
Metalearning to Continually Learn In Context

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 General-purpose learning systems should improve themselves in open-ended fash-
2 ion in ever-changing environments. Conventional learning algorithms for neural
3 networks, however, suffer from catastrophic forgetting (CF)—previously acquired
4 skills are forgotten when a new task is learned. Instead of hand-crafting new
5 algorithms for avoiding CF, we propose Automated Continual Learning (ACL) to
6 train self-referential neural networks to meta-learn their own *in-context* continual
7 (meta-)learning algorithms. ACL encodes continual learning desiderata—good
8 performance on both old and new tasks—into its meta-learning objectives. Our
9 experiments demonstrate that, in general, in-context learning algorithms also suffer
10 from CF but ACL effectively solves such “in-context catastrophic forgetting”. Our
11 ACL-learned algorithms outperform hand-crafted ones and popular meta-continual
12 learning methods on the Split-MNIST benchmark in the replay-free setting, and
13 enables continual learning of diverse tasks consisting of multiple few-shot and stan-
14 dard image classification datasets. Going beyond, we also highlight the limitations
15 of in-context continual learning, by investigating the possibilities to extend ACL to
16 the realm of state-of-the-art CL methods which leverage pre-trained models.¹

17 1 Introduction

18 Enemies of memories are other memories [1]. Continually-learning artificial neural networks (NNs)
19 are memory systems in which their *weights* store memories of task-solving skills or programs, and
20 their *learning algorithm* is responsible for memory read/write operations. Conventional learning
21 algorithms—used to train NNs in the standard scenarios where all training data is available *at once*—
22 are known to be inadequate for continual learning (CL) of multiple tasks where data for each task
23 is available *sequentially and exclusively*, one at a time. They suffer from “catastrophic forgetting”
24 (CF; [2–5]); the NNs forget, or rather, the learning algorithm erases, previously acquired skills, in
25 exchange of learning to solve a new task. Naturally, a certain degree of forgetting is unavoidable
26 when the memory capacity is limited, and the amount of things to remember exceeds such an upper
27 bound. In general, however, capacity is not the fundamental cause of CF; typically, the same NNs,
28 suffering from CF when trained on two tasks sequentially, can perform well on both tasks when they
29 are jointly trained on the two tasks at once instead (see, e.g., [6]).

30 The real root of CF lies in the learning algorithm as a memory mechanism. A “good” CL algorithm
31 should preserve previously acquired knowledge while also leveraging previous learning experiences
32 to improve future learning, by maximally exploiting the limited memory space of model parameters.
33 All of this is the *decision-making problem of learning algorithms*. In fact, we can not blame the
34 conventional learning algorithms for causing CF, since they are not aware of such a problem. They
35 are designed to train NNs for a given task at hand; they treat each learning experience independently
36 (they are stationary up to certain momentum parameters in certain optimizers), and ignore any

¹Here we’ll add a link to our public GitHub code repository.

37 potential influence of current learning on past or future learning experiences. Effectively, more
38 sophisticated algorithms previously proposed against CF [7, 8], such as elastic weight consolidation
39 [9, 10] or synaptic intelligence [11], often introduce manually-designed constraints as regularization
40 terms to explicitly penalize current learning for deteriorating knowledge acquired in past learning.

41 Here, instead of hand-crafting learning algorithms for continual learning, we train self-referential
42 neural networks [12, 13] to meta-learn their own “in-context” continual learning algorithms. We
43 train them through gradient descent on learning objectives that reflect desiderata for continual learn-
44 ing algorithms—good performance on both old and new tasks, including forward and backward
45 transfer. In fact, by extending the standard settings of few-shot or meta-learning based on sequence-
46 processing NNs [14–18], the continual learning problem can also be formulated as a long-span
47 sequence processing task [19]. Corresponding CL sequences can be obtained by concatenating multi-
48 ple few-shot/meta-learning sub-sequences, where each sub-sequence consists of input/target examples
49 corresponding to the task to be learned in-context. As we’ll see in Sec. 3, this setting also allows us
50 to seamlessly express classic desiderata for CL as part of objective functions of the meta-learner.

51 Once formulated as such a sequence-learning task, we let gradient descent search for CL algorithms
52 achieving the desired CL behaviors in the program space of NN weights. In principle, all typical
53 challenges of CL—such as the stability-plasticity dilemma [20]—are automatically discovered and
54 handled by the gradient-based program search process. Once trained, CL is automated through
55 recursive self-modification dynamics of the trained NN, without requiring any human intervention
56 such as adding extra regularization or setting hyper-parameters for CL. Therefore, we call our
57 method, Automated Continual Learning (ACL).

58 Our experiments focus on supervised image classification, making use of standard few-shot learning
59 datasets for meta-training, namely, Mini-ImageNet [21, 22], Omniglot [23], and FC100 [24], while
60 we also meta-test on other datasets including MNIST [25], FashionMNIST [26] and CIFAR-10 [27].

61 **Our core contribution** is a set of focused experiments showing various facets of in-context CL: (1)
62 We first reveal the “in-context catastrophic forgetting” problem using two-task settings (Sec. 4.1) and
63 analyse its emergence (Sec. 4.2). We are not aware of any prior work discussing this problem. (2)
64 We show very promising results of our ACL-trained learning algorithm on the classic Split-MNIST
65 [6, 28] benchmark, outperforming hand-crafted learning algorithms and prior meta-continual learning
66 methods [29–31]. (3) We experimentally illustrate the limitations of ACL on 5-datasets [32] and
67 Split-CIFAR100 by comparing to more recent prompt-based state-of-the-art CL methods [33, 34].

68 2 Background

69 2.1 Continual Learning

70 The main focus of this work is on continual learning [35, 36] in *supervised* learning settings even
71 though high-level principles we discuss here also transfer to reinforcement learning settings [37].
72 In addition, we focus on the realm of CL methods that keep model sizes constant (unlike certain
73 CL methods that incrementally add more parameters as more tasks are presented; see, e.g., [38]),
74 and do not make use of any external replay memory (used in other CL methods; see, e.g., [39–43]).

75 Classic desiderata for a CL system (see, e.g., [44, 45]) are typically summarized as good performance
76 on three metrics: *classification accuracies* on each dataset (their average), *backward transfer* (i.e., im-
77 pact of learning a new task on the model’s performance on previous tasks; e.g., catastrophic forgetting
78 is a negative backward transfer), and *forward transfer* (impact of learning a task for the model’s per-
79 formance on a future task). From a broader perspective of meta-learning systems, we may also measure
80 other effects such as *learning acceleration* (i.e., whether the system leverages previous learning ex-
81 periences to accelerate future learning); here our primary focus remains the classic CL metrics above.

82 2.2 Few-shot/meta-learning via Sequence Learning

83 In Sec. 3, we’ll formulate continual learning as a long-span sequence processing task. This is a direct
84 extension of the classic few-shot/meta learning formulated as a sequence learning problem. In fact,
85 since the seminal works [14–17] (see also [46]), many sequence processing neural networks (see,
86 e.g., [47–58] including Transformers [59, 18]) have been trained as a meta-learner [13, 12] that learn
87 by observing sequences of training examples (i.e., pairs of inputs and their labels) in-context.

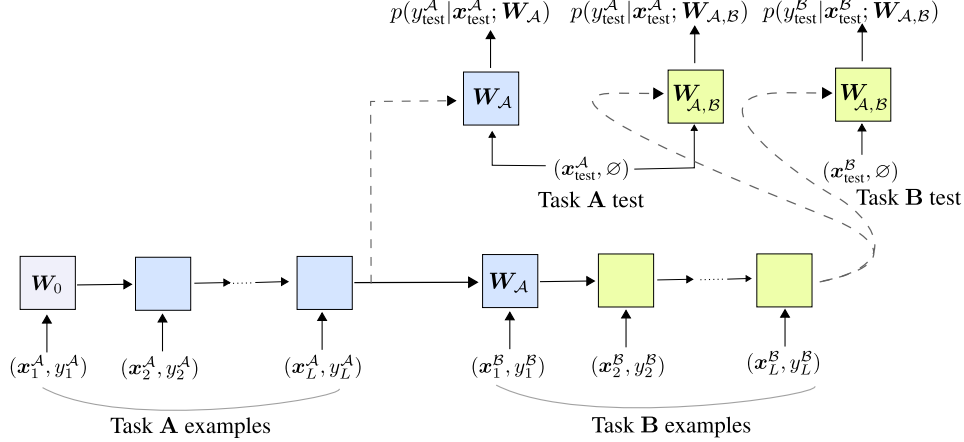


Figure 1: An illustration of meta-training in Automated Continual Learning (ACL) for a self-referential/modifying weight matrix \mathbf{W}_0 . Weights \mathbf{W}_A obtained by observing examples for Task A (blue) are used to predict a test example for Task A. Weights $\mathbf{W}_{A,B}$ obtained by observing examples for Task A then those for Task B (yellow) are used to predict a test example for Task A (backward transfer) as well as a test example for Task B (forward transfer).

88 Here we briefly review such a formulation. Let d, N, K, P be positive integers. In sequential
 89 N -way K -shot classification settings, a sequence processing NN with a parameter vector $\theta \in \mathbb{R}^P$
 90 observes a pair (\mathbf{x}_t, y_t) where $\mathbf{x}_t \in \mathbb{R}^d$ is the input and $y_t \in \{1, \dots, N\}$ is its label at each step
 91 $t \in \{1, \dots, N \cdot K\}$, corresponding to K examples for each one of N classes. After the presentation of
 92 these $N \cdot K$ examples (often called the *support set*), one extra input $\mathbf{x} \in \mathbb{R}^d$ (often called the *query*)
 93 is fed to the model without its true label but with an “unknown label” token \emptyset (number of input labels
 94 accepted by the model is thus $N + 1$). The model is trained to predict its true label, i.e., the parameters
 95 of the model θ are optimized to maximize the probability $p(y | (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{N \cdot K}, y_{N \cdot K}), (\mathbf{x}, \emptyset); \theta)$
 96 of the correct label $y \in \{1, \dots, N\}$ of the input query \mathbf{x} . Since class-to-label associations are
 97 randomized and unique to each sequence $((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{N \cdot K}, y_{N \cdot K}), (\mathbf{x}, \emptyset))$, each such a sequence
 98 represents a new (few-shot or meta) learning example to train the model. To be more specific, this
 99 is the *synchronous* label setting of Mishra et al. [18] where the learning phase (observing examples,
 100 (\mathbf{x}_1, y_1) etc.) is separated from the prediction phase (predicting label y given (\mathbf{x}, \emptyset)). We opt for
 101 this variant in our experiments as we empirically find this (at least in our specific settings) more
 102 stable than the *delayed* label setting [14] where the model has to make a prediction for every input,
 103 and the label is fed to the model with a delay of one time step.

104 2.3 Self-Referential Weight Matrices

105 Our method (Sec. 3) can be applied to any sequence-processing NN architectures in principle.
 106 Nevertheless, certain architectures naturally fit better to parameterize a self-improving continual
 107 learner. Here we use the *modern self-referential weight matrix* (SRWM; [19, 60]) to build a generic
 108 self-modifying NN. An SRWM is a weight matrix that sequentially modifies itself as a response
 109 to a stream of input observations [12, 61]. The modern SRWM belongs to the family of linear
 110 Transformers (LTs) a.k.a. Fast Weight Programmers (FWPs; [62–68]). Linear Transformers and
 111 FWPs are an important class of the now popular Transformers [59]: unlike the standard ones whose
 112 computational requirements grow quadratically and whose state size grows linearly with the context
 113 length, LTs/FWPs’ complexity is linear and the state size is constant w.r.t. sequence length (like
 114 in the standard RNNs). This is an important property for in-context CL, since, conceptually, we
 115 want such a CL system to continue to learn for an arbitrarily long, lifelong time span. Moreover,
 116 the duality between linear attention and FWPs [67]—and likewise, between linear attention and
 117 gradient descent-trained linear layers [69, 70]—have played a key role in certain theoretical analyses
 118 of in-context learning capabilities of Transformers [71, 72].

119 The dynamics of an SRWM [19] are described as follows. Let $d_{\text{in}}, d_{\text{out}}, t$ be positive integers, and \otimes
 120 denote outer product. At each time step t , an SRWM $\mathbf{W}_{t-1} \in \mathbb{R}^{(d_{\text{out}}+2*d_{\text{in}}+1) \times d_{\text{in}}}$ observes an input

121 $\mathbf{x}_t \in \mathbb{R}^{d_{\text{in}}}$, and outputs $\mathbf{y}_t \in \mathbb{R}^{d_{\text{out}}}$, while also updating itself to \mathbf{W}_t as:

$$[\mathbf{y}_t, \mathbf{k}_t, \mathbf{q}_t, \beta_t] = \mathbf{W}_{t-1} \mathbf{x}_t \quad (1)$$

$$\mathbf{v}_t = \mathbf{W}_{t-1} \phi(\mathbf{q}_t); \bar{\mathbf{v}}_t = \mathbf{W}_{t-1} \phi(\mathbf{k}_t) \quad (2)$$

$$\mathbf{W}_t = \mathbf{W}_{t-1} + \sigma(\beta_t)(\mathbf{v}_t - \bar{\mathbf{v}}_t) \otimes \phi(\mathbf{k}_t) \quad (3)$$

122 where $\mathbf{v}_t, \bar{\mathbf{v}}_t \in \mathbb{R}^{(d_{\text{out}}+2*d_{\text{in}}+1)}$ are value vectors, $\mathbf{q}_t \in \mathbb{R}^{d_{\text{in}}}$ and $\mathbf{k}_t \in \mathbb{R}^{d_{\text{in}}}$ are query and key vectors,
 123 and $\sigma(\beta_t) \in \mathbb{R}$ is the learning rate. σ and ϕ denote sigmoid and softmax functions respectively. ϕ
 124 is typically also applied to \mathbf{x}_t in Eq. 1; here we follow Irie et al. [19]’s few-shot image classification
 125 setting, and use the variant without it. Eq. 3 corresponds to a rank-one update of the SRWM, from
 126 \mathbf{W}_{t-1} to \mathbf{W}_t , through the *delta learning rule* [73, 67] where the self-generated patterns, \mathbf{v}_t , $\phi(\mathbf{k}_t)$,
 127 and $\sigma(\beta_t)$, play the role of *target*, *input*, and *learning rate* of the learning rule respectively. The
 128 delta rule is crucial for the performance of LTs [67, 68, 74, 75].

129 The initial weight matrix \mathbf{W}_0 is the only trainable parameters of this layer, that encodes the initial
 130 self-modification algorithm. We use the layer above as a direct replacement to the self-attention layer
 131 in the Transformer architecture [59]; and use the multi-head version of the computation above [19].

132 3 Method

133 **Task Formulation.** We formulate continual learning as a long-span sequence learning task. Let
 134 D, N, K, L denote positive integers. Consider two N -way classification tasks **A** and **B** to be
 135 learned sequentially (as we’ll see, this can be straightforwardly extended to more tasks). The
 136 formulation here applies to both “meta-training” and “meta-test” phases (see Appendix A.1 for more
 137 on this terminology). We denote the respective training datasets as \mathcal{A} and \mathcal{B} , and test sets as \mathcal{A}'
 138 and \mathcal{B}' . We assume that each datapoint in these datasets consists of one input feature $\mathbf{x} \in \mathbb{R}^D$ of
 139 dimension D (generically denoted as vector \mathbf{x} , but it is an image in all our experiments) and one label
 140 $y \in \{1, \dots, N\}$. We consider two sequences of L training examples $((\mathbf{x}_1^A, y_1^A), \dots, (\mathbf{x}_L^A, y_L^A))$ and
 141 $((\mathbf{x}_1^B, y_1^B), \dots, (\mathbf{x}_L^B, y_L^B))$ sampled from the respective training sets \mathcal{A} and \mathcal{B} . In practice, $L = NK$
 142 where K is the number of training examples for each class. By concatenating these two sequences,
 143 we obtain one long sequence representing CL examples to be presented as an input sequence to
 144 a (left-to-right) auto-regressive model. At the end of the sequence, the model is tasked to make
 145 predictions on test examples sampled from both \mathcal{A}' and \mathcal{B}' ; we assume a single test example for
 146 each task (hence, without index): $(\mathbf{x}^{A'}, y^{A'})$ and $(\mathbf{x}^{B'}, y^{B'})$ respectively; which we simply denote as
 147 $(\mathbf{x}_{\text{test}}^A, y_{\text{test}}^A)$ and $(\mathbf{x}_{\text{test}}^B, y_{\text{test}}^B)$ instead.

148 Our model is a self-referential NN that modifies its own weight matrices as a function of input
 149 observations. To simplify the notation, we denote the *state* of our self-referential NN as a
 150 *single* SRWM \mathbf{W}_* (even though it may have many of them in practice) where we’ll replace $*$
 151 by various symbols representing the context/inputs it has observed. Given a training sequence
 152 $((\mathbf{x}_1^A, y_1^A), \dots, (\mathbf{x}_L^A, y_L^A), (\mathbf{x}_1^B, y_1^B), \dots, (\mathbf{x}_L^B, y_L^B))$, our model auto-regressively consumes one input
 153 at a time, from left to right, in the auto-regressive fashion. Let $\mathbf{W}_{\mathcal{A}}$ denote the state of the SRWM that
 154 has consumed the first part of the sequence, i.e., the examples from Task **A**, $(\mathbf{x}_1^A, y_1^A), \dots, (\mathbf{x}_L^A, y_L^A)$,
 155 and let $\mathbf{W}_{\mathcal{A}, \mathcal{B}}$ denote the state of our SRWM having observed the entire sequence.

156 **ACL Meta-Training Objectives.** The ACL meta-training objective function tasks the model to
 157 correctly predict the test examples of all tasks learned so far at each task boundaries. That is, in the
 158 case of two-task scenario described above (learning Task **A** then Task **B**), we use the weight matrix
 159 $\mathbf{W}_{\mathcal{A}}$ to predict the label y_{test}^A from input $(\mathbf{x}_{\text{test}}^A, \emptyset)$, and we use the weight matrix $\mathbf{W}_{\mathcal{A}, \mathcal{B}}$ to predict the
 160 label y_{test}^B from input $(\mathbf{x}_{\text{test}}^B, \emptyset)$ as well as the label y_{test}^A from input $(\mathbf{x}_{\text{test}}^A, \emptyset)$. By letting $p(y|\mathbf{x}; \mathbf{W}_*)$
 161 denote the model’s output probability for label $y \in \{1, \dots, N\}$ given input \mathbf{x} and model weights/state
 162 \mathbf{W}_* , the ACL objective can be expressed as:

$$\text{minimize}_{\theta} - (\log(p(y_{\text{test}}^A | \mathbf{x}_{\text{test}}^A; \mathbf{W}_{\mathcal{A}})) + \log(p(y_{\text{test}}^B | \mathbf{x}_{\text{test}}^B; \mathbf{W}_{\mathcal{A}, \mathcal{B}})) + \log(p(y_{\text{test}}^A | \mathbf{x}_{\text{test}}^A; \mathbf{W}_{\mathcal{A}, \mathcal{B}}))) \quad (4)$$

163 for an arbitrary input meta-training sequence $((\mathbf{x}_1^A, y_1^A), \dots, (\mathbf{x}_L^A, y_L^A), (\mathbf{x}_1^B, y_1^B), \dots, (\mathbf{x}_L^B, y_L^B))$
 164 (which is extensible to mini-batches with multiple such sequences), where θ denotes the model
 165 parameters (for the SRWM layer, it is the initial weights \mathbf{W}_0). Figure 1 illustrates the overall
 166 meta-training process of ACL.

Table 1: 5-way classification accuracies using 15 (meta-test training) examples for each class in the context. Each row is a single model. **Bold** numbers highlight cases where in-context catastrophic forgetting is avoided through ACL.

Meta-Training Tasks			Meta-Test Tasks: Context/Train (top) & Test (bottom)					
			ACL	A	A \rightarrow B		B	B \rightarrow A
Task A	Task B		A	B	A	B	A	B
Omniglot	Mini-ImageNet	No	97.6 \pm 0.2	52.8 \pm 0.7	22.9 \pm 0.7	52.1 \pm 0.8	97.8 \pm 0.3	20.4 \pm 0.6
		Yes	98.3 \pm 0.2	54.4 \pm 0.8	98.2 \pm 0.2	54.8 \pm 0.9	98.0 \pm 0.3	54.6 \pm 1.0
FC100	Mini-ImageNet	No	49.7 \pm 0.7	55.0 \pm 1.0	21.3 \pm 0.7	55.1 \pm 0.6	49.9 \pm 0.8	21.7 \pm 0.8
		Yes	53.8 \pm 1.7	52.5 \pm 1.2	46.2 \pm 1.3	59.9 \pm 0.7	45.5 \pm 0.9	53.0 \pm 0.6

Table 2: Similar to Table 1 above but using MNIST and CIFAR-10 (unseen domains) for meta-testing.

Meta-Training Tasks			Meta-Test Tasks: Context/Train (top) & Test (bottom)					
			ACL	MNIST	MNIST \rightarrow CIFAR-10		CIFAR-10	CIFAR-10 \rightarrow MNIST
Task A	Task B		MNIST	CIFAR-10	MNIST	CIFAR-10	MNIST	CIFAR-10
Omniglot	Mini-ImageNet	No	71.1 \pm 4.0	49.4 \pm 2.4	43.7 \pm 2.3	51.5 \pm 1.4	68.9 \pm 4.1	24.9 \pm 3.2
		Yes	75.4 \pm 3.0	50.8 \pm 1.3	81.5 \pm 2.7	51.6 \pm 1.3	77.9 \pm 2.3	51.8 \pm 2.0
FC100	Mini-ImageNet	No	60.1 \pm 2.0	56.1 \pm 2.3	17.2 \pm 3.5	54.4 \pm 1.7	58.6 \pm 1.6	21.2 \pm 3.1
		Yes	70.0 \pm 2.4	51.0 \pm 1.0	68.2 \pm 2.7	59.2 \pm 1.7	66.9 \pm 3.4	52.5 \pm 1.3

167 The ACL objective function above (Eq. 4) is simple but encapsulates desiderata for continual learning
 168 (Sec. 2.1). The last term of Eq. 4 with $p(y_{\text{test}}^A | \mathbf{x}_{\text{test}}^A; \mathbf{W}_{\mathcal{A}, \mathcal{B}})$ or schematically $p(\mathcal{A}' | \mathcal{A}, \mathcal{B})$, optimizes
 169 for *backward transfer*: (1) remembering the first task **A** after learning **B** (combatting catastrophic
 170 forgetting), and (2) leveraging learning of **B** to improve performance on the past task **A**. The
 171 second term of Eq. 4, $p(y_{\text{test}}^B | \mathbf{x}_{\text{test}}^B; \mathbf{W}_{\mathcal{A}, \mathcal{B}})$ or schematically $p(\mathcal{B}' | \mathcal{A}, \mathcal{B})$, optimizes *forward transfer*
 172 leveraging the past learning experience of **A** to improve predictions in the second task **B**, in addition
 173 to simply learning to solve Task **B** from the corresponding training examples. To complete, the first
 174 term of Eq. 4 is the single-task meta-learning objective for Task **A**.

175 **Overall Model Architecture.** As we mention in Sec. 2, in our NN architecture, the core sequential
 176 dynamics of CL are learned by the self-referential layers. However, as an image-processing NN, our
 177 model makes use of a vision backend. We use the ‘‘Conv-4’’ architecture [21] (typically used in the
 178 context of few-shot learning) in all our experiments, except in the last one where we use a pre-trained
 179 vision Transformer [76]. Overall, the model takes an image as input, process it through a feedforward
 180 vision NN, whose output is fed to the SRWM-layer block. Note that this is one of the limitations of
 181 this work: more general ACL should also learn to modify the vision components.²

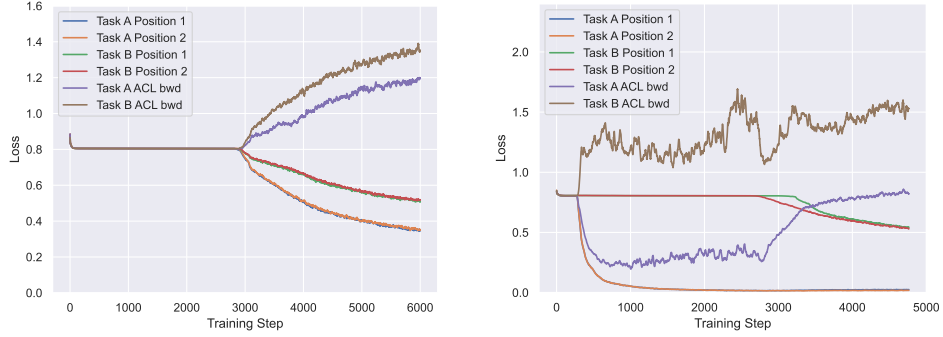
182 Another crucial architectural choice that is specific to continual/multi-task image processing is
 183 normalization layers (see also Bronskill et al. [78]). Typical NNs used in few-shot learning (e.g.,
 184 Vinyals et al. [21]) contain batch normalization (BN; [79]) layers. All our models use instance
 185 normalization (IN; [80]) instead of BN because in our preliminary experiments, we expectably found
 186 IN to generalize much better than BN layers in the CL setting.

187 4 Experiments

188 4.1 Two-Task Setting: Comprehensible Study

189 We first reveal the problem of ‘‘in-context catastrophic forgetting’’ and show how our ACL method
 190 (Sec. 3) can overcome it. As a minimum setting for this, we focus on the two-task ‘‘domain-

²One ‘‘straightforward’’ architecture fitting the bill is an MLP-mixer architecture (Tolstikhin et al. [77]; built of several linear layers), where all linear layers are replaced by the self-referential linear layers of Sec. 2.3. While we implemented such a model, it turned out to be too slow for us to conduct corresponding experiments. Our public code will include a ‘‘self-referential MLP-mixer’’ implementation, but for further experiments, we leave the future work on such an architecture using more efficient CUDA kernels.



(a) Case: Two tasks are learned simultaneously. (b) Case: One task is learned first (here Task A).

Figure 2: **ACL/No**-case meta-training curves displaying 6 individual meta-training loss terms, when the last term of the ACL objective (the backward transfer loss; “*Task A ACL bwd*” and “*Task B ACL bwd*” in the legend) is **not** minimized (**ACL/No** case in Tables 1 and 2). Here Task A is Omniglot and Task B is Mini-ImageNet. We observe that, in both cases, without explicit minimization, backward transfer capability (*purple* and *brown* curves) of the learned learning algorithm gradually degrades as it learns to learn a new task (all other colors), causing in-context catastrophic forgetting. Note that *blue/orange* and *green/red* curve pairs almost overlap; indicating that when a task is learned, the model can learn it whether it is in the first or second segment of the continual learning sequence.

191 incremental” CL setting (see Appendix A.1). We consider two meta-training task combinations:
 192 Omniglot [23] and Mini-ImageNet [21, 22] or FC100 [24] (which is based on CIFAR100 [27]) and
 193 Mini-ImageNet. The order of appearance of two tasks within meta-training sequences is alternated
 194 for every batch. Appendix A.2 provides further details. We compare systems trained with or without
 195 the backward transfer term in the ACL loss (the last term in Eq. 4).

196 Unless otherwise indicated (e.g, later for classic Split-MNIST; Sec. 4.3), all tasks are configured
 197 to be a 5-way classification task. This is one of the classic configurations for few-shot learning tasks,
 198 and also allows us to evaluate the principle of ACL with reasonable computational costs (like any
 199 sequence learning-based meta-learning methods, scaling this to many more classes is challenging; we
 200 also discuss this in Sec. 5). For standard datasets such as MNIST, we split the dataset into sub-datasets
 201 of disjoint classes [81]: for example for MNIST which is originally a 10-way classification task, we
 202 split it into two 5-way tasks, one consisting of images of class ‘0’ to ‘4’ (‘MNIST-04’), and another
 203 one made of class ‘5’ to ‘9’ images (‘MNIST-59’). When we refer to a dataset without specifying
 204 the class range, we refer to the first sub-set. Unless stated otherwise, we concatenate 15 examples
 205 from each class for each task in the context for both meta-training and meta-testing (resulting in
 206 sequences of length 75 for each task). All images are resized to 32×32 -size 3-channel images, and
 207 normalized according to the original dataset statistics. We refer to Appendix A for further details.

208 Table 1 shows the results when the models are meta-tested on the test sets of the corresponding
 209 few-shot learning datasets used for meta-training. We observe that for both pairs of meta-training
 210 tasks, the models without the ACL loss *catastrophically forget* the first task after learning the second
 211 one: the accuracy on the first task is at the chance level of about 20% for 5-way classification after
 212 learning the second task in-context (see rows with “ACL No”). The ACL loss clearly addresses this
 213 problem: the ACL-learned CL algorithms preserve the performance of the first task. This effect is
 214 particularly pronounced in the Omniglot/Mini-ImageNet case (involving two very different domains).

215 Table 2 shows evaluations of the same models but using two standard datasets, 5-way MNIST and
 216 CIFAR-10, for meta-testing. Again, ACL-trained models better preserve the memory of the first
 217 task after learning the second one. In the Omniglot/Mini-ImageNet case, we even observe certain
 218 positive backward transfer effects: in particular, in the “MNIST-then-CIFAR10” continual learning
 219 case, the performance on MNIST noticeably improves after learning CIFAR10 (possibly leveraging
 220 ‘more data’ provided in-context).

221 4.2 Analysis: Emergence of In-Context Catastrophic Forgetting

222 Now we closely look at the emergence of “in-context catastrophic forgetting” during meta-training
 223 for the baseline models trained **without** the backward transfer term (the last/third term in Eq. 4) in

Table 3: Classification accuracies (%) on the **Split-MNIST** domain-incremental (DIL) and class-incremental learning (CIL) settings [6]. Both tasks are 5-task CL problems. For the CIL case, we also report the 2-task case for which we can directly evaluate our out-of-the-box ACL meta-learner of Sec. 4.1 (trained with a 5-way output and the 2-task ACL loss) which, however, is not applicable (N.A.) to the 5-task CIL requiring a 10-way output. Mean/std over 10 training/meta-testing runs. **No method here requires replay memory**. See Appendix A.7 & B for further details and discussions.

Method	Domain Incremental	Class Incremental	
	5-task	2-task	5-task
Plain Stochastic Gradient Descent (SGD)	63.2 ± 0.4	48.8 ± 0.1	19.5 ± 0.1
Adam	55.2 ± 1.4	49.7 ± 0.1	19.7 ± 0.1
Adam + L2	66.0 ± 3.7	51.8 ± 1.9	22.5 ± 1.1
Elastic Weight Consolidation (EWC)	58.9 ± 2.6	49.7 ± 0.1	19.8 ± 0.1
Online EWC	57.3 ± 1.4	49.7 ± 0.1	19.8 ± 0.1
Synaptic Intelligence (SI)	64.8 ± 3.1	49.4 ± 0.2	19.7 ± 0.1
Memory Aware Synapses (MAS)	68.6 ± 6.9	49.6 ± 0.1	19.5 ± 0.3
Learning w/o Forgetting (LwF)	71.0 ± 1.3	-	24.2 ± 0.3
Online-aware Meta Learning (OML)	69.9 ± 2.8	46.6 ± 7.2	24.9 ± 4.1
+ optimized # meta-testing iterations	73.6 ± 5.3	62.1 ± 7.9	34.2 ± 4.6
Generative Meta-Continual Learning (GeMCL)	63.8 ± 3.8	91.2 ± 2.8	79.0 ± 2.1
ACL (Out-of-the-box, DIL, 2-task ACL model; Sec. 4.1)	72.2 ± 0.9	71.5 ± 5.9	N.A.
+ meta-finetuned with 5-task ACL loss, Omniglot	84.5 ± 1.6	96.0 ± 1.0	84.3 ± 1.2

224 the ACL objective loss (corresponding to the **ACL/No** cases in Tables 1 and 2). We focus on the
 225 Omniglot/Mini-ImageNet case, but similar trends can also be observed in the FC100/Mini-ImageNet
 226 case. Figures 2a and 2b show two representative cases we typically observe. These figures show an
 227 evolution of six individual meta-training loss terms (the lower the better), reported separately for
 228 the cases where Task A (here Omniglot) or Task B (here Mini-ImageNet) appears at the first (1) or
 229 second (2) position in the 2-task CL meta-training training sequences. 4 out of 6 curves correspond to
 230 the learning progress, showing whether the model becomes capable of in-context learning the given
 231 task (A or B) at the given position (1 or 2). The 2 remaining curves are the ACL backward transfer
 232 losses, also measured for Task A and B separately here.

233 Figure 2a shows the case where two tasks are learned about at the same time. We observe that when
 234 the learning curves go down, the ACL losses go up, indicating that more the model learns, more it
 235 tends to forget the task in-context learned previously. We also find this same trend when one task
 236 is learned before the other one as is the case in Figure 2b. Here Task A alone is learned first; while
 237 Task B is not learned, both learning and ACL curves go down for Task A (essentially, as the model
 238 does not learn the second task, there is no force that encourages forgetting). After around 3000 steps,
 239 the model also starts learning Task B. From this point, the ACL loss for Task A also starts to go
 240 up, indicating again an *opposing force effect* between learning a new task and remembering a past
 241 task. These observations clearly indicate that, without explicitly taking into account the backward
 242 transfer loss as part of learning objectives, our gradient descent search tends to find solutions/CL
 243 algorithms that prefer to erase previously learned knowledge (this is rather intuitive; it seems easier to
 244 find such algorithms that ignore any influence of the current learning to past learning than those that
 245 also preserve prior knowledge). In all cases, we find our ACL objective to be crucial for the learned
 246 CL algorithms to be capable of remembering the old task while also learning the new one.

247 4.3 General Evaluation

248 **Evaluation on Standard Split-MNIST.** Here we evaluate ACL on the standard Split-MNIST task in
 249 domain-incremental and class-incremental settings [6, 28], and compare its performance to existing
 250 CL and meta-CL algorithms (see Appendix A.7 for full references of these methods). Our comparison
 251 focuses on methods that do not require replay memory. Table 3 shows the results. Since our
 252 ACL-trained models are general-purpose learners, they can be directly evaluated (meta-tested) on
 253 a new task, here Split-MNIST. The second-to-last row of Table 3, “ACL (Out-of-the-box model)”,
 254 corresponds to our model from Sec. 4.1 meta-trained on Omniglot and Mini-ImageNet using the
 255 2-task ACL objective. It performs very competitively with the best existing methods in the domain-

256 incremental setting, while it largely outperforms them (all but another meta-CL method, GeMCL) in
 257 the 2-task class-incremental setting. The same model can be further meta-finetuned using the 5-task
 258 version of the ACL loss (here we only used Omniglot as the meta-training data). The resulting model
 259 (the last row of Table 3) outperforms all other methods in all settings studied here. Note that on
 260 the ‘in-domain’ Omniglot test set, ACL and GeMCL perform similarly (see Appendix B.2/Table 9).
 261 We are not aware of any existing hand-crafted CL algorithms that can achieve ACL’s performance
 262 without any replay memory. We refer to Appendix A.7/B for further discussions and ablation studies.

263 **Evaluation on diverse task domains.** Using the setting of Sec. 4.1, we also evaluate our ACL-trained
 264 models for CL involving more tasks/domains; using meta-test sequences made of MNIST, CIFAR-10,
 265 and Fashion MNIST. We also evaluate the impact of the number of tasks in the ACL objective: in
 266 addition to the model meta-trained on Omniglot/Mini-ImageNet (Sec. 4.1), we also meta-train a model
 267 (with the same architecture and hyper-parameters) using 3 tasks, Omniglot, Mini-ImageNet, and
 268 FC100, using the 3-task ACL objective (see Appendix A.5); which is meta-trained not only on longer
 269 CL sequences but also on more data. The full results of this experiment can be found in Appendix
 270 B.4. We find that the two ACL-trained models are indeed capable of retaining the knowledge without
 271 catastrophic forgetting for multiple tasks during meta-testing, while the performance on prior tasks
 272 gradually degrades as the model learns new tasks, and performance on new tasks becomes moderate
 273 (see also Sec. 5 on limitations). The 3-task version outperforms the 2-task one overall, encouragingly
 274 indicating a potential for further improvements even with a fixed parameter count.

275 **Going beyond: limitations and outlook.** The experiments presented above effectively demonstrate
 276 the possibility to encode a continual learning algorithm into self-referential weight matrices, that
 277 outperforms handcrafted learning algorithms and existing metalearning approaches for CL. While
 278 we consider this as an important result for metalearning and in-context learning in general, we note
 279 that current state-of-the-art CL methods use neither regularization-based CL algorithms nor meta-
 280 continual learning methods we mention above, but the so-called *learning to prompt* (L2P)-family
 281 of methods [33, 34] that leverage pre-trained models, namely a vision Transformer (ViT) pre-trained
 282 on ImageNet [76]. A natural question we should ask is whether we could foresee ACL beyond the
 283 scope considered so far, and evaluate it in such a setting. To study this, we take a pre-trained (frozen)
 284 vision model, and add self-referential layers (to be meta-trained from scratch) on top of it to build a
 285 continual learner. This allows us to highlight an important challenge of in-context CL in what follows.

286 We use two tasks from the L2P works above [33, 34]: 5-datasets [32] and Split-CIFAR-100, in the
 287 class-incremental setting, but we focus on a “mini” versions thereof: we only use the two first classes
 288 within each task (i.e., 2-way version) and for Split-CIFAR100, we only use the 5 first tasks; as we’ll
 289 see, this setting is enough to illustrate an important limitation of in-context CL. Again following
 290 L2P [33, 34], we use ViT-B/16 [76] (available via PyTorch) as the pre-trained vision model, which
 291 we keep frozen. We use the same configuration for the self-referential component from the Split-
 292 MNIST experiment. We meta-train the resulting model using Mini-ImageNet and Omniglot with the
 293 5-task ACL loss. Table 4 shows the results. Even in this simple “mini” version of the tasks, ACL’s
 294 performance is far behind that of L2P methods. Notably, the frozen ImageNet-pre-trained features
 295 with the meta-learner trained on Mini-ImageNet and Omniglot are not enough to perform well on the
 296 5-th task of Split-CIFAR100, and SVHN and notMNIST of 5-datasets. This shows the necessity to
 297 meta-train on more diverse tasks for in-context CL to be possibly successful in more general settings.

Table 4: Experiments with “mini” Split-CIFAR100 and 5-datasets tasks. Meta-training is done using **Mini-ImageNet** and **Omniglot**. All meta-evaluation images are therefore from unseen domains. Numbers marked with * are *reference* numbers (evaluated in the more challenging, original version of these tasks) which can not be directly compared to ours.

	Split-CIFAR100	5-datasets
L2P [34]	83.9* \pm 0.3	81.1* \pm 0.9
DualPrompt [34]	86.5* \pm 0.3	88.1* \pm 0.4
ACL (Individual Task)	Task 1 95.9 \pm 0.9 Task 2 85.6 \pm 3.6 Task 3 93.4 \pm 1.4 Task 4 97.0 \pm 0.7 Task 5 67.6 \pm 7.0	CIFAR10 91.3 \pm 1.2 MNIST 98.9 \pm 0.3 Fashion 93.5 \pm 2.0 SVHN 66.1 \pm 9.4 notMNIST 76.3 \pm 6.7
ACL	68.3 \pm 2.0	61.5 \pm 2.1

298 5 Discussion

299 **Other Limitations.** In addition to the limitations already mentioned above, here we discuss others.
300 First of all, as an in-context/learned learning algorithm, there are challenges in terms of both domain
301 and length generalization (we qualitatively observe these to some extent in Sec. 4; further discussion
302 and experimental results are presented in Appendix B.3 & B.5). Regarding the length generalization,
303 we note that unlike the standard “quadratic” Transformers, linear Transformers/FWPs-like SRWMs
304 can be trained by *carrying over states* across two consecutive batches for arbitrarily long sequences.
305 Such an approach has been successfully applied to language modeling with FWPs [67]. This
306 possibility, however, has not been investigated here, and is left for future work. Also, directly scaling
307 ACL for real-world tasks requiring many more classes does not seem straightforward: it would
308 require very long training sequences. That said, it may be possible that ACL could be achieved
309 without exactly following the process we propose; as we discuss below for the case of LLMs, certain
310 real-world data may naturally give rise to an ACL-like objective. This work is also limited to the
311 task of image classification, which can be solved by feedforward NNs. Future work may investigate
312 the possibility to extend ACL to continual learning of sequence learning tasks, such as continually
313 learning new languages. Finally, ACL learns CL algorithms that are specific to the pre-specified
314 model architecture; more general meta-learning algorithms may aim at achieving learning algorithms
315 that are applicable to any model, as is the case for many classic learning algorithms.

316 **Related work.** There are several recent works that are categorized as ‘meta-continual learning’ or
317 ‘continual meta-learning’ (see, e.g., [29, 30, 82–84, 51]). For example, Javed and White [29], Beaulieu
318 et al. [30] use “model-agnostic meta-learning” (MAML; [85, 86]) to meta-learn *representations* for
319 CL while still making use of classic learning algorithms for CL; this requires tuning of the learning
320 rate and number of iterations for optimal performance during CL at meta-test time (see, e.g., Appendix
321 A.7). In contrast, our approach learn *learning algorithms* in the spirit of Hochreiter et al. [14], Younger
322 et al. [15]; this may be categorized as ‘in-context continual learning.’ Several recent works (see, e.g.,
323 [87, 88]) mention the possibility of such in-context CL but existing works [19, 89, 90] that learn mul-
324 tiple tasks sequentially in-context do not focus on catastrophic forgetting which is one of the central
325 challenges of CL. Here we show that in-context learning also suffers from catastrophic forgetting in
326 general (Sec. 4.1-4.2) and propose ACL to address this problem. We also note that the use of SRWM is
327 relevant to ‘continual meta-learning’ since with a regular sequence processor with slow weights, there
328 remains the question of how to continually learn the slow weights (meta-parameters). In principle, re-
329 cursive self-modification as in SRWM is an answer to this question as it collapses such meta-levels into
330 single self-reference [12]. We also refer to [91–93] for other prior work on meta-continual learning.

331 **Artificial v. Natural ACL in Large Language Models?** Recently, “on-the-fly” few-shot/meta
332 learning capability of sequence processing NNs has attracted broader interests in the context of large
333 language models (LLMs; [94]). In fact, the task of language modeling itself has a form of *sequence*
334 *processing with error feedback* (essential for meta-learning [95]): the correct label to be predicted is
335 fed to the model with a delay of one time step in an auto-regressive manner. Trained on a large amount
336 of text covering a wide variety of credit assignment paths, LLMs exhibit certain sequential few-shot
337 learning capabilities in practice [96]. This was rebranded as *in-context learning*, and has been the
338 subject of numerous recent studies (e.g., [97–103, 71, 72]). Here we explicitly/artificially construct
339 ACL meta-training sequences and objectives, but in modern LLMs trained on a large amount of data
340 mixing a large diversity of dependencies using a large backpropagation span, it is conceivable that
341 some ACL-like objectives may naturally appear in the data.

342 6 Conclusion

343 Our Automated Continual Learning (ACL) trains sequence-processing self-referential neural networks
344 (SRNNs) to learn their own in-context continual (meta-)learning algorithms. ACL encodes classic
345 desiderata for continual learning (e.g., forward and backward transfer) into the objective function of
346 the meta-learner. ACL uses gradient descent to deal with classic challenges of CL, to automatically
347 discover CL algorithms with good behavior. Once trained, our SRNNs autonomously run their
348 own CL algorithms without requiring any human intervention. Our experiments reveal the original
349 problem of in-context catastrophic forgetting, and demonstrate the effectiveness of the proposed
350 approach to combat it. We demonstrate very promising results on the classic Split-MNIST benchmark
351 where existing hand-crafted algorithms fail, while also discussing its limitations in more general
352 scenarios. We believe this comprehensive study to be an important step for in-context CL research.

References

- 353 [1] David Eagleman. *Livewired: The inside story of the ever-changing brain*. 2020.
- 354 [2] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks:
355 The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages
356 109–165. 1989.
- 357 [3] Roger Ratcliff. Connectionist models of recognition memory: constraints imposed by learning
358 and forgetting functions. *Psychological review*, 97(2):285, 1990.
- 359 [4] Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive*
360 *sciences*, 3(4):128–135, 1999.
- 361 [5] James L McClelland, Bruce L McNaughton, and Randall C O’Reilly. Why there are comple-
362 mentary learning systems in the hippocampus and neocortex: insights from the successes and
363 failures of connectionist models of learning and memory. *Psychological review*, 102(3):419,
364 1995.
- 365 [6] Yen-Chang Hsu, Yen-Cheng Liu, Anita Ramasamy, and Zsolt Kira. Re-evaluating continual
366 learning scenarios: A categorization and case for strong baselines. In *NeurIPS Workshop on*
367 *Continual Learning*, Montréal, Canada, December 2018.
- 368 [7] Chris A Kortge. Episodic memory in connectionist networks. In *12th Annual Conference. CSS*
369 *Pod*, pages 764–771, 1990.
- 370 [8] Robert M French. Using semi-distributed representations to overcome catastrophic forgetting
371 in connectionist networks. In *Proc. Cognitive science society conference*, volume 1, pages
372 173–178, 1991.
- 373 [9] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins,
374 Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska,
375 et al. Overcoming catastrophic forgetting in neural networks. *Proc. National academy of*
376 *sciences*, 114(13):3521–3526, 2017.
- 377 [10] Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska,
378 Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable frame-
379 work for continual learning. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 4535–
380 4544, Stockholm, Sweden, July 2018.
- 381 [11] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic
382 intelligence. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 3987–3995, Sydney,
383 Australia, August 2017.
- 384 [12] Jürgen Schmidhuber. Steps towards “self-referential” learning. Technical Report CU-CS-627-
385 92, Dept. of Comp. Sci., University of Colorado at Boulder, November 1992.
- 386 [13] Jürgen Schmidhuber. *Evolutionary principles in self-referential learning, or on learning how*
387 *to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987.
- 388 [14] Sepp Hochreiter, A. Steven Younger, and Peter R. Conwell. Learning to learn using gradient
389 descent. In *Proc. Int. Conf. on Artificial Neural Networks (ICANN)*, volume 2130, pages
390 87–94, Vienna, Austria, August 2001.
- 391 [15] A Steven Younger, Peter R Conwell, and Neil E Cotter. Fixed-weight on-line learning. *IEEE*
392 *Transactions on Neural Networks*, 10(2):272–283, 1999.
- 393 [16] Neil E Cotter and Peter R Conwell. Learning algorithms and fixed dynamics. In *Proc. Int.*
394 *Joint Conf. on Neural Networks (IJCNN)*, pages 799–801, Seattle, WA, USA, July 1991.
- 395 [17] Neil E Cotter and Peter R Conwell. Fixed-weight networks can learn. In *Proc. Int. Joint Conf.*
396 *on Neural Networks (IJCNN)*, pages 553–559, San Diego, CA, USA, June 1990.
- 397 [18] Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive
398 meta-learner. In *Int. Conf. on Learning Representations (ICLR)*, Vancouver, Canada, 2018.
- 399

- 400 [19] Kazuki Irie, Imanol Schlag, Róbert Csordás, and Jürgen Schmidhuber. A modern self-
401 referential weight matrix that learns to modify itself. In *Proc. Int. Conf. on Machine Learning*
402 (*ICML*), pages 9660–9677, Baltimore, MA, USA, July 2022.
- 403 [20] Stephen T Grossberg. *Studies of mind and brain: Neural principles of learning, perception,*
404 *development, cognition, and motor control*. Springer, 1982.
- 405 [21] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Koray Kavukcuoglu, and Daan Wierstra.
406 Matching networks for one shot learning. In *Proc. Advances in Neural Information Processing*
407 *Systems (NIPS)*, pages 3630–3638, Barcelona, Spain, December 2016.
- 408 [22] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *Int. Conf.*
409 *on Learning Representations (ICLR)*, Toulon, France, April 2017.
- 410 [23] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept
411 learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- 412 [24] Boris N. Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. TADAM: task dependent
413 adaptive metric for improved few-shot learning. In *Proc. Advances in Neural Information*
414 *Processing Systems (NeurIPS)*, pages 719–729, Montréal, Canada, December 2018.
- 415 [25] Yann LeCun, Corinna Cortes, and Christopher JC Burges. The MNIST database of handwritten
416 digits. URL <http://yann.lecun.com/exdb/mnist>, 1998.
- 417 [26] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for
418 benchmarking machine learning algorithms. *Preprint arXiv:1708.07747*, 2017.
- 419 [27] Alex Krizhevsky. Learning multiple layers of features from tiny images. Master’s thesis,
420 Computer Science Department, University of Toronto, 2009.
- 421 [28] Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. In *NeurIPS*
422 *Workshop on Continual Learning*, Montréal, Canada, December 2018.
- 423 [29] Khurram Javed and Martha White. Meta-learning representations for continual learning.
424 In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, pages 1818–1828,
425 Vancouver, BC, Canada, December 2019.
- 426 [30] Shawn Beaulieu, Lapo Frati, Thomas Miconi, Joel Lehman, Kenneth O. Stanley, Jeff Clune,
427 and Nick Cheney. Learning to continually learn. In *Proc. European Conf. on Artificial*
428 *Intelligence (ECAI)*, pages 992–1001, August 2020.
- 429 [31] Mohammadamin Banayeeanzade, Rasoul Mirzaiezhadeh, Hosein Hasani, and Mahdiah So-
430 leymani. Generative vs. discriminative: Rethinking the meta-continual learning. In *Proc.*
431 *Advances in Neural Information Processing Systems (NeurIPS)*, pages 21592–21604, Virtual
432 only, December 2021.
- 433 [32] Sayna Ebrahimi, Franziska Meier, Roberto Calandra, Trevor Darrell, and Marcus Rohrbach.
434 Adversarial continual learning. In *Proc. European Conf. on Computer Vision (ECCV)*, pages
435 386–402, Glasgow, UK, August 2020.
- 436 [33] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su,
437 Vincent Perot, Jennifer G. Dy, and Tomas Pfister. Learning to prompt for continual learning.
438 In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 139–149,
439 New Orleans, LA, USA, June 2022.
- 440 [34] Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi
441 Ren, Guolong Su, Vincent Perot, Jennifer G. Dy, and Tomas Pfister. Dualprompt: Complemen-
442 tary prompting for rehearsal-free continual learning. In *Proc. European Conf. on Computer*
443 *Vision (ECCV)*, pages 631–648, Tel Aviv, Israel, October 2022.
- 444 [35] Sebastian Thrun. Lifelong learning algorithms. In *Learning to learn*, pages 181–209. 1998.
- 445 [36] Rich Caruana. Multitask learning. *Machine learning*, 28:41–75, 1997.

- 446 [37] Mark B. Ring. *Continual Learning in Reinforcement Environments*. PhD thesis, University of
447 Texas at Austin, Austin, TX, USA, 1994.
- 448 [38] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick,
449 Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *Preprint*
450 *arXiv:1606.04671*, 2016.
- 451 [39] Anthony Robins. Catastrophic forgetting, rehearsal and pseudorehearsal. *Connection Science*,
452 7(2):123–146, 1995.
- 453 [40] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep
454 generative replay. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pages
455 2990–2999, Long Beach, CA, USA, December 2017.
- 456 [41] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P. Lillicrap, and Gregory Wayne.
457 Experience replay for continual learning. In *Proc. Advances in Neural Information Processing*
458 *Systems (NeurIPS)*, pages 348–358, Vancouver, Canada, December 2019.
- 459 [42] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and
460 Gerald Tesaro. Learning to learn without forgetting by maximizing transfer and minimizing
461 interference. In *Int. Conf. on Learning Representations (ICLR)*, New Orleans, LA, USA, May
462 2019.
- 463 [43] Yaqian Zhang, Bernhard Pfahringer, Eibe Frank, Albert Bifet, Nick Jin Sean Lim, and Yunzhe
464 Jia. A simple but strong baseline for online continual learning: Repeated augmented rehearsal.
465 In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, New Orleans, LA,
466 USA, December 2022.
- 467 [44] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning.
468 In *Proc. Advances in Neural Information Processing Systems (NIPS)*, pages 6467–6476, Long
469 Beach, CA, USA, December 2017.
- 470 [45] Tom Veniat, Ludovic Denoyer, and Marc’Aurelio Ranzato. Efficient continual learning with
471 modular networks and task-driven priors. In *Int. Conf. on Learning Representations (ICLR)*,
472 Virtual only, May 2021.
- 473 [46] Devang K Naik and Richard J Mammone. Meta-neural networks that learn by learning. In
474 *Proc. International Joint Conference on Neural Networks (IJCNN)*, volume 1, pages 437–442,
475 Baltimore, MD, USA, June 1992.
- 476 [47] Tom Bosc. Learning to learn neural networks. In *NIPS Workshop on Reasoning, Attention,*
477 *Memory*, Montreal, Canada, December 2015.
- 478 [48] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy P. Lillicrap.
479 Meta-learning with memory-augmented neural networks. In *Proc. Int. Conf. on Machine*
480 *Learning (ICML)*, pages 1842–1850, New York City, NY, USA, June 2016.
- 481 [49] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RL^2 :
482 Fast reinforcement learning via slow reinforcement learning. *Preprint arXiv:1611.02779*,
483 2016.
- 484 [50] Jane Wang, Zeb Kurth-Nelson, Hubert Soyer, Joel Z. Leibo, Dhruva Tirumala, Rémi Munos,
485 Charles Blundell, Dharshan Kumaran, and Matt M. Botvinick. Learning to reinforcement
486 learn. In *Proc. Annual Meeting of the Cognitive Science Society (CogSci)*, London, UK, July
487 2017.
- 488 [51] Tsendsuren Munkhdalai and Hong Yu. Meta networks. In *Proc. Int. Conf. on Machine*
489 *Learning (ICML)*, pages 2554–2563, Sydney, Australia, August 2017.
- 490 [52] Tsendsuren Munkhdalai and Adam Trischler. Metalearning with Hebbian fast weights. *Preprint*
491 *arXiv:1807.05076*, 2018.
- 492 [53] Thomas Miconi, Kenneth Stanley, and Jeff Clune. Differentiable plasticity: training plastic
493 neural networks with backpropagation. In *Proc. Int. Conf. on Machine Learning (ICML)*,
494 pages 3559–3568, Stockholm, Sweden, July 2018.

- 495 [54] Thomas Miconi, Aditya Rawal, Jeff Clune, and Kenneth O. Stanley. Backpropamine: training
496 self-modifying neural networks with differentiable neuromodulated plasticity. In *Int. Conf. on*
497 *Learning Representations (ICLR)*, New Orleans, LA, USA, May 2019.
- 498 [55] Tsendsuren Munkhdalai, Alessandro Sordani, Tong Wang, and Adam Trischler. Metalearned
499 neural memory. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*,
500 pages 13310–13321, Vancouver, Canada, December 2019.
- 501 [56] Louis Kirsch and Jürgen Schmidhuber. Meta learning backpropagation and improving it. In
502 *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, pages 14122–14134,
503 Virtual only, December 2021.
- 504 [57] Mark Sandler, Max Vladymyrov, Andrey Zhmoginov, Nolan Miller, Tom Madams, Andrew
505 Jackson, and Blaise Agüera y Arcas. Meta-learning bidirectional update rules. In *Proc. Int.*
506 *Conf. on Machine Learning (ICML)*, pages 9288–9300, Virtual only, July 2021.
- 507 [58] Mike Huisman, Thomas M Moerland, Aske Plaat, and Jan N van Rijn. Are LSTMs good
508 few-shot learners? *Machine Learning*, pages 1–28, 2023.
- 509 [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
510 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proc. Advances in Neural*
511 *Information Processing Systems (NIPS)*, pages 5998–6008, Long Beach, CA, USA, December
512 2017.
- 513 [60] Kazuki Irie, Róbert Csordás, and Jürgen Schmidhuber. Practical computational power of linear
514 transformers and their recurrent and self-referential extensions. In *Proc. Conf. on Empirical*
515 *Methods in Natural Language Processing (EMNLP)*, Sentosa, Singapore, 2023.
- 516 [61] Jürgen Schmidhuber. A self-referential weight matrix. In *Proc. Int. Conf. on Artificial Neural*
517 *Networks (ICANN)*, pages 446–451, Amsterdam, Netherlands, September 1993.
- 518 [62] Jürgen Schmidhuber. Learning to control fast-weight memories: An alternative to recurrent
519 nets. Technical Report FKI-147-91, Institut für Informatik, Technische Universität München,
520 March 1991.
- 521 [63] Jürgen Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic
522 recurrent networks. *Neural Computation*, 4(1):131–139, 1992.
- 523 [64] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers
524 are RNNs: Fast autoregressive transformers with linear attention. In *Proc. Int. Conf. on*
525 *Machine Learning (ICML)*, Virtual only, July 2020.
- 526 [65] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane,
527 Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking
528 attention with performers. In *Int. Conf. on Learning Representations (ICLR)*, Virtual only,
529 2021.
- 530 [66] Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah A Smith, and Lingpeng
531 Kong. Random feature attention. In *Int. Conf. on Learning Representations (ICLR)*, Virtual
532 only, 2021.
- 533 [67] Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear Transformers are secretly fast
534 weight programmers. In *Proc. Int. Conf. on Machine Learning (ICML)*, Virtual only, July
535 2021.
- 536 [68] Kazuki Irie, Imanol Schlag, Róbert Csordás, and Jürgen Schmidhuber. Going beyond linear
537 transformers with recurrent fast weight programmers. In *Proc. Advances in Neural Information*
538 *Processing Systems (NeurIPS)*, Virtual only, December 2021.
- 539 [69] Kazuki Irie, Róbert Csordás, and Jürgen Schmidhuber. The dual form of neural networks
540 revisited: Connecting test time predictions to training patterns via spotlights of attention. In
541 *Proc. Int. Conf. on Machine Learning (ICML)*, Baltimore, MD, USA, July 2022.

- 542 [70] Mark A. Aizerman, Emmanuil M. Braverman, and Lev I. Rozonoer. Theoretical foundations
543 of potential function method in pattern recognition. *Automation and Remote Control*, 25(6):
544 917–936, 1964.
- 545 [71] Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander
546 Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by
547 gradient descent. In *Proc. Int. Conf. on Machine Learning (ICML)*, Honolulu, HI, USA, July
548 2023.
- 549 [72] Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can
550 GPT learn in-context? language models secretly perform gradient descent as meta-optimizers.
551 In *Proc. Findings Association for Computational Linguistics (ACL)*, pages 4005–4019, Toronto,
552 Canada, July 2023.
- 553 [73] Bernard Widrow and Marcian E Hoff. Adaptive switching circuits. In *Proc. IRE WESCON*
554 *Convention Record*, pages 96–104, Los Angeles, CA, USA, August 1960.
- 555 [74] Kazuki Irie, Francesco Faccio, and Jürgen Schmidhuber. Neural differential equations for
556 learning to program neural nets through continuous learning rules. In *Proc. Advances in Neural*
557 *Information Processing Systems (NeurIPS)*, New Orleans, LA, USA, December 2022.
- 558 [75] Kazuki Irie and Jürgen Schmidhuber. Images as weight matrices: Sequential image generation
559 through synaptic learning rules. In *Int. Conf. on Learning Representations (ICLR)*, Kigali,
560 Rwanda, May 2023.
- 561 [76] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
562 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
563 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
564 recognition at scale. In *Int. Conf. on Learning Representations (ICLR)*, Virtual only, May
565 2021.
- 566 [77] Ilya O. Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas
567 Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic,
568 and Alexey Dosovitskiy. MLP-Mixer: An all-MLP architecture for vision. In *Proc. Advances*
569 *in Neural Information Processing Systems (NeurIPS)*, pages 24261–24272, Virtual only,
570 December 2021.
- 571 [78] John Bronskill, Jonathan Gordon, James Requeima, Sebastian Nowozin, and Richard E. Turner.
572 TaskNorm: Rethinking batch normalization for meta-learning. In *Proc. Int. Conf. on Machine*
573 *Learning (ICML)*, pages 1153–1164, Virtual only, 2020.
- 574 [79] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training
575 by reducing internal covariate shift. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages
576 448–456, Lille, France, July 2015.
- 577 [80] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing
578 ingredient for fast stylization. *Preprint arXiv:1607.08022*, 2016.
- 579 [81] Rupesh Kumar Srivastava, Jonathan Masci, Sohrob Kazerounian, Faustino J. Gomez, and Jür-
580 gen Schmidhuber. Compete to compute. In *Proc. Advances in Neural Information Processing*
581 *Systems (NIPS)*, pages 2310–2318, Lake Tahoe, NV, USA, December 2013.
- 582 [82] Massimo Caccia, Pau Rodríguez, Oleksiy Ostapenko, Fabrice Normandin, Min Lin, Lucas
583 Page-Caccia, Issam Hadj Laradji, Irina Rish, Alexandre Lacoste, David Vázquez, and Laurent
584 Charlin. Online fast adaptation and knowledge accumulation (OSAKA): a new approach to
585 continual learning. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*,
586 Virtual only, December 2020.
- 587 [83] Xu He, Jakub Sygnowski, Alexandre Galashov, Andrei A Rusu, Yee Whye Teh, and Razvan
588 Pascanu. Task agnostic continual learning via meta learning. *Preprint arXiv:1906.05201*,
589 2019.

- 590 [84] Pau Ching Yap, Hippolyt Ritter, and David Barber. Addressing catastrophic forgetting in
591 few-shot problems. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 11909–11919,
592 Virtual only, July 2021.
- 593 [85] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast
594 adaptation of deep networks. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 1126–
595 1135, Sydney, Australia, August 2017.
- 596 [86] Chelsea Finn and Sergey Levine. Meta-learning and universality: Deep representations
597 and gradient descent can approximate any learning algorithm. In *Int. Conf. on Learning
598 Representations (ICLR)*, Vancouver, Canada, April 2018.
- 599 [87] Kazuki Irie and Jürgen Schmidhuber. Accelerating neural self-improvement via bootstrapping.
600 In *ICLR Workshop on Mathematical and Empirical Understanding of Foundation Models*,
601 Kigali, Rwanda, May 2023.
- 602 [88] Johannes von Oswald, Eyvind Niklasson, Maximilian Schlegel, Seijin Kobayashi, Nicolas
603 Zucchet, Nino Scherrer, Nolan Miller, Mark Sandler, Max Vladymyrov, Razvan Pascanu, et al.
604 Uncovering mesa-optimization algorithms in Transformers. *Preprint arXiv:2309.05858*, 2023.
- 605 [89] Julian Coda-Forno, Marcel Binz, Zeynep Akata, Matthew Botvinick, Jane X Wang, and Eric
606 Schulz. Meta-in-context learning in large language models. *Preprint arXiv:2305.12907*, 2023.
- 607 [90] Soochan Lee, Jaehyeon Son, and Gunhee Kim. Recasting continual learning as sequence
608 modeling. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, New
609 Orleans, LA, USA, December 2023.
- 610 [91] Jürgen Schmidhuber. On learning how to learn learning strategies. Technical Report FKI-198-
611 94, Institut für Informatik, Technische Universität München, November 1994.
- 612 [92] Jürgen Schmidhuber. Beyond “genetic programming”: Incremental self-improvement. In *Proc.
613 Workshop on Genetic Programming at ML95*, pages 42–49, 1995.
- 614 [93] Jürgen Schmidhuber, Jieyu Zhao, and Marco Wiering. Shifting inductive bias with success-
615 story algorithm, adaptive Levin search, and incremental self-improvement. *Machine Learning*,
616 28(1):105–130, 1997.
- 617 [94] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Lan-
618 guage models are unsupervised multitask learners. [Online]. : [https://blog.openai.com/better-
619 language-models/](https://blog.openai.com/better-language-models/), 2019.
- 620 [95] Jürgen Schmidhuber. Making the world differentiable: On using fully recurrent self-supervised
621 neural networks for dynamic reinforcement learning and planning in non-stationary environ-
622 ments. *Institut für Informatik, Technische Universität München. Technical Report FKI-126*,
623 90, 1990.
- 624 [96] Tom B Brown et al. Language models are few-shot learners. In *Proc. Advances in Neural
625 Information Processing Systems (NeurIPS)*, Virtual only, December 2020.
- 626 [97] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of
627 in-context learning as implicit bayesian inference. In *Int. Conf. on Learning Representations
628 (ICLR)*, Virtual only, April 2022.
- 629 [98] Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and
630 Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning
631 work? In *Proc. Conf. on Empirical Methods in Natural Language Processing (EMNLP)*, pages
632 11048–11064, Abu Dhabi, UAE, December 2022.
- 633 [99] Kang Min Yoo, Junyeob Kim, Hyuhng Joon Kim, Hyunsoo Cho, Hwiyeol Jo, Sang-Woo
634 Lee, Sang-goo Lee, and Taeuk Kim. Ground-truth labels matter: A deeper look into input-
635 label demonstrations. In *Proc. Conf. on Empirical Methods in Natural Language Processing
636 (EMNLP)*, pages 2422–2437, Abu Dhabi, UAE, December 2022.

- 637 [100] Stephanie CY Chan, Adam Santoro, Andrew Kyle Lampinen, Jane X Wang, Aaditya K Singh,
638 Pierre Harvey Richemond, James McClelland, and Felix Hill. Data distributional properties
639 drive emergent in-context learning in transformers. In *Proc. Advances in Neural Information
640 Processing Systems (NeurIPS)*, New Orleans, LA, USA, November 2022.
- 641 [101] Stephanie CY Chan, Ishita Dasgupta, Junkyung Kim, Dharshan Kumaran, Andrew K
642 Lampinen, and Felix Hill. Transformers generalize differently from information stored in
643 context vs in weights. In *NeurIPS Workshop on Memory in Artificial and Real Intelligence
644 (MemARI)*, New Orleans, LA, USA, November 2022.
- 645 [102] Louis Kirsch, James Harrison, Jascha Sohl-Dickstein, and Luke Metz. General-purpose in-
646 context learning by meta-learning transformers. In *NeurIPS Workshop on Memory in Artificial
647 and Real Intelligence (MemARI)*, New Orleans, LA, USA, November 2022.
- 648 [103] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning
649 algorithm is in-context learning? investigations with linear models. In *Int. Conf. on Learning
650 Representations (ICLR)*, Kigali, Rwanda, May 2023.
- 651 [104] Gido M Van de Ven and Andreas S Tolias. Generative replay with feedback connections as a
652 general strategy for continual learning. *Preprint arXiv:1809.10635*, 2018.
- 653 [105] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al.
654 Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on
655 deep learning and unsupervised feature learning*, Granada, Spain, December 2011.
- 656 [106] Yaroslav Bulatov. Notmnist dataset. *Google (Books/OCR), Tech. Rep.[Online]. Available:
657 http://yaroslavvb.blogspot.it/2011/09/notmnist-dataset.html*, 2011.
- 658 [107] Tristan Deleu, Tobias Würfl, Mandana Samiei, Joseph Paul Cohen, and Yoshua Bengio.
659 Torchmeta: A meta-learning library for PyTorch. *Preprint arXiv:1909.06576*, 2019.
- 660 [108] Jerry A Fodor and Zenon W Pylyshyn. Connectionism and cognitive architecture: A critical
661 analysis. *Cognition*, 28(1-2):3–71, 1988.
- 662 [109] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuyte-
663 laars. Memory aware synapses: Learning what (not) to forget. In *Proc. European Conf. on
664 Computer Vision (ECCV)*, pages 144–161, Munich, Germany, September 2018.
- 665 [110] Zhizhong Li and Derek Hoiem. Learning without forgetting. In *Proc. European Conf. on
666 Computer Vision (ECCV)*, pages 614–629, Amsterdam, Netherlands, October 2016.
- 667 [111] Adam Paszke et al. Pytorch: An imperative style, high-performance deep learning library.
668 In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, pages 8026–8037,
669 Vancouver, Canada, December 2019.
- 670 [112] Róbert Csordás, Kazuki Irie, and Jürgen Schmidhuber. The devil is in the detail: Simple tricks
671 improve systematic generalization of transformers. In *Proc. Conf. on Empirical Methods in
672 Natural Language Processing (EMNLP)*, Punta Cana, Dominican Republic, November 2021.
- 673 [113] Kazuki Irie, Imanol Schlag, Róbert Csordás, and Jürgen Schmidhuber. Improving baselines in
674 the wild. In *Workshop on Distribution Shifts, NeurIPS*, Virtual only, 2021.
- 675 [114] James Requeima, Jonathan Gordon, John Bronskill, Sebastian Nowozin, and Richard E. Turner.
676 Fast and flexible multi-task classification using conditional neural adaptive processes. In *Proc.
677 Advances in Neural Information Processing Systems (NeurIPS)*, pages 7957–7968, Vancouver,
678 Canada, December 2019.
- 679 [115] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross
680 Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle.
681 Meta-dataset: A dataset of datasets for learning to learn from few examples. In *Int. Conf. on
682 Learning Representations (ICLR)*, Addis Ababa, Ethiopia, April 2020.
- 683 [116] Jürgen Schmidhuber. One big net for everything. *Preprint arXiv:1802.08864*, 2018.

684 [117] Alex Graves, Marc G. Bellemare, Jacob Menick, Rémi Munos, and Koray Kavukcuoglu.
685 Automated curriculum learning for neural networks. In *Proc. Int. Conf. on Machine Learning*
686 (*ICML*), pages 1311–1320, Sydney, Australia, August 2017.

687 A Experimental Details

688 A.1 Continual and Meta-learning Terminologies

689 We review the following classic terminologies of continual learning and meta-learning used through-
690 out this paper.

691 **Continual learning.** “Domain-incremental learning (DIL)” and “class-incremental learning (CIL)”
692 are two classic settings in continual learning [104, 28, 6]. They differ as follows. Let M and N
693 denote positive integers. Consider continual learning of M tasks where each task is an N -way
694 classification. In the DIL case, a model has an N -way output classification layer, i.e., the class ‘0’ of
695 the first task shares the same weights as the class ‘0’ of the second task, and so on. In the CIL case, a
696 model’s output dimension is $N * M$; the class indices of different tasks are not shared, neither are the
697 corresponding weights in the output layer. In our experiments, all CIL models have the $(N * M)$ -way
698 output from the first task (instead of progressively increasing the output size). In this work, we skip
699 the third variant called “task-incremental learning” which assumes that we have access to the task
700 identity as an extra input, as it makes the CL problem almost trivial. CIL is typically reported to be
701 the hardest setting among them.

702 **Meta-learning.** We need to introduce “meta-training” and “meta-test” terminologie since each of
703 these phases involve “training/test” processes within itself. Each of them requires the corresponding
704 training and test examples. We refer to these as “meta-training training/test examples”, and “meta-test
705 training/test examples” following the terminology of Beaulieu et al. [30]. While these are rather
706 “heavy” terminologies, they are unambiguous and help avoid potential confusions. In both phases,
707 our sequence-processing neural net observes a sequence of (meta-training or meta-test) training
708 examples—each consisting of input features and a correct label—, and the resulting states of the
709 sequence processor (i.e., weights in the case of SRWM) are used to make predictions on (meta-
710 training or meta-test) test examples—input features presented to the model without its label. During
711 the meta-training phase, we modify the trainable parameters of the meta-learner through gradient
712 descent minimizing the meta-learning loss function (using backpropagation through time). During
713 meta-testing, no human-designed optimization for weight modification is used anymore; the SRWMs
714 modify their own weights following their own learning rules defined as their forward pass (Eqs. 1-3).
715 In connection with the now-popular in-context learning [96], we also refer to a (meta-training or
716 meta-test) training-example sequence as *context*.

717 A.2 Datasets

718 For classic image classification datasets such as MNIST [25], CIFAR10 [27], and FashionMNIST
719 (FMNIST; Xiao et al. [26]) we refer to the original references for details.

720 For Omniglot [23], we use Vinyals et al. [21]’s 1028/172/432-split for the train/validation/test set, as
721 well as their data augmentation methods using rotation of 90, 180, and 270 degrees. Original images
722 are grayscale hand-written characters from 50 different alphabets. There are 1632 different classes
723 with 20 examples for each class.

724 Mini-ImageNet contains color images from 100 classes with 600 examples for each class. We use the
725 standard train/valid/test class splits of 64/16/20 following [22].

726 FC100 is based on CIFAR100 [27]. 100 color image classes (600 images per class, each of size
727 32×32) are split into train/valid/test classes of 60/20/20 [24].

728 The “5-datasets” dataset [32] consists of 5 datasets: CIFAR10, MNIST, FashionMNST, SVNH [105],
729 and notMNIST [106].

730 Split-CIFAR100 is also based on CIFAR100. The standard setting splits CIFAR100 into 10 10-way
731 classification tasks.

732 **Meta-train/test sequence construction procedure.** We use `torchmeta` [107] which provides
 733 common few-shot/meta learning settings for these datasets to sample and construct their meta-
 734 train/test datasets. The construction of “meta-training training” sequences for an N -way classification,
 735 using a dataset containing C classes works as follows; for each sequence, we sample N random
 736 but distinct classes out of C ($N < C$). The resulting classes are re-labelled such that each class is
 737 assigned to one out of N distinct random label index which is unique to the sequence. For each of
 738 these N classes, we sample K examples. We randomly order these $N * K$ examples to obtain a
 739 sequence. Each such a sequence “simulates” an unknown task the model has to learn.

740 A.3 Training Details & Hyper-Parameters

741 We use the same model and training hyper-parameters in all our experiments. All hyper-parameters
 742 are summarized in Table 5. We use the Adam optimizer with the standard Transformer learning rate
 743 warmup scheduling [59]. The vision backend is the classic 4-layer convolutional NN of Vinyals
 744 et al. [21]. Most configurations follow those of Irie et al. [19]; except that we initialize the ‘query’
 745 sub-matrix in the self-referential weight matrix using a normal distribution with a mean value of 0
 746 and standard deviation of $0.01/\sqrt{d_{\text{head}}}$ while other sub-matrices use an std of $1/\sqrt{d_{\text{head}}}$ (motivated
 747 by the fact that a generated query vector is immediately multiplied with the same SRWM to produce
 748 a value vector). For any further details, we’ll refer the readers to our public code we’ll release upon
 749 acceptance. We conduct our experiments using a single V100-32GB, 2080-12GB or P100-16GB
 750 GPUs, and the longest single training run takes about one day.

Table 5: Hyper-parameters.

Parameters	Values
Number of SRWM layers	2
Total hidden size	256
Feedforward block multiplier	2
Number of heads	16
Batch size	16 or 32

751 A.4 Evaluation Procedure

752 For evaluation on few-shot learning datasets (i.e., Omniglot, Mini-Imagenet and FC100), we use 5
 753 different sets consisting of 32 K random test episodes each, and report mean and standard deviation.

754 For evaluation on standard datasets, we use 5 different random support sets for in-context learning,
 755 and evaluate on the entire test set. We report the corresponding mean and standard deviation across
 756 these 5 evaluation runs.

757 For the Split-MNIST experiment, we do 10 meta-testing runs to compute the mean and standard
 758 deviation as the baseline models are also trained for 10 runs in Hsu et al. [6] (see other details in
 759 Appendix A.7).

760 A.5 ACL Objectives with More Tasks

761 We can straightforwardly extend the 2-task version of ACL presented in Sec. 3 to more tasks. In
 762 the 3-task case (we denote the three tasks as **A**, **B**, and **C**) used in Sec. 4.3, the objective function
 763 contains six terms. Following three terms are added to Eq. 4:

$$- \left(\log(p(y_{\text{test}}^C | \mathbf{x}_{\text{test}}^C; \mathbf{W}_{A,B,C})) + \log(p(y_{\text{test}}^B | \mathbf{x}_{\text{test}}^B; \mathbf{W}_{A,B,C})) + \log(p(y_{\text{test}}^A | \mathbf{x}_{\text{test}}^A; \mathbf{W}_{A,B,C})) \right)$$

764 This also naturally extends to the 5-task loss used in the Split-MNIST experiment (Table 3). As
 765 one can observe, the number of terms rapidly/quadratically increases with the number of tasks.
 766 Nevertheless, computing these loss terms isn’t immediately impractical because they essentially just
 767 require forwarding the network for one step, for many independent inputs/images. This can be heavily
 768 parallelized as a batch operation. While this can be a concern when scaling up more, a natural open
 769 research question is whether we really need all these terms in the case we have many more tasks.

Table 6: Impact of the choice of meta-validation datasets. Classification accuracies (%) on three datasets: **Split-CIFAR-10**, **Split-Fashion MNIST** (Split-FMNIST), and **Split-MNIST** in the **domain-incremental** setting (we omit “Split-” in the second column). “OOB” denotes “out-of-the-box”. “mImageNet” here refers to mini-ImageNet.

Meta-Finetune Datasets	Meta-Validation Sets	Meta-Test on Split-		
		MNIST	FMNIST	CIFAR-10
None (OOB: 2-task ACL; Sec. 4.1)	Omniglot + mImageNet	72.2 ± 0.9	75.6 ± 0.7	65.3 ± 1.6
Omniglot	MNIST	84.3 ± 1.2	78.1 ± 1.9	55.8 ± 1.2
	FMNIST	81.6 ± 1.3	90.4 ± 0.5	59.5 ± 2.1
	CIFAR10	75.2 ± 2.3	78.2 ± 0.9	63.4 ± 1.4
Omniglot + mImageNet	MNIST	76.6 ± 1.4	85.3 ± 1.1	66.2 ± 1.1
	FMNIST	73.2 ± 2.3	89.9 ± 0.6	66.6 ± 0.7
	CIFAR10	76.3 ± 3.0	88.1 ± 1.3	68.6 ± 0.5

770 Ideally, we want these models to ‘systematically generalize’ to more tasks even when they are trained
771 with only a handful of them [108]. This is an interesting research question on generalization to be
772 studied in a future work.

773 A.6 Auxiliary 1-shot Learning Objective

774 In practice, instead of training the models only for “15-shot learning,” we also add an auxiliary loss
775 for 1-shot learning. This naturally encourages the models to learn in-context from the first examples.

776 A.7 Details of the Split-MNIST experiment

777 Here we provide details of the Split-MNIST experiments presented in Sec. 4 and Table 3.

778 Split-MNIST is obtained by transforming the classic 10-class single-task MNIST dataset into a
779 sequence of 5 tasks by partitioning the 10 classes into 5 groups/pairs of two classes each, in a fixed
780 order from 0 to 9 (i.e., grouping 0/1, 2/3, 4/5, 6/7, and 8/9). Regarding the difference between
781 domain/class-incremental settings, we refer to Appendix A.1.

782 The baseline methods presented in Table 3 include: standard SGD and Adam optimizers, Adam
783 with the L2 regularization, elastic weight consolidation [9] and its online variant [10], synaptic
784 intelligence [11], memory aware synapses [109], learning without forgetting (LwF; Li and Hoiem
785 [110]). For these methods, we directly take the numbers reported in Hsu et al. [6] for the 5-task
786 domain/class-incremental settings.

787 For the 2-task class incremental setting, we use Hsu et al. [6]’s code to train the correspond models
788 (the number for LwF is currently missing as it is not implemented in their code base; we plan to add
789 the corresponding/missing entry in Table 3 for the final version of this paper).

790 Finally we also evaluate two meta-CL baselines: Online-aware Meta-Learning (OML; Javed and
791 White [29]) and Generative Meta-Continual Learning (GemCL; Banayeezade et al. [31]). OML is
792 a MAML-based meta-learning approach. We note that as reported by Javed and White [29] in their
793 public code repository; after some critical bug fix, the performance of their OML matches that of
794 Beaulieu et al. [30] (which is a direct application of OML to another model architecture). Therefore,
795 we focus on OML as our main MAML-based baseline. We take the out-of-the-box model (meta-
796 trained for Omniglot, with a 1000-way output) made publicly available by Javed and White [29]. We
797 evaluate the corresponding model in two ways. In the first, ‘out-of-the-box’ case, we take the meta-
798 pre-trained model and only tune its meta-testing learning rate (which is done by Javed and White [29]
799 even for meta-testing in Omniglot). We find that this setting does not perform very well; in the other
800 case (‘optimized # meta-testing iterations’), we additionally tune the number of meta-test training
801 iterations. We’ve done a grid search of the meta-test learning rate in $3 * \{1e^{-2}, 1e^{-3}, 1e^{-4}, 1e^{-5}\}$
802 and the number of meta-test training steps in $\{1, 2, 5, 8, 10\}$ using a meta-validation set based on an
803 MNIST validation set (5 K held-out images from the training set); we found the learning rate of $3e^{-4}$
804 and 8 steps to consistently perform the best in all our settings. We’ve also tried it ‘with’ and ‘without’

Table 7: Impact of the number of in-context examples. Classification accuracies (%) on **Split-MNIST** in the 2-task and 5-task class-incremental learning (CIL) settings and the 5-task domain-incremental learning (DIL) setting. For ACL models, we use the same number of examples for meta-validation as for meta-training. According to Banayeezade et al. [31], GeMCL is meta-trained with the 5-shot setting but meta-validated in the 15-shot setting.

Number of Examples		DIL		CIL 2-task		CIL 5-task	
Meta-Train/Valid	Meta-Test	GeMCL	ACL	GeMCL	ACL	GeMCL	ACL
5	5	-	84.1 ± 1.2	-	93.4 ± 1.2	-	74.6 ± 2.3
	15	-	83.8 ± 2.8	-	94.3 ± 1.9	-	65.5 ± 4.0
15	5	62.2 ± 5.2	83.9 ± 1.0	87.3 ± 2.5	93.6 ± 1.7	71.7 ± 2.5	76.7 ± 3.6
	15	63.8 ± 3.8	84.5 ± 1.6	91.2 ± 2.8	96.0 ± 1.0	79.0 ± 2.1	84.3 ± 1.2

805 the standard mean/std normalization of the MNIST dataset; better performance was achieved without
 806 such normalization (which is in fact consistent as they do not normalize the Omniglot dataset for
 807 their meta-training/testing). Their performance on the 5-task class-incremental setting is somewhat
 808 surprising/disappointing (since generalization from Omniglot to MNIST is typically straightforward,
 809 at least, in common non-continual few-shot learning settings; see, e.g., Munkhdalai and Yu [51]). At
 810 the same time, to the best of our knowledge, OML-trained models have not been tested in such a
 811 condition in prior work; from what we observe, the publicly available out-of-the-box model might
 812 be overtuned for Omniglot/Mini-ImageNet or the frozen ‘representation network’ is not ideal for
 813 generalization. We note that the sensitivity of these MAML-based methods [29, 30] w.r.t. meta-test
 814 hyper-parameters has been also noted by Banayeezade et al. [31]; these are characteristics of
 815 hand-crafted learning algorithms that we want to avoid with learned learning algorithms.

816 We use code and a pre-trained model (trained on Omniglot) made public by Banayeezade et al.
 817 [31] for the GeMCL baseline (see also Table 7); like our method, GeMCL also do not require any
 818 special tuning at test-time.

819 Our out-of-the-box ACL models (trained on Omniglot and Mini-ImageNet) do not require any
 820 tuning at meta-test time. Nevertheless, we’ve checked the effect of the number of meta-test training
 821 examples (5 vs. 15; 15 is the number used in meta-training); we found the consistent number, i.e., 15,
 822 to work better than 5. For the version that is meta-finetuned using the 5-task ACL objective (using
 823 only the Omniglot dataset), we use 5 or 15 examples for both meta-train and meta-test training (see an
 824 ablation study in Table 7). To obtain a sequence of 5 tasks, we simply sample 5 tasks from Omniglot
 825 (in principle, we should make sure that different tasks in the same sequence have no class overlap;
 826 in practice, our current implementation simply randomly draws 5 independent tasks from Omniglot).

827 A.8 Details of the Split-CIFAR100 and 5-datasets experiment using ViT

828 As we described in Sec. 4, for the experiments on Split-CIFAR100 and 5-datasets, following
 829 Wang et al. [33, 34], we use ViT-B/16 pre-trained on ImageNet [76] which is available through
 830 torchvision [111]. In this experiments, we resize all images to 3x224x224 and feed them to the
 831 ViT. We remove the output layer of the ViT, and use its 768-dimensional feature from the penultimate
 832 layer as the image encoding. The self-referential component which is added to this encoder has the
 833 same architecture (2 layers, 16 heads) as the rest of the paper (see all hyper-parameters in Table 5)
 834 All ViT parameters are frozen during meta-training.

835 B Extra Experimental Results

836 B.1 Ablation Studies on the Meta-validation Dataset

837 Here we conduct ablation studies on the choice of meta-validation sets to select model checkpoints. In
 838 general, when dealing with out-of-domain generalization, the choice of validation procedures to select
 839 final model checkpoints plays a crucial role in the evaluation of the corresponding method [112, 113].
 840 The out-of-the-box models are chosen based on the average meta-validation performance on the
 841 validation set corresponding to the few-shot learning datasets used in meta-training: Omniglot and

Table 8: Meta-testing on sequences that are longer than those from meta-training. Classification accuracies (%) on 5-task **Split-FMNIST** and 5-task **Split-MNIST** in the **domain-incremental** settings. The model is the one finetuned with 5-task ACL loss using Omniglot as the meta-finetuning set and FMNIST as the meta-validation set (i.e., the numbers in the top part of the table are taken from Table 6). In the first column, “Split-FMNIST, Split-MNIST” indicates continual learning of 5 Split-FMNIST tasks followed by 5 tasks of Split-MNIST (and “Split-MNIST, Split-FMNIST” is the opposite order). Performance is measured at the end of the entire sequence.

Meta-Test Training Task Sequence	# Tasks	Meta-Test Test Tasks	
		Split-FMNIST	Split-MNIST
Split-FMNIST	5	90.4 ± 0.5	-
Split-MNIST	5	-	81.6 ± 1.3
Split-FMNIST, Split-MNIST	10	79.3 ± 2.7	74.3 ± 0.9
Split-MNIST, Split-FMNIST	10	78.1 ± 3.1	78.5 ± 1.7

Table 9: Classification accuracies (%) on 5-task 2-way Split-Omniglot. Mean/std is computed over 10 meta-test runs.

Method	Domain Incremental	Class Incremental
GeMCL	64.6 ± 9.2	97.4 ± 2.7
ACL	92.3 ± 0.4	96.8 ± 0.8

842 mini-ImageNet (or Omniglot, mini-ImageNet, and FC100 in the case of 3-task ACL), independently of
843 any potential meta-test datasets. In contrast, in the meta-finetuning process of Table 3, we selected our
844 model checkpoint by meta-validation on the MNIST validation dataset (we held out 5 K images from
845 the training set). Here we evaluate ACL models meta-finetuned for the “5-task domain-incremental
846 binary classification” on three Split-‘X’ tasks where ‘X’ is MNIST, FashionMNIST (FMNIST) or
847 CIFAR-10 for various choices of meta-validation sets (in each case we hold out 5 K images from
848 the corresponding training set). In addition, we also evaluate the effect of meta-finetuning datasets
849 (Omniglot only v. Omniglot and mini-ImageNet). Table 6 shows the results (we use 15 meta-training
850 and meta-testing examples except for the Omniglot-finetuned/MNIST-validated model from Table 3
851 which happens to be configured with 5 examples; this will be fixed in the final version). Effectively,
852 meta-validation on the matching validation set is useful. Also, meta-finetuning only on Omniglot is
853 beneficial for the performance on MNIST when meta-validated on MNIST or FMNIST. However,
854 importantly, we emphasize that our ultimate goal is not to obtain a model that is specifically tuned for
855 certain datasets; we aim at building models that generally work well across a wide range of tasks
856 (ideally on any tasks); in fact, several existing works in the few-shot learning literature evaluate
857 their methods in such settings (see, e.g., Requeima et al. [114], Bronskill et al. [78], Triantafillou
858 et al. [115]). This also goes hand-in-hand with scaling up ACL (our current model is tiny; see
859 hyper-parameters in Table 5; the vision component is also a shallow ‘Conv-4’ net) and various other
860 considerations on self-improving continual learners (see, e.g., Schmidhuber [116]), such as automated
861 curriculum learning [117].

862 B.2 Performance on Split-Omniglot

863 Here we report the performance of the models used in the Split-MNIST experiment (Sec. 4.3) on
864 “in-domain” 5-task 2-way Split-Omniglot. Table 9 shows the result. Performance is very similar
865 between our ACL and the baseline GeMCL on this task in the class incremental setting, unlike on
866 Split-MNIST (Table 3) where we observe a larger performance gap between these same models. Here
867 we also include the “domain incremental” setting for the sake of completeness but note that GeMCL
868 is not originally trained for this setting.

Table 10: 5-way classification accuracies using 15 examples for each class for each task in the context. 2-task models are meta-trained on Omniglot and Mini-ImageNet, while 3-task models are in addition meta-trained on FC100. ‘A, B’ in ‘Context/Train’ column indicates that models sequentially observe meta-test training examples of Task A then B; evaluation is only done at the end of the sequence. “no ACL” is the baseline 2-task models trained without the ACL loss.

Meta-Testing Tasks		Number of Meta-Training Tasks		
Context/Train	Test	2 (no ACL)	2	3
A: MNIST-04	A	71.1 ± 4.0	75.4 ± 3.0	89.7 ± 1.6
B: CIFAR10-04	B	51.5 ± 1.4	51.6 ± 1.3	55.3 ± 0.9
C: MNIST-59	C	65.9 ± 2.4	63.0 ± 3.3	76.1 ± 2.0
D: FMNIST-04	D	52.8 ± 3.4	54.8 ± 1.3	59.2 ± 4.0
	Average	60.3	61.2	70.1
A, B	A	43.7 ± 2.3	81.5 ± 2.7	88.0 ± 2.2
	B	49.4 ± 2.4	50.8 ± 1.3	52.9 ± 1.2
	Average	46.6	66.1	70.5
A, B, C	A	26.5 ± 3.2	64.5 ± 6.0	82.2 ± 1.7
	B	32.3 ± 1.7	50.8 ± 1.2	50.3 ± 2.0
	C	56.5 ± 8.1	33.7 ± 2.2	44.3 ± 3.0
	Average	38.4	49.7	58.9
A, B, C, D	A	24.6 ± 2.7	64.3 ± 4.8	78.9 ± 2.3
	B	20.6 ± 2.3	47.5 ± 1.0	49.2 ± 1.3
	C	38.5 ± 4.4	32.7 ± 1.9	45.4 ± 3.9
	D	36.1 ± 2.5	31.2 ± 4.9	30.1 ± 5.8
	Average	30.0	43.9	50.9

869 B.3 Effect of Number of In-Context Examples

870 Table 7 shows an ablation study on the number of examples used for meta-training and meta-testing
871 on the Split-MNIST task. We observe that for an ACL model trained only with 5 examples during
872 meta-training, more examples (15 examples) provided during meta-testing is not beneficial. In fact,
873 they even largely hurt in certain cases (see the last column); this is one form of “length generalization”
874 problem. When the number of meta-training examples is consistent with the one used during
875 meta-testing, the 15-example case consistently outperforms the 5-example one.

876 B.4 Effect of Number of Tasks in the ACL Loss

877 Table 10 provides the complete results discussed in Sec. 4.3 under “Evaluation on diverse task
878 domains”.

879 B.5 Further Discussion on Limitations

880 Here we provide further discussion and experimental results on the limitations of our approach as a
881 learned algorithm.

882 **Domain generalization.** As a data-driven learned algorithm, the domain generalization capability
883 is a typical limitation as it depends on the meta-trained data. Certain results we presented above
884 are representative of this limitation. In particular, in Table 6, the model meta-trained/finetuned on
885 Omniglot using Split-MNIST as meta-validation set do not perform well on Split-CIFAR10. While
886 meta-training and meta-validating on a larger/diverse set of datasets may be an immediate remedy
887 to obtain more robust ACL models, we note that since ACL is also a “continual meta-learning”
888 algorithm (Sec. 5), an ideal ACL model should also continually incorporate and learn from more data
889 during potentially lifelong meta-testing; we leave such an investigation for future work.

890 **Length generalization.** We already qualitatively observed the limited length generalization capabil-
891 ity in Table 10 (meta-trained with up to 3 tasks and meta-tested with up to 4 tasks). Here we provide
892 one more experiment evaluating ACL models meta-trained for 5 tasks on a concatenation of two
893 5-task Split-MNIST and Split-FMNIST tasks (resulting in 10 tasks). Table 8 shows the results. Again,

894 while the model does not completely break, increasing the number of tasks to 10 rapidly degrades the
895 performance compared to the 5-task setting the model is meta-trained for. Similarly, its performance
896 on the Split-Omniglot domain incremental setting (Sec. B.2) degrades with increased numbers of
897 tasks: accuracies for 5, 10 and 20 tasks are $92.3\% \pm 0.4$, $82.0\% \pm 0.4$ and $67.6\% \pm 1.1$ respectively.
898 As noted in Sec. 5, this is a general limitation of sequence processing neural networks, and there is a
899 potential remedy for this limitation (meta-training on more tasks and “context carry-over”) which we
900 leave for future work.

901 **B.6 A Comment on Meta-Generalization**

902 We also note that in general, “unseen” datasets do not necessarily imply that they are harder tasks than
903 “in-domain” test sets; when meta-trained on Omniglot and mini-ImageNet, meta-generalization on
904 “unseen” MNIST is easier (the accuracy is higher) than on the “in-domain” test set of mini-ImageNet
905 with heldout/unseen classes (compare Tables 1 and 2).

906 **NeurIPS Paper Checklist**

907 **1. Claims**

908 Question: Do the main claims made in the abstract and introduction accurately reflect the
909 paper's contributions and scope?

910 Answer: [Yes]

911 Justification: We accurately state contributions and scope of the work in the abstract and
912 introduction.

913 Guidelines:

- 914 • The answer NA means that the abstract and introduction do not include the claims
915 made in the paper.
- 916 • The abstract and/or introduction should clearly state the claims made, including the
917 contributions made in the paper and important assumptions and limitations. A No or
918 NA answer to this question will not be perceived well by the reviewers.
- 919 • The claims made should match theoretical and experimental results, and reflect how
920 much the results can be expected to generalize to other settings.
- 921 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
922 are not attained by the paper.

923 **2. Limitations**

924 Question: Does the paper discuss the limitations of the work performed by the authors?

925 Answer: [Yes]

926 Justification: We discuss limitations of our method in Sec. 4 and 5.

927 Guidelines:

- 928 • The answer NA means that the paper has no limitation while the answer No means that
929 the paper has limitations, but those are not discussed in the paper.
- 930 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 931 • The paper should point out any strong assumptions and how robust the results are to
932 violations of these assumptions (e.g., independence assumptions, noiseless settings,
933 model well-specification, asymptotic approximations only holding locally). The authors
934 should reflect on how these assumptions might be violated in practice and what the
935 implications would be.
- 936 • The authors should reflect on the scope of the claims made, e.g., if the approach was
937 only tested on a few datasets or with a few runs. In general, empirical results often
938 depend on implicit assumptions, which should be articulated.
- 939 • The authors should reflect on the factors that influence the performance of the approach.
940 For example, a facial recognition algorithm may perform poorly when image resolution
941 is low or images are taken in low lighting. Or a speech-to-text system might not be
942 used reliably to provide closed captions for online lectures because it fails to handle
943 technical jargon.
- 944 • The authors should discuss the computational efficiency of the proposed algorithms
945 and how they scale with dataset size.
- 946 • If applicable, the authors should discuss possible limitations of their approach to
947 address problems of privacy and fairness.
- 948 • While the authors might fear that complete honesty about limitations might be used by
949 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
950 limitations that aren't acknowledged in the paper. The authors should use their best
951 judgment and recognize that individual actions in favor of transparency play an impor-
952 tant role in developing norms that preserve the integrity of the community. Reviewers
953 will be specifically instructed to not penalize honesty concerning limitations.

954 **3. Theory Assumptions and Proofs**

955 Question: For each theoretical result, does the paper provide the full set of assumptions and
956 a complete (and correct) proof?

957 Answer: [NA]

958 Justification: This is not a theoretical paper.

959 Guidelines:

- 960 • The answer NA means that the paper does not include theoretical results.
- 961 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 962 referenced.
- 963 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 964 • The proofs can either appear in the main paper or the supplemental material, but if
- 965 they appear in the supplemental material, the authors are encouraged to provide a short
- 966 proof sketch to provide intuition.
- 967 • Inversely, any informal proof provided in the core of the paper should be complemented
- 968 by formal proofs provided in appendix or supplemental material.
- 969 • Theorems and Lemmas that the proof relies upon should be properly referenced.

970 4. Experimental Result Reproducibility

971 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

972 perimental results of the paper to the extent that it affects the main claims and/or conclusions

973 of the paper (regardless of whether the code and data are provided or not)?

974 Answer: [Yes]

975 Justification: We provide experimental details in the main text and details in Appendix A.

976 We also provide our code in the supplemental material.

977 Guidelines:

- 978 • The answer NA means that the paper does not include experiments.
- 979 • If the paper includes experiments, a No answer to this question will not be perceived
- 980 well by the reviewers: Making the paper reproducible is important, regardless of
- 981 whether the code and data are provided or not.
- 982 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 983 to make their results reproducible or verifiable.
- 984 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 985 For example, if the contribution is a novel architecture, describing the architecture fully
- 986 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 987 be necessary to either make it possible for others to replicate the model with the same
- 988 dataset, or provide access to the model. In general, releasing code and data is often
- 989 one good way to accomplish this, but reproducibility can also be provided via detailed
- 990 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 991 of a large language model), releasing of a model checkpoint, or other means that are
- 992 appropriate to the research performed.
- 993 • While NeurIPS does not require releasing code, the conference does require all submis-
- 994 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 995 nature of the contribution. For example
- 996 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 997 to reproduce that algorithm.
- 998 (b) If the contribution is primarily a new model architecture, the paper should describe
- 999 the architecture clearly and fully.
- 1000 (c) If the contribution is a new model (e.g., a large language model), then there should
- 1001 either be a way to access this model for reproducing the results or a way to reproduce
- 1002 the model (e.g., with an open-source dataset or instructions for how to construct
- 1003 the dataset).
- 1004 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 1005 authors are welcome to describe the particular way they provide for reproducibility.
- 1006 In the case of closed-source models, it may be that access to the model is limited in
- 1007 some way (e.g., to registered users), but it should be possible for other researchers
- 1008 to have some path to reproducing or verifying the results.

1009 5. Open access to data and code

1010 Question: Does the paper provide open access to the data and code, with sufficient instruc-

1011 tions to faithfully reproduce the main experimental results, as described in supplemental

1012 material?

1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064

Answer: [Yes]

Justification: We provide experimental details in the main text and details in Appendix A. We also provide our code in the supplemental material. The data we use are classic datasets which are publicly available.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide experimental details in the main text and details in Appendix A. We also provide our code in the supplemental material.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All our results are mean/std computed using 10 evaluation seeds.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- 1065 • The assumptions made should be given (e.g., Normally distributed errors).
- 1066 • It should be clear whether the error bar is the standard deviation or the standard error
- 1067 of the mean.
- 1068 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 1069 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 1070 of Normality of errors is not verified.
- 1071 • For asymmetric distributions, the authors should be careful not to show in tables or
- 1072 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 1073 error rates).
- 1074 • If error bars are reported in tables or plots, The authors should explain in the text how
- 1075 they were calculated and reference the corresponding figures or tables in the text.

1076 8. Experiments Compute Resources

1077 Question: For each experiment, does the paper provide sufficient information on the com-
1078 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1079 the experiments?

1080 Answer: [Yes]

1081 Justification: We provide compute resource related information in Appendix A.

1082 Guidelines:

- 1083 • The answer NA means that the paper does not include experiments.
- 1084 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1085 or cloud provider, including relevant memory and storage.
- 1086 • The paper should provide the amount of compute required for each of the individual
1087 experimental runs as well as estimate the total compute.
- 1088 • The paper should disclose whether the full research project required more compute
1089 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1090 didn't make it into the paper).

1091 9. Code Of Ethics

1092 Question: Does the research conducted in the paper conform, in every respect, with the
1093 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1094 Answer: [NA]

1095 Justification: We do not have anything to report.

1096 Guidelines:

- 1097 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1098 • If the authors answer No, they should explain the special circumstances that require a
1099 deviation from the Code of Ethics.
- 1100 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1101 eration due to laws or regulations in their jurisdiction).

1102 10. Broader Impacts

1103 Question: Does the paper discuss both potential positive societal impacts and negative
1104 societal impacts of the work performed?

1105 Answer: [NA]

1106 Justification: Our work does not have any such impacts.

1107 Guidelines:

- 1108 • The answer NA means that there is no societal impact of the work performed.
- 1109 • If the authors answer NA or No, they should explain why their work has no societal
1110 impact or why the paper does not address societal impact.
- 1111 • Examples of negative societal impacts include potential malicious or unintended uses
1112 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1113 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1114 groups), privacy considerations, and security considerations.

- 1115 • The conference expects that many papers will be foundational research and not tied
1116 to particular applications, let alone deployments. However, if there is a direct path to
1117 any negative applications, the authors should point it out. For example, it is legitimate
1118 to point out that an improvement in the quality of generative models could be used to
1119 generate deepfakes for disinformation. On the other hand, it is not needed to point out
1120 that a generic algorithm for optimizing neural networks could enable people to train
1121 models that generate Deepfakes faster.
- 1122 • The authors should consider possible harms that could arise when the technology is
1123 being used as intended and functioning correctly, harms that could arise when the
1124 technology is being used as intended but gives incorrect results, and harms following
1125 from (intentional or unintentional) misuse of the technology.
- 1126 • If there are negative societal impacts, the authors could also discuss possible mitigation
1127 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1128 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1129 feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

1130 Question: Does the paper describe safeguards that have been put in place for responsible
1131 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1132 image generators, or scraped datasets)?

1133 Answer: [NA]

1134 Justification: Our work does not imply any such risks.

1135 Guidelines:

- 1136
- 1137 • The answer NA means that the paper poses no such risks.
 - 1138 • Released models that have a high risk for misuse or dual-use should be released with
1139 necessary safeguards to allow for controlled use of the model, for example by requiring
1140 that users adhere to usage guidelines or restrictions to access the model or implementing
1141 safety filters.
 - 1142 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1143 should describe how they avoided releasing unsafe images.
 - 1144 • We recognize that providing effective safeguards is challenging, and many papers do
1145 not require this, but we encourage authors to take this into account and make a best
1146 faith effort.

12. Licenses for existing assets

1147 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1148 the paper, properly credited and are the license and terms of use explicitly mentioned and
1149 properly respected?

1150 Answer: [Yes]

1151 Justification: Our codebase includes certain publicly available code. The corresponding
1152 license files are included in the supplemental material.

1153 Guidelines:

- 1154
- 1155 • The answer NA means that the paper does not use existing assets.
 - 1156 • The authors should cite the original paper that produced the code package or dataset.
 - 1157 • The authors should state which version of the asset is used and, if possible, include a
1158 URL.
 - 1159 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - 1160 • For scraped data from a particular source (e.g., website), the copyright and terms of
1161 service of that source should be provided.
 - 1162 • If assets are released, the license, copyright information, and terms of use in the
1163 package should be provided. For popular datasets, paperswithcode.com/datasets
1164 has curated licenses for some datasets. Their licensing guide can help determine the
1165 license of a dataset.
 - 1166 • For existing datasets that are re-packaged, both the original license and the license of
1167 the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The documentations of our code are included in the readme file in the supplemental material.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: We do not have such experiments.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: We do not have such experiments.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.