots. **(c)** Two attention patterns of single-step Neurocoder. The binary decision tree derived from single-step Neurocoder’s attention patterns. The two patterns across components represent the decisions going up and down across the binary tree. Visualisation for **(d)** multi-step ($J = 5$, 20 memory slots) and **(e)** single-step ($J = 1$, 10 memory slots) cases showing while processing a sequence of the polynomial auto-regression task. The Neurocoder’s attentions to \mathbf{M}_U that form the first component of the active program are shown over sequence timesteps (*upper*) with Neurocoder’s y_t prediction (orange) and ground truth (blue) (*lower*). The vertical dash green lines separate polynomial chunks. Each chunk represents a local pattern, and thus ideally requires a specific active program to compute the input x_t . Although both predict well, only the multi-step Neurocoder discovers the chunk boundaries, assigning program attention to the first component in accordance with sequence changes.

177 classifier Main Network, the execution is $y_t = x_t W_t$. Appendix’s Table 2 summarises the notations
 178 used for important parameters of Neurocoder.

179 3 Results

180 To demonstrate the flexibility of Neurocoder framework, we consider different learning paradigms:
 181 instance-based, sequential, multi-task and continual learning. We do not focus on breaking perfor-
 182 mance records by augmenting state-of-the-art models with Neurocoder. Rather our inquiry is on
 183 re-coding feed-forward layers with the Neurocoder’s programs and testing on varied data types to
 184 demonstrate its intrinsic properties. For some experiments, we include ablation studies.

185 We compare the performance of diverse Main Networks (MN) with and without Neurocoder. We
 186 also augment the Main Networks with other recent conditional computing methods, either modular
 187 (sparse Mixture of Experts, Neural Stored-program Memory) or monolithic (HyperNets, FiLM) to
 188 form stronger baselines across our experiments. In our experiments, we always apply Neurocoder
 189 to all layers of multi-layer perceptrons (MLP) or just the final feed-forward layer of deep CNN
 190 networks (LeNet, DenseNet, ResNet), RNNs (GRU, LSTM), MANN (NTM). Other competitors
 191 such as MOE, NSM, HyperNet and FiLM are applied to the Main Networks in the same manner.

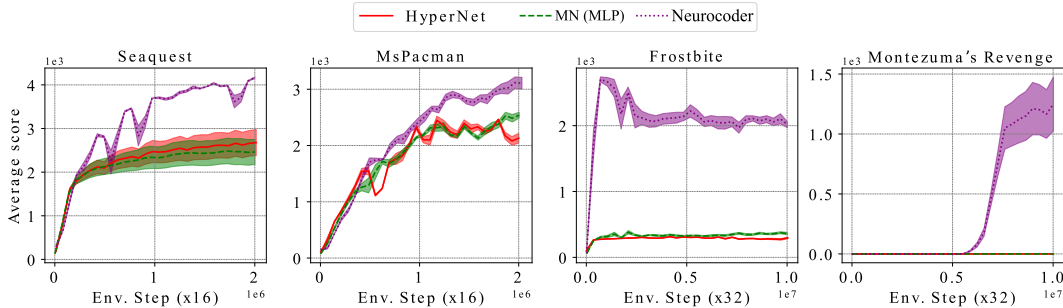


Figure 3: Learning curves (mean and std. over 5 runs) on representative Atari 2600 games. All baselines are applied to the actor/critic networks in the A3C agent.

192 3.1 Instance-based learning - Object Recognition

193 We tested Neurocoder on instance-based learning through classical image classification tasks using
 194 MNIST [24] and CIFAR [21] datasets. The first experiment interpreted Neurocoder’s behaviour
 195 in classifying digits into 10 classes (0 – 9) using linear classifier Main Network. With equivalent
 196 model size, Neurocoder using the novel recurrent attention surpasses the performance of the linear
 197 classifier [24] by up to 5% (Fig. 2 (a)).

198 To differentiate the input, Neurocoder attends to different components of the active program to
 199 guide the decision-making process. Fig. 2 (b) shows single-step and multi-step attention to the first
 200 3 singular vectors for each digit across memory slots. Multi-step attention produces richer patterns
 201 compared to single-step Neurocoder that manages only 2 attention weight patterns.

202 Fig. 2 (c) illustrates how Neurocoder performs modular learning by showing the attention assign-
 203 ment for top 3 singular vectors as a binary decision tree. Digits under the same parental node share
 204 similar attention paths, and thereby similar active programs. Some digits look unique (e.g. 7) result-
 205 ing in active programs composed of unique attention paths, discriminating themselves early in the
 206 decision tree. Some digits (e.g. 0 and 9) share the same attention pattern for the first 3 components
 207 and are thus unclassifiable. They can only be distinguished by considering more singular vectors.

208 We integrated Neurocoder with deep networks - 5-layer *LeNet* and 100-layer *DenseNet* - and tested
 209 on CIFAR datasets. Neurocoder significantly outperformed the original Main Networks with perfor-
 210 mance gain 1 – 5%. Compared with recent conditional computing models such as sparse Mixture of
 211 Experts (MOE [35]) and Neural Stored-program Memory (NSM [22]), Neurocoder required a tenth
 212 of the number of parameters and performed better by up to 8 – 10% (see Appendix’s Table 3).

213 3.2 Sequential learning - Adaption to sequence changes and game playing using 214 reinforcement learning

215 Recurrent neural networks (RNN) can learn from sequential data by updating the hidden states of
 216 the networks. However, this does not suffice when local patterns shift, as is often the case. We now
 217 demonstrate that Neurocoder helps RNNs overcome this limitation by composing diverse programs
 218 to handle sequence changes.

219 **Synthetic polynomial auto-regression** We created a simple auto-regression task in which data
 220 points are sampled from polynomial function chunks that change over time. The Main Network is
 221 a strong *RNN-Gated Recurrent Unit (GRU)* [5]. We found that GRU integrated with a single-step
 222 or multi-step Neurocoder converged much faster than all other baselines. The other conditional
 223 computing counterparts (HyperNet [13], FiLM [28]) adapt by re-scaling weights or activation of the
 224 GRU, which were shown inferior to our modular approach (Appendix’s Fig. 6).

225 Visualising the first singular vector attention weights in \mathbf{M}_U , we find that the multi-step attention
 226 Neurocoder changes its attention following polynomial changes - it attends to the same singular pro-
 227 gram when processing data from the same polynomial and alters attention for data from a different
 228 polynomial (Fig. 2(d)). In contrast, the single-step Neurocoder only changes its attention when
 229 there is a remarkable change in y -coordinate values (Fig. 2(e)). Although single-step Neurocoder

Method	MN (MLP [17])	MN (MLP ours)	NSM	Neurocoder
Adam	55.16±1.38	53.55±1.27	54.85±2.81	58.46±0.46
Adagrad	58.08±1.06	57.83±2.74	58.42±1.87	62.28±4.03
L2	66.00±3.73	64.37±2.40	62.83±7.21	69.89±1.72
SI	64.76±3.09	64.41±3.36	64.36±2.99	67.96±3.22
EWC	58.85±2.59	58.41±2.37	58.12±3.24	65.66±1.25
O-EWC	57.33±1.44	57.78±1.84	58.55±3.40	73.97±1.50

Table 1: Incremental domain continual learning with Split MNIST. Final test accuracy (mean and std.) over 10 runs.

230 converges well, it did not discover the underlying structure of the data, and thus underperformed
 231 the multi-step Neurocoder. We hypothesise that when recurrence is employed, usage-based atten-
 232 tion takes effect, stipulating better memory utilisation and diverse attentions over timesteps. We ran
 233 multi-step Neurocoder without usage-based attention. The results were worse than the full multi-
 234 step Neurocoder, which confirms our hypothesis (Appendix’s Fig. 6).

235 **Atari game reinforcement learning** We used reinforcement learning as a further testbed to show
 236 the ability to adapt to environmental changes. We performed experiments on several Atari 2600
 237 games [3] wherein the agent was implemented as the *Asynchronous Advantage Actor-Critic (A3C*
 238 *[26])*. In the Atari platform, agents are allowed to observe the screen snapshot of the games and act
 239 to earn the highest score. We augmented the A3C by employing Neurocoder’s working programs
 240 for feed-forward layers of the actor and critic networks, aiming to decompose the policy and value
 241 function into singular programs that were selected depending on the game state.

242 *Frostbite and Montezuma’s Revenge*. These games are known to be challenging for A3C and other
 243 algorithms [26]. We trained A3C and HyperNet-based A3C for over 300 million steps, yet these
 244 models did not show any sign of learning, performing equivalently to random agents. For such com-
 245 plicated environments with sparse rewards, both the monolithic neural networks and the HyperNet’s
 246 unstored fast-weights fail to learn (almost zero scores). In contrast, Neurocoder enabled A3C to
 247 achieve from 1,500 to 3,000 scores on these environments (Fig. 3), confirming the importance of
 248 decomposing a complex solution to smaller, simple stored programs.

249 3.3 Multi-task learning - Solving multiple algorithms simultaneously

250 Here we explore the modular learning capability of Neurocoder in multi-task setting. Inspired by al-
 251 gorithmic sequencing tasks [22], we created a challenging sequential multi-task benchmark wherein
 252 the input sequence is a series of sub-sequences from 4 algorithms: Copy, Repeat Copy, Associative
 253 Recall and Priority Sort [11]. Each sub-sequence, following a task identification vector, represents
 254 the input for each task. In each input sequence, n tasks were sampled from the set of 4 algorithms
 255 randomly with replacement and the output sequences were created correspondingly.

256 We trained a *MANN-Neural Turing Machine (NTM [11])* Main Network with FiLM, HyperNet and
 257 our Neurocoder augmentation on sequences of $n = 4$ tasks, and tested with sequences of $n = 4$ and
 258 $n = 8$ tasks. Appendix’s Fig. 7 demonstrates that Neurocoder was performant in both test settings,
 259 not only achieving lowest error on $n = 4$, but also being the only one generalised well to $n = 8$
 260 scenario, which was unseen during training.

261 3.4 Continual learning - Learning tasks sequentially without catastrophic forgetting

262 In continual learning, standard neural networks often suffer from “catastrophic forgetting” in which
 263 they cannot retain knowledge acquired from old tasks upon learning new ones [10]. Our Neurocoder
 264 offers natural mitigation of such catastrophic forgetting in neural networks by attending to different
 265 singular programs whilst learning different tasks.

266 In this case, in addition to the Main Network, we examine several continual learning algorithms
 267 with and without Neurocoder. These algorithms, including Elastic Weight Consolidation (EWC
 268 [41]) and Synaptic Intelligence (SI [41]), work by regularising the loss function and thus can be
 269 easily combined with Neurocoder by modifying the loss \mathcal{L}_{task} . We demonstrate that Neurocoder

270 can improve these continual learning algorithms without requiring additional assumptions as in other
271 approaches [25, 36, 34] that either utilise task embedding or replay memory.

272 **Split MNIST** We first considered the split MNIST dataset—a standard continual learning bench-
273 mark wherein the original MNIST was split into a 5 2-way classification tasks, consecutively pre-
274 sented to a *Multi-layer Perceptron* Main Network (MLP). We followed the benchmarking as in [17]
275 in which various optimisers and state-of-the-art continual learning methods were examined under in-
276 cremental task and domain scenarios. We measured the performance of the MLP versus Neurocoder
277 and NSM under each continual learning method. In both scenarios, Neurocoder was compatible
278 with all continual learning methods, demonstrating superior performance over MLP and NSM with
279 performance gain between 1 to 16% (see Appendix’s Table 5 and 1).

280 **Split CIFAR** We verified the scalability of Neurocoder to more challenging datasets. We split
281 CIFAR datasets as in the split MNIST, resulting in 5-task 2-way split CIFAR10 and a 20-task 5-way
282 split CIFAR100. We used Main Network *ResNet* [15]—a very deep CNN architecture.

283 When we stressed the orthogonal loss ($a = 10$) and used bigger program memory (100 slots), Neu-
284 rocoder improved ResNet classification by 15% and 10% on CIFAR10 and CIFAR100, respectively.
285 When we integrated Neurocoder with Synaptic Intelligence (SI [41]), the performance was further
286 improved, maintaining a stable performance above 80% accuracy for CIFAR10 and outperforming
287 using SI alone by 10% for CIFAR100 (see Appendix’s Fig. 8).

288 4 Discussion

289 Our experiments demonstrate that Neurocoder is capable of re-coding Neural Programs in distinctive
290 neural networks, amplifying their capabilities in diverse learning scenarios: instance-based, sequen-
291 tial, multi-task and continual learning. This consistently results in significant performance increase,
292 and further creates novel robustness to pattern shift and catastrophic forgetting. This ability for each
293 architecture to re-code itself is made possible without changing the way it is trained, or majorly
294 increasing the number of parameters it needs to learn (see Appendix Table 7).

295 The MNIST problem illustrates the reasoning process of Neurocoder when classifying digit images
296 wherein its singular program assignment resembles a binary tree decision-making process - it shows
297 how some singular programs are shared, others are not. The polynomial auto-regression problem
298 highlights the importance of efficient memory utilisation in re-constructing the working program
299 enabling discovery of hidden structures in sequential data. Training our framework with reinforce-
300 ment learning, we enable neural agents to solve complex games wherein traditional methods fail
301 or learn slowly. Neurocoder also works well with multi-task setting, as shown in the challenging
302 multi-algorithm benchmark. Finally, continual learning problems show that Neurocoder mitigates
303 catastrophic forgetting efficiently under different learning settings/algorithms.

304 Our solution offers a single framework that is scalable and adaptable to various problems and learn-
305 ing paradigms. Unlike previous attempts to employ a bank of separate big programs [20, 35, 22],
306 Neurocoder maintains only shareable, smaller components that can reconstruct the whole program
307 space, thereby heavily utilising the parameters and preventing the model from proliferating. We
308 note that Neurocoder is orthogonal to approaches employing tensor decomposition to reduce the
309 number of parameters or hasten the computation [27, 23]. Neurocoder composes rather than decom-
310 pose the neural weights. Our aim is not only to enable efficient parameter usage, but also achieve
311 general-purpose computing power, outperforming other methods in numerous learning problems.

312 One limitation of this work is the number of additional hyperparameters, which prevents us from
313 fully tuning Neurocoder. Our research aims to add new capabilities to current neural networks to
314 improve their performance and make them robust in different learning scenario. Hence, we do
315 not see any intermediate negative societal impact. In future work, we will extend Neurocoder’s
316 application beyond feed-forward layers. It would be interesting to efficiently replace all neural layers
317 including CNN or Transformer by Neurocoder’s programs. We can also further extend Neurocoder’s
318 ability by allowing a growing Program Memory, in which the model decides to add or erase memory
319 slots as the number of data patterns grows or shrinks beyond the current program space’s capacity.
320 Such a system represents a more flexible general-purpose computer that can dynamically allocate
321 computing resources by itself without human pre-specification.

References

- 322
- 323 [1] Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Neural module networks. In
324 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 39–48,
325 2016.
- 326 [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by
327 jointly learning to align and translate. In *International Conference on Learning Representations*, 2015.
328
- 329 [3] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning
330 environment: An evaluation platform for general agents. *Journal of Artificial Intelligence*
331 *Research*, 47:253–279, 2013.
- 332 [4] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradi-
333 ents through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*,
334 2013.
- 335 [5] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares,
336 Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-
337 decoder for statistical machine translation. In *Conference on Empirical Methods in Natural*
338 *Language Processing (EMNLP)*, pages 1724–1734. Association for Computational Linguistics,
339 October 2014.
- 340 [6] JC Eccles. The modular operation of the cerebral neocortex considered as the material basis of
341 mental events. *Neuroscience*, 6(10):1839–1855, 1981.
- 342 [7] Gerald M Edelman. Neural darwinism: selection and reentrant signaling in higher brain func-
343 tion. *Neuron*, 10(2):115–125, 1993.
- 344 [8] Gerald M Edelman and Vernon B Mountcastle. *The mindful brain: cortical organization and*
345 *the group-selective theory of higher brain function*. Massachusetts Inst of Technology Pr, 1978.
- 346 [9] Richard SJ Frackowiak. *Human brain function*. Elsevier, 2004.
- 347 [10] Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive*
348 *sciences*, 3(4):128–135, 1999.
- 349 [11] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint*
350 *arXiv:1410.5401*, 2014.
- 351 [12] Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka
352 Grabska-Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John
353 Agapiou, et al. Hybrid computing using a neural network with dynamic external memory.
354 *Nature*, 538(7626):471–476, 2016.
- 355 [13] David Ha, Andrew M. Dai, and Quoc V. Le. Hypernetworks. In *International Conference on*
356 *Learning Representations*, 2017.
- 357 [14] Bart LM Happel and Jacob MJ Murre. Design and evolution of modular neural network archi-
358 tectures. *Neural networks*, 7(6-7):985–1004, 1994.
- 359 [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
360 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recogni-*
361 *tion*, pages 770–778, 2016.
- 362 [16] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*,
363 9(8):1735–1780, 1997.
- 364 [17] Yen-Chang Hsu, Yen-Cheng Liu, Anita Ramasamy, and Zsolt Kira. Re-evaluating continual
365 learning scenarios: A categorization and case for strong baselines. In *NeurIPS Continual*
366 *learning Workshop*, 2018.

- 367 [18] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely con-
368 nected convolutional networks. In *Proceedings of the IEEE conference on computer vision and*
369 *pattern recognition*, pages 4700–4708, 2017.
- 370 [19] D.H. Hubel. *Eye, Brain, and Vision*. Scientific American Library series. Scientific American
371 Library, 1988.
- 372 [20] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mix-
373 tures of local experts. *Neural computation*, 3(1):79–87, 1991.
- 374 [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
375 *TR-2009*, 2009.
- 376 [22] Hung Le, Truyen Tran, and Svetha Venkatesh. Neural stored-program memory. In *Internat-*
377 *ional Conference on Learning Representations*, 2020.
- 378 [23] V Lebedev, Y Ganin, M Rakhuba, I Oseledets, and V Lempitsky. Speeding-up convolutional
379 neural networks using fine-tuned cp-decomposition. In *International Conference on Learning*
380 *Representations*, 2015.
- 381 [24] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning
382 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 383 [25] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learn-
384 ing. In *Advances in neural information processing systems*, pages 6467–6476, 2017.
- 385 [26] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap,
386 Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforce-
387 ment learning. In *International conference on machine learning*, pages 1928–1937, 2016.
- 388 [27] Alexander Novikov, Dmitrii Podoprikin, Anton Osokin, and Dmitry P Vetrov. Tensorizing
389 neural networks. In *Advances in neural information processing systems*, pages 442–450, 2015.
- 390 [28] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. FiLM:
391 Visual Reasoning with a General Conditioning Layer. In *AAAI Conference on Artificial Intel-*
392 *ligence*, New Orleans, United States, February 2018.
- 393 [29] Clemens Rosenbaum, Tim Klinger, and Matthew Riemer. Routing networks: Adaptive selec-
394 tion of non-linear functions for multi-task learning. In *International Conference on Learning*
395 *Representations*, 2018.
- 396 [30] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by
397 back-propagating errors. *Nature*, 323(6088):533–536, 1986.
- 398 [31] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap.
399 Meta-learning with memory-augmented neural networks. In *International conference on ma-*
400 *chine learning*, pages 1842–1850, 2016.
- 401 [32] Jürgen Schmidhuber. Making the world differentiable: On using self-supervised fully recurrent
402 neural networks for dynamic reinforcement learning and planning in non-stationary environm-
403 nts. *TR FKI-126-90*, 1990.
- 404 [33] Jürgen Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic
405 recurrent networks. *Neural Computation*, 4(1):131–139, 1992.
- 406 [34] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic
407 forgetting with hard attention to the task. In *International Conference on Machine Learning*,
408 pages 4548–4557, 2018.
- 409 [35] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey E.
410 Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-
411 experts layer. In *International Conference on Learning Representations*, 2017.

- 412 [36] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep
413 generative replay. In *Advances in Neural Information Processing Systems*, pages 2990–2999,
414 2017.
- 415 [37] A.M Turing. On computable numbers, with an application to the entscheidungsproblem. In
416 *Proceedings of the London Mathematical Society*, 1936.
- 417 [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
418 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural informa-*
419 *tion processing systems*, pages 5998–6008, 2017.
- 420 [39] Christoph von der Malsburg. The correlation theory of brain function, 1981.
- 421 [40] John Von Neumann. First draft of a report on the edvac. *IEEE Annals of the History of*
422 *Computing*, 15(4):27–75, 1993.
- 423 [41] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intel-
424 ligence. *Proceedings of machine learning research*, 70:3987, 2017.

425

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) See Discussion and Appendix's "Training procedure and hyper-parameter selections."
 - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See Discussion.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[No\]](#) Data is public, provided with link. Code will be available after published. All training details are available and can be used to implement and reproduce the results.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See Appendix's "Training procedure and hyper-parameter selections."
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) See Experimental Results.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Appendix's "Training procedure and hyper-parameter selections."
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)
 - (b) Did you mention the license of the assets? [\[No\]](#) All assets are public. We will mention the license detail after the paper is published.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[No\]](#)
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [\[No\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[No\]](#)
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)