# LIFELONG LEARNING THE TASK-PARAMETER RELA-TIONSHIPS FOR KNOWLEDGE TRANSFER

#### **Anonymous authors**

Paper under double-blind review

## Abstract

In this work, we focus on knowledge transfer in the lifelong learning setting. We propose a lifelong learner that exploits the similarities between the optimal weight spaces of tasks, thereby enabling knowledge transfer across tasks in a continual learning setting. To characterize the "*task-parameter relationships*", we propose a metric called adaptation rate integral (ARI) that measures the expected adaptation rate over a finite number of steps for a (task, parameter) pair. These task-parameter relationships are learned by an auxiliary network trained on guided explorations of parameter space. The learned auxiliary network is then used to heuristically select the best parameter sets on seen tasks, which are consolidated using a hypernetwork. We show that the proposed approach can transfer knowledge to new tasks without any increase in overall model capacity, while naturally mitigating catastrophic forgetting.

## **1** INTRODUCTION

An embodied agent that operates in the unbounded partially-observable real-world with its diversity of tasks, must possess the ability to acquire knowledge and skills continually. An embodied and therefore finite agent, cannot however feasibly grow its skill set as a set of disjoint abilities learnt afresh for each task. For complex tasks, this form of learning would be prohibitively expensive. Conversely, learning a *single* skill that can be modulated and maintained to be effective on the large diversity of tasks is unfeasible (a broad proof is provided in the No Free Lunch Theorems; Wolpert & Macready (1997)). A *skill <u>set</u>* is therefore necessary - one that allows the agent to benefit from the *redundancies* between similar tasks that require similar skills. Skills optimized for previous tasks may reasonably be leveraged to improve learning quality and efficiency on unseen yet related tasks.

Such mechanisms are widely reflected in fundamental neuro-cognitive processes in human learning and memory (Tomita et al. (2021); Caramazza & Shelton (1998). The re-use of neural circuitry across diverse cognitive tasks appears to be a central organizational principle of the brain Anderson (2010). Shared structures and properties across tasks and environments are exploited. Similar tasks are solved rapidly and more effectively by re-using acquired skills, while novel experiences are prioritized within the learning bandwidths (Caramazza & Shelton (1998); Coutanche & Thompson-Schill (2015)). A continual consolidation of knowledge occurs (Ritvo et al. (2019); Wilson & McNaughton (1994); Alvarez & Squire (1994)) aimed at retaining salient, widely and frequently re-usable knowledge/skills. This is particularly evident in memory organization, knowledge consolidation, and the role of novelty and forgetting in the memory (Tomita et al. (2021); Ritvo et al. (2019)). The aim of our work is to incorporate these mechanisms within the continual learning (CL) setting.

Modern machine learning (ML) algorithms still struggle to replicate this human ability to learn continually. The phenomenon of *catastrophic forgetting* (McCloskey & Cohen (1989); French (1999)) - the difficulty in the retention of old information when new information is acquired, is commonly observed in training continually across multiple domains. This is a fundamental consequence of the transfer-interference trade-off (Riemer et al. (2018)) - for a *singular* finite network continually adapted to a shifting distribution, catastrophic forgetting is inevitable. The field of *continual learning* (CL) has therefore focused largely on mitigating catastrophic forgetting - with limited success in enabling knowledge transfer (Rebuffi et al. (2017); Lopez-Paz & Ranzato (2017); Kirkpatrick et al. (2017); Zenke et al. (2017)) due in part to the single network constraints generally applied in CL (Ramesh & Chaudhari (2021)).



Figure 1: Illustration of shared structure in the optimal weight manifolds of tasks: Parameters are shown to be adapted from the initialization set  $\Theta_0 \sim p(\Theta)$  (centered, in grey) for tasks  $\tau_1, \tau_2, \tau_3$  in sequence. We illustrate how parameters trained on tasks with a greater vicinity of their optimal weight manifolds (W\*), may be transferred with a higher adaptation rate (illustrated as a shorter path). In this work, we aim to model the relationship between parameters in the weight space and the adaptation path lengths (and vicinity) to the optimal weight manifold of tasks.

In contrast, an important characteristic of the human learning process is this elevated 'quality of convergence' onto new tasks that share observed structures - a knowledge transfer that enables a higher rate and quality of learning (varying with the degree and type of shared task structures). In a lifelong learning setting, an efficient use of continually acquired knowledge necessitates a benefit to such a quality of adaptation on new, albeit related tasks. In the limit, a continually learning agent that most rapidly converges on a new task, effectively already understands the task. This manner of optimization benefits an agent's supervision/data requirements in learning a new task.

We motivate this re-use of acquired skills for improved adaptation quality on new tasks, as an important trait of a generally learning agent. A faster 'rate of adaptation' (or convergence) in ML can be formulated as a shorter adaptation trajectory to the basin of convergence. This can be achieved by an initialization with a higher affinity to the basin or by enforcing a more direct trajectory (Flennerhag et al. (2018)) through auxiliary feedback/constraints. In contrast, work on mitigating catastrophic forgetting is focused on retaining parameters within the (optimal) weight manifolds learned for a task (Riemer et al. (2018)). In this work we begin first, with the hypothesis that the shared structure between tasks may be captured in the optimal weight spaces of the (networks trained on the) tasks. And second, to allow skill re-use for a new task, a heuristic to search the explored (or previously learned) weights set for the most promising weights is required.

To this end, we make the following contributions: i) To measure 'adaptation quality' for a (parameter  $\theta$ , task  $\tau$ ), we formulate a metric - *adaptation rate integral* (ARI) that captures the convergence rate and performance of a parameter  $\theta$  trained on a task  $\tau$ . ii) We then develop a heuristic that can efficiently search the *observed* space of parameters (base model weights trained on previous tasks) by estimating the adaptation quality of any (parameter  $\theta$ , task  $\tau$ ) pair. iii) Finally, in order to efficiently store the parameters learned from previous tasks, we employ a meta model that stores all observed parameters in its representation. This approach does not require any additional model capacity compared to a single (base) model. We leverage these contributions to incrementally learn and explore the observed space of parameters, improving the the degree of knowledge transfer as well as the retained accuracy of the continual learner. We show that the benefits to knowledge transfer come with no increase in overall model capacity, while mitigating catastrophic forgetting naturally.

## 2 PRELIMINARIES

In a supervised learning setting, a hypothesis  $h : \mathcal{X} \to \mathcal{Y}, h \in \mathcal{H}$  is learned on the input and label spaces  $\mathcal{X}$ , and  $\mathcal{Y}$  ( $\mathcal{H}$  is the hypothesis space). The learner  $h \equiv f(\cdot; \theta) : \mathcal{X} \to \mathcal{Y}$  can be defined as a ML model f parameterized by  $\theta$ . The input and label spaces are related by a joint probability distribution  $P(x, y) \mid x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . The empirical risk minimization (ERM) principle formulates this problem of learning h as a minimization of the population risk  $e_P(h) = \mathcal{P}(h(x) \neq y)$ . For a



Figure 2: Overall Method: During training, a set of parameters are trained on the given task and their corresponding ARIs are calculated. After filtering using the ARI Maximization algorithm, the selected parameters are stored in the hypernetwork for future retrieval. At Inference, from the parameter set  $\Theta_M$  stored in the hypernetwork, the parameter  $\theta_i$  with the maximum predicted ARI for the given inference task  $\tau_k$  is fetched and evaluated.

finite sample set  $S \equiv (X, Y) = \{x_i, y_i\}_{i \in [1,N]}$  of size N seen by the learner h, the empirical risk is  $\hat{e}_{S}(h) = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}_{h(x_i) \neq y_i}$  where  $\mathbf{1}_z$  is an indicator function. The constrained ERM optimization problem is therefore defined as:

ERM Optimization: 
$$\min_{\forall h \in \mathcal{H}} \hat{e}_{\boldsymbol{S}}(h).$$
 (1)

While in a typical continual learning (CL) setting (Gupta et al. (2020b)), the learner observes a sequence of M tasks  $[\tau_1, \tau_2, ..., \tau_M]$ ;  $\tau_i \sim P(\tau)$  sampled from a distribution of tasks  $P(\tau)$ . Here, a task  $\tau_i$  is defined as a set of  $N_i$  input, label pairs:  $S_i \equiv (X_i, Y_i) = \{x_n, y_n\}_{n \in [1, N_i]}$ . Typically, the objective is to find parameters  $\theta$  that minimize the *cumulative loss* on all m seen tasks  $\tau_{[1:m]}$ , while having limited access to data  $S_i$  from previous tasks  $\tau_i$  (i < m). The CL objective is:

$$\min_{\theta} \sum_{i=1}^{m} \mathbb{E}_{\boldsymbol{S}_{i}} \left[ l_{i}(f(\boldsymbol{X}_{i};\theta),\boldsymbol{Y}_{i}) \right] = \min_{\theta} \mathbb{E}_{\boldsymbol{S}_{[1:m]}} \left[ L_{m}(f(\boldsymbol{X}_{[1:m]};\theta),\boldsymbol{Y}_{[1:m]}) \right]$$
(2)

where  $l_i$  is the loss on task  $\tau_i$  and  $L_m$  is the sum of all m task-specific losses  $L_m = \sum_{i=1}^m l_i$ .

#### 3 Methodology

Our proposed continual learning algorithm operates in the meta space of parameters (where we refer to the weights of a 'base' model trained on each task as a parameter ( $\theta$ ). It uses a heuristic to incrementally explore the meta space and consolidate parameters in a knowledge base that are estimated to perform well on the distribution of observed tasks.

#### 3.1 PROPOSED CL ALGORITHM

The proposed method relies on i) a meta model  $F_{\Phi}$  that efficiently stores parameters as a knowledge base ( $\Theta_{\mathcal{M}}$ ) and ii) a small auxiliary network  $\hat{q}$  that estimates the 'quality of convergence' of a parameter  $\theta$  trained on a task  $\tau$ .

**Training** consists of an Exploration and a Consolidation phase. Let's take  $\Theta_0$  as a set of randomly initialized parameters. i) *Exploration*: A combined training set  $\Theta_{Tr}$  of parameters from the knowledge base and initialization set ( $\Theta_{Tr} \subseteq \Theta_M \cup \Theta_0$ ) is selected and adapted on task  $\tau$ . A measure of the convergence quality of all  $\Theta_{Tr}$  parameters on task  $\tau$  is calculated and retained. ii) *Consolidation*: The calculated convergence qualities from all observed (parameter, task) pairs is used to train the auxiliary network  $\hat{q}$ . Finally, the subset of the combined knowledge base and adapted training set ( $\Theta_M \cup \Theta_{Tr}$ ), that has the highest estimated quality of convergence on the observed distribution of tasks, is retained in the updated knowledge base  $\Theta_M$ .

**Testing** consists of selecting the parameter from the knowledge base, with the highest estimated quality of convergence for the task, fine-tuning and evaluating it.

A training budget Q can be enforced by simply using the auxiliary network  $\hat{q}$  to select the top candidate networks for training on a task. Further, as the architectures are the same, data streams (identical &) independent and only the initializations differ, training of the parameter set is fully parallelizable with negligible overhead to training time complexity.

#### 3.2 CONVERGENCE QUALITY MEASURE: ADAPTATION RATE INTEGRAL (ARI)

To formalize the 'quality of convergence' of a given parameter to the task space, we introduce the Adaptation Rate Integral (ARI) metric. In a typical supervised learning setting, solving the ERM optimization (eq. 1) involves iteratively modifying an initial hypothesis h(0) in discrete steps  $h(t+1) = h(t) + \alpha(t).\delta(t)$ . Here,  $t \in [0, T]$  denotes the step index, h(t) represents the hypothesis at step t,  $\alpha(t)$  denotes the learning rate, and  $\delta(t)$  is typically an estimate of the gradient of the empirical risk. We define the *adaptation rate integral* or ARI, as simply the step-averaged area under the curve of  $(1 - \hat{e}_{\mathbf{S}}(h(t))) \forall t \in [0, T]$ . For a continuous step space with infinitesimally small step sizes dt,

$$\psi (\text{ARI}) = 1 - \frac{1}{T} \int_0^T \hat{e}_{\boldsymbol{S}}(h(t)) dt$$
(3)

If the optimizer is allowed a maximum of T steps to solve eq. 1, the adaptation rate integral (ARI) is maximized when the step averaged empirical risk on sample set S over T steps is minimum. Ideally, a learning algorithm should converge to a good quality solution (one that achieves global minimum of the empirical risk) in the fewest possible steps. The ARI value attempts to measure both the quality of the converged solution and how fast this solution can be reached. The challenge lies in how to estimate the value of  $\psi$  given an initial hypothesis h(0) and the optimizer, without explicitly constructing the complete adaptation trajectory h.

#### 3.2.1 PROPOSED CL OBJECTIVE: ARI Maximization

Let the initial parameters of function  $f(\cdot; \theta)$  adapting on task  $\tau$  with a loss function l be  $\theta_0$ . Then, an SGD operator  $U(\theta_0)$  acting on parameter  $\theta_0$  is defined as follows:  $U(\theta_0) = \theta_1 = \theta_0 - \alpha \nabla_{\theta_0} l(\theta_0)$ . Thus, the parameters  $\theta_T$  obtained after T adaptation steps can be composed as  $U_T(\theta_0) = U \circ U \cdots \circ U(\theta_0) = \theta_T$ . In the CL setting, for a model f with initial parameters  $\theta_0$  adapted on task  $\tau_i$  over T steps, the *adaptation rate integral* is defined as:

$$\psi(\theta_0, \tau_i, T) = 1 - \frac{1}{T} \sum_{t=0}^{T} \hat{e}_{\boldsymbol{S}_i}(f(\boldsymbol{X}_i; \theta_t), \boldsymbol{Y}_i) ; U_t(\theta_0) = \theta_t$$
(4)

Now, in order to learn the optimal initial parameters  $\theta_0^* \in \Theta$  that maximizes the rate of adaptation on a distribution of tasks  $\tau \sim \mathcal{P}(\tau)$ , the learning objective using *ARI (adaptation rate integral)* maximization becomes:

$$\max_{\theta_0 \in \Theta} \mathbb{E}_{\tau_i \sim \mathcal{P}(\tau)} \left[ \psi(\theta_0, \tau_i, T) \right] = \max_{\theta_0 \in \Theta} \mathbb{E}_{\tau_i \sim \mathcal{P}(\tau)} \left[ 1 - \frac{1}{T} \sum_{t=0}^T \hat{e}_{\boldsymbol{S}_i}(f(\boldsymbol{X}_i; \theta_t), \boldsymbol{Y}_i) \right]$$
(5)

For the continual learning setting, the ARI maximization objective can thus be expressed as:

ARI-Maximization Objective for CL: 
$$\max_{\theta_0 \in \Theta} \sum_{i=1}^m \psi(\theta_0, \tau_i, T) = \max_{\theta_0 \in \Theta} \Psi(\theta_0, \tau_{[1:m]}, T)$$
(6)

where the learner has a budget of T iterations/steps to converge to each task  $\tau$ , and  $\Psi(\theta_0, \tau_{[1:m]}, T) = \sum_{i=1}^{m} \psi(\theta_0, \tau_i, T)$ . In simple terms, given a new task  $\tau$  where a learner's initial parameter  $\theta_0$  is evolved across T adaptation steps  $\theta_T = U_T(\theta_0)$ , we would like to find the optimal  $\theta_0 \in \Theta$  such that expected step-averaged empirical risk on  $\tau$  across the T steps is minimized (Eq. 5). The proposed objective explicitly attempts to maximize the quality of adaptation of the learner to the observed task set as well as new tasks, which are sampled from the same task distribution.

#### 3.3 OVERALL LIFELONG LEARNER & COMPONENTS

As a lifelong learner, our proposed approach operates as a compressed knowledge base  $\Theta_{\mathcal{M}} = \{\theta^1, \theta^2, \cdots, \theta^{\mathcal{N}}\}\$  of explored base model parameters that are optimized for maximal 'adaptation quality' to the distribution of tasks  $\tau \sim P(\tau)$  (while requiring the parameter budget of a single base model). Typically, models are often over-parameterized in continual learning with wider layers to mitigate the degree of interference. We circumvent the problem of interference, by searching the space of base model parameters and maintaining the subset of parameters with maximal expected ARI on the observed task distribution.

Algorithm 1 Lifelong Learning Algorithm

**Input**: Observed Tasks:  $\tau_{[1:i-1]} = \{\tau_1, \tau_2, \cdots, \tau_{i-1}\}$ , memory budget  $\mathcal{B}$ , training budget  $\mathcal{Q}$ **Require:** ARI Estimator:  $\hat{q}_{\phi}$ , meta model  $F_{\Phi}$  and embedding vectors  $E_{\mathcal{M}}$ , model and task encoders:  $\mathcal{E}_{\text{param}}, \mathcal{E}_{\text{task}}$  that generate model, task encodings  $\eta_{\theta}, \eta_{\tau}$  resp., Generated initialization set  $\Theta_0 =$  $\{\theta_0^1, \theta_0^2, .., \theta_0^\mathcal{N}\}$ **Ensure:**  $|\Theta_M| \leq \mathcal{B}$  and  $|\Theta_{Tr}| \leq \mathcal{Q}$ Knowledge Base $\Theta_{\mathcal{M}}$ , Training Set  $\Theta_{Tr}$ , Buffer  $S = \{\}$ 1: for Training on task  $\tau_i$  do  $\Theta_{\mathcal{M}} \leftarrow \{F_{\Phi}(e_i) \forall e_i \in E_{\mathcal{M}}\}$  Generate from meta-model 2: ▷ Exploration Select training set  $\Theta_{Tr} \mid \Theta_{Tr} \subseteq (\Theta_{\mathcal{M}} \cup \Theta_0)$  as, 3:  $\Theta_{Tr} \leftarrow \arg\min_{\Theta_{Tr}, |\Theta_{Tr}| \leq \mathcal{Q}} \sum_{\theta_0 \in \Theta_{Tr}} \hat{q}_{\phi}(\boldsymbol{\eta}_{\theta_0}, \boldsymbol{\eta}_{\tau_i})$ Train  $\Theta_{Tr}$  using SGD to obtain  $\hat{\Theta_{Tr}} \mid \hat{\Theta_{Tr}} \leftarrow \{U_T(\theta_{Tr}^j)\}, \forall \theta_{Tr}^j \in \Theta_{Tr}$ 4:  $\psi(\theta_{Tr}, \tau_i, T) \leftarrow \text{Calculate true adaptation rate integral } \forall \theta_{Tr} \in \Theta_{Tr} \text{ on } \tau_i \text{ (Eq. 4)}$ 5: Add true ARIs to buffer  $S = S \cup \{\psi(\theta_{Tr}, \tau_i, T), \eta_{\theta_{Tr}}, \eta_{\tau_i}\} \forall \theta_{Tr} \in \Theta_{Tr}$ 6: ▷ Consolidation Train ARI Estimator  $\hat{q}_{\phi}$  on S, 7:  $\phi \leftarrow \phi - \nu \nabla_{\phi} \sum_{\{\psi, \eta_{\theta}, \eta_{\tau}\} \in S} ||\psi - \hat{q}_{\phi}(\eta_{\theta}, \eta_{\tau})||_{2}^{2} \quad (\text{Eq. 7})$ Consolidate  $\Theta_{\mathcal{M}}$  to retain parameter subset that optimizes: 8:  $\begin{aligned} \Theta_{\mathcal{M}} &\leftarrow \arg\min_{\Theta_{\mathcal{M}} \subseteq (\hat{\Theta}_{r} \cup \Theta_{\mathcal{M}}), |\Theta_{\mathcal{M}}| \leq \mathcal{B}} \sum_{\theta_{j} \in \Theta_{\mathcal{M}}} \hat{q}_{\phi}(\theta_{j}, \tau_{[1:i]}, T) \\ \text{Train meta-model } F_{\Phi} \text{ on } \Theta_{\mathcal{M}} : \\ \Phi &\leftarrow \Phi - \nu' \nabla_{\Phi} \sum_{\theta_{k} \in \Theta_{\mathcal{M}}} ||\theta_{k} - F_{\Phi}(e_{k})||_{2}^{2} \\ e_{k} &\leftarrow e_{k} - \nu' \nabla_{e_{k}} ||\theta_{k} - F_{\Phi}(e_{k})||_{2}^{2}; \forall k \in [1, n(\Theta_{\mathcal{M}})] \end{aligned}$ 9:  $E_{\mathcal{M}} = \{e_k\}_{k=1}^{n(\Theta_{\mathcal{M}})}$ 10: 11: end for ▷ Inference 12: for Inference on task  $\tau_m$  do  $\Theta_{\mathcal{M}} \leftarrow \{F_{\Phi}(e_i) \forall e_i \in E_{\mathcal{M}}\}$  Generate from meta-model 13:  $\theta_* \leftarrow \arg \max_{\theta_* \in \Theta_M} \hat{q}_{\phi}(\theta_*, \tau_m, T)$ 14:  $\theta_* \leftarrow$  Finetune  $\theta_*$  using exemplar memory D, then infer on  $\tau_m$ 15: 16: end for

The proposed overall lifelong learning algorithm is detailed in Algorithm 1. Below, we detail the various components of the CL learner:

#### 3.3.1 ARI ESTIMATOR

As searching exhaustively through a large space of parameters is prohibitively expensive, we require a heuristic to search the parameter space efficiently for parameters that most effectively adapt to observed and new tasks. Such a heuristic would effectively characterize the model-task relationships.

We learn a small-capacity auxiliary network  $\hat{q}_{\phi}$  that estimates the true ARI  $\psi(\theta, \tau, T)$  for a (parameter, task) pair  $(\theta, \tau), \theta \in \Theta, \tau \sim P(\tau)$ . Thus, when a new task  $\tau$  is observed, an estimation of adaptation rate of the parameter  $\theta$  to the task  $\tau$  can be generated off hand by the auxiliary network, without requiring a full training pass over the task. However, this first requires a projection of parameters and tasks into a shared space  $\gamma$ , which is informative of a parameter  $\theta$ 's adaptation rates for a task and a task  $\tau$ 's characteristics. We learn two auxiliary encoders - a parameter encoder  $\mathcal{E}_{\text{param}}$  and a task encoder  $\mathcal{E}_{\text{task}}$ . Following the approach in Jeong et al. (2021b), we extract functional signatures of the parameters  $g_{\alpha}(\theta)$  and the task  $g_{\beta}(\tau)$  and pass them to the encoders to generate the respective embeddings  $\eta_{\theta} = \mathcal{E}_{\text{param}}(g_{\alpha}(\theta))$  and  $\eta_{\tau} = \mathcal{E}_{\text{task}}(g_{\beta}(\tau))$ . These embeddings are then used by the ARI estimator network  $\hat{q}_{\phi}$ , which learns to regress the true ARI measured when parameter  $\theta$  is adapted on task  $\tau$  for T steps via solving the following optimization objective.

$$\min_{\phi} ||\psi(\theta_j, \tau_i, T) - \hat{q}_{\phi}(\boldsymbol{\eta}_{\theta_j}, \boldsymbol{\eta}_{\tau_i})||_2^2 \quad \forall \ \theta_j \in \Theta_{\mathcal{M}}, \tau_i \in \tau_{[1:t]},$$
(7)

where  $\Theta_{\mathcal{M}}$  is the set of explored parameters (knowledge base) that has been adapted on the m observed tasks  $\tau[1:m]$ . Note that  $\phi$  includes the parameters of the two encoders  $\mathcal{E}_{\text{param}}$  and  $\mathcal{E}_{\text{task}}$  as well as the regression model  $\hat{q}$ . Given an embedding space that captures degrees of structural/domain similarity across the task distribution, the auxiliary network  $\hat{q}_{\phi}$  generalizes to predict the rate of adaptability to new tasks similar to the observed set. The training objective (Eq. 7) is continually 'consolidated' and naturally becomes more stable as number of observed tasks increases.

#### 3.4 KNOWLEDGE BASE REPRESENTATION USING HYPERNETWORKS

In order to efficiently maintain the set of explored parameters (knowledge base)  $\Theta_{\mathcal{M}}$  without explicitly storing them, we leverage a hypernetwork meta model  $F_{\Phi}$ . Hypernetworks (Ha et al. (2016)) are meta models that learn to map embedding vectors to parameters (Von Oswald et al. (2019)), and can be thought of as weight generators. They have been show to efficiently retain a large parameter set with no decrease in the evaluated performance of the parameters (Von Oswald et al. (2019)).

We train the hypernetwork  $F_{\Phi}$  to map the knowledge base  $\Theta_{\mathcal{M}}$  to a set of learned embedding vectors  $\{e_i\}_{i=1}^{n(\Theta)}$  that can be thought of as indices for parameters in  $\Theta_{\mathcal{M}}$ . Thus, given a parameter index *i*, the parameter  $\theta_i$  can be be readily generated using the hypernetwork as  $\theta_i = F_{\Phi}(e_i)$ . This effectively reduces the storage complexity of the knowledge base from  $|\Theta_{\mathcal{M}}|$  to  $|F_{\Phi}|$ , where  $|F_{\Phi}| \cong |\theta|, \theta \in \Theta_{\mathcal{M}}$ . We discuss the details of the hyper-network architecture and hyperparameters in Section 4.4.

#### 4 EXPERIMENTS

Following earlier continual learning literature Lopez-Paz & Ranzato (2017); Chaudhry et al. (2018b); Rebuffi et al. (2017), and owing to compute restrictions (involved in training a number of base models parallely) we conduct experiments and ablations using 4 smaller scale continual learning benchmarks - Split-MNIST Chaudhry et al. (2018a), Permuted-MNIST Zenke et al. (2017), Split CIFAR-10 Zenke et al. (2017), Split CIFAR-100 Rebuffi et al. (2017). We evaluate against prominent baselines - GEM Lopez-Paz & Ranzato (2017), iCaRL Rebuffi et al. (2017), EWC Kirkpatrick et al. (2017).

#### 4.1 DATASETS

We briefly describe the datasets employed in our experiments:

**Split-MNIST** and **Permuted-MNIST**: Consecutive classes from the MNIST dataset (LeCun (1998)) are paired and presented as 5 incremental tasks in Split-MNIST (Chaudhry et al. (2018a)). In Permuted-MNIST, each task is a unique spatial permutation of the original MNIST data. 10 permutations of MNIST resulting in 10 tasks, are generated. While Split-MNIST represents a set of incremental tasks where the number of classes is expanded within a shared data distribution, the Permuted-MNIST dataset represents a set of tasks that may be sampled from disjoint data distributions.

Following the protocol in (Lopez-Paz & Ranzato (2017)), 1000 images per task are considered for training, while the entire test set is considered for evaluation. **Split CIFAR-10** and **Split CIFAR-100**: Classes are grouped from CIFAR-10 (Krizhevsky et al. (2009)) and CIFAR-100 (Krizhevsky et al. (2009)) to generate independent incremental tasks. 5 and 10 classes are grouped within CIFAR-10, CIFAR-100 to generate 10 and 20 tasks respectively. As followed in (Lopez-Paz & Ranzato (2017)), 2500 images per task are considered for training, while the entire test set is considered for evaluation.

#### 4.2 BASELINES

We briefly describe the baseline formulations below: Single model - A naive baseline where a single model is trained continually on all tasks. GEM (Lopez-Paz & Ranzato (2017)) - an approach that leverages the exemplar memory to explicitly bound the model's loss on previous task samples. iCaRL (Rebuffi et al. (2017)) - a class-incremental learning approach that uses a nearest-mean-of-exemplars classification strategy, with a knowledge distillation loss over feature representations of past tasks to limit catastrophic forgetting. EWC (Kirkpatrick et al. (2017)) - A regularization based approach to mitigate catastrophic forgetting, where the parameters crucial to the performance on previous tasks, as measured by the Fisher Information Matrix (FIM), are not modified. Specifically, we consider the approach in EWC++ (Chaudhry et al. (2018a)) and online EWC (Schwarz et al. (2018)) to calculating FIM as a moving average.

#### 4.3 EVALUATION METRICS

Following previous works (Lopez-Paz & Ranzato (2017); Chaudhry et al. (2018a)), we consider the metrics of Average Accuracy (ACC) and Backward Transfer (BWT) for all experiments. After the model is trained on task  $t_i$ , it is evaluated on the test-sets of all tasks in task-set T, resulting in a matrix  $R \in \mathcal{R}^{T \times T}$ . Each element  $R_{ij}$  is the test-classification accuracy of the model on task  $t_j$  after learning on examples from task  $t_i$ . Average accuracy (ACC  $\in [0, 100]$ ) after learning task T can be defined as: ACC =  $\frac{1}{T} \sum_{i=1}^{T} R_{T,i}$ . Average Forgetting (F  $\in [-100, 100]$ ) after learning the  $T^{th}$  task can be defined as: F =  $\frac{1}{T-1} \sum_{i=1}^{T-1} (\max_{l \in \{1,...,T-1\}} R_{l,i} - R_{T,i})$ .

#### 4.4 EXPERIMENTAL SETTINGS

We follow the protocols in Lopez-Paz & Ranzato (2017); Riemer et al. (2018) for our choice of experimental settings and build on the implementation provided by Gupta et al. (2020b) and Von Oswald et al. (2019) to implement our baselines and hypernetwork meta-models respectively. For MNIST experiments, we follow Lopez-Paz & Ranzato (2017) and use a two layer, 100-neuron each, fully-connnected neural network with ReLU activation for the MNIST datasets. As a meta model, we use a fully-connected two-hidden layer ([100, 100]) chunked hypernetwork (Von Oswald et al. (2019)) with a chunk size of 200 and embedding vectors of size 8. The meta model contains 59,668 weights in comparison to a single base model with 89,400 weights. All baselines including ours are implemented using a single-headed base network. For CIFAR, we use a modified version of the ResNet18 (He et al. (2016)) with one-third the feature maps across all layers, as in Lopez-Paz & Ranzato (2017). As a meta model, we use a larger hypernetwork with structured chunking that internally maintains 6 smaller composite two hidden layer ([25, 25]) hypernetworks, for a total of 166,610 weights (compared to 181,495 for a single base model). The baselines for CIFAR are all multi-headed and task-aware, while base models of our method are trained as single-headed networks (for each task). The hypernetworks are all trained with embedding vectors of size 8. Similar to baselines Lopez-Paz & Ranzato (2017), we maintain a small exemplar memory D of size 200 for MNIST experiments and 400 for CIFAR. For our ARI estimator, we follow Jeong et al. (2021a) in generating task or parameter signatures by the activations of a pre-trained ResNet18 on samples per class (of each task) or the activations of the parameters on random gaussian noise, respectively. The activations are normalized and padded to a size of 2048, before being projected to 128 dimensions. The ARI estimator used is a simple two hidden layer ([100, 100]) network with ReLU activations. We enforce a maximum size of 20 base parameters in our knowledge base, and a training budget of 10 base models trained per task. Our initialization set contains 3 random generated parameters.

$\text{Datasets} \rightarrow$	Split MNIST		Permuted MNIST		Split CIFAR-10		Split CIFAR-100	
Methods $\downarrow$	ACC (†)	F (↓)	ACC (†)	$F\left(\downarrow ight)$	ACC (†)	F (↓)	ACC (†)	F (↓)
Single EWC Kirkpatrick et al. (2017)	$\overline{30.2 \pm 1.3}$ $43.9 \pm 1.1$	$96.1 \pm 1.1$ $94.4 \pm 1.2$	$63.4 \pm 1.2$ $72.36 \pm 3.8$	$14.2 \pm 1.1$ $11.3 \pm 1.1$	$71.3 \pm 1.9$ $59.1 \pm 2.2$	$\overline{13.5 \pm 3.9}$ $10.1 \pm 1.9$	$34.3 \pm 1.2$ $34.1 \pm 0.1$	$21.4 \pm 1.9$ 19.6 ± 1.6
GEM Lopez-Paz & Ranzato (2017) iCaRL Rebuffi et al. (2017) Ours	$\begin{array}{c} 81.1 \pm 3.2 \\ 85.4 \pm 0.9 \\ 88.7 \pm 0.1 \end{array}$	$19.5 \pm 2.2$ $1.1 \pm 0.3$ $0.5 \pm 0.1$	$82.8 \pm 1.2$ - $84.2 \pm 1.3$	$1.4 \pm 1.2$ $0.2 \pm 0.0$	$75.3 \pm 1.1$ $66.2 \pm 1.2$ $81.9 \pm 0.2$	$4.9 \pm 1.3$ $3.9 \pm 2.1$ $0.3 \pm 0.2$	$44.3 \pm 3.1$ $31.9 \pm 1.9$ $48.8 \pm 1.2$	$\begin{array}{c} 0.9 \pm 0.3 \\ 2.9 \pm 0.9 \\ 0.3 \pm 1.1 \end{array}$

Table 1: Average Accuracy (A) and Average Forgetting (F) for baselines across four datasets.

Table 2: Ablation Study: Evaluating the ARI Estimator

$Datasets \rightarrow$	Split MNIST		Permuted	d MNIST	Split CIFAR-10		Split CIFAR-100	
Methods $\downarrow$	ACC (†)	F (↓)	ACC (†)	F (↓)	ACC $(\uparrow)$	F (↓)	ACC (†)	F (↓)
Random Selection ARI Estimator (Ours*)	$\left  \begin{array}{c} \overline{60.1 \pm 0.3} \\ 88.7 \pm 0.1 \end{array} \right $	$40.1 \pm 5.1$ $0.5 \pm 0.1$	$45.4 \pm 8.1$ $84.2 \pm 1.3$	$68.2 \pm 3.9 \\ 0.2 \pm 0.0$	$49.4 \pm 4.1$ $81.9 \pm 0.2$	$20.2 \pm 3.4$ $0.3 \pm 0.2$	$\overline{31.4 \pm 3.3}$ $48.8 \pm 1.2$	$     \begin{array}{r}       24.4 \pm 2.9 \\       0.3 \pm 1.1     \end{array}   $

### 4.5 TRAINING DETAILS

To train the base model across all baselines we use the Adam (Kingma & Ba (2014)) optimizer set with an initial learning rate of 0.001, weight decay of 0.001 and a batch size of 10, similar to baseline methods (Lopez-Paz & Ranzato (2017); Rebuffi et al. (2017); Kirkpatrick et al. (2017)). We also utilize a small buffer S to store the collected (ARI, task & parameter encoding) tuples. For our method, the ARI values for each parameter, task pair are calculated based on empirical steps required till convergence. The ARI estimator is trained for 1000 epochs with a batch size of 250 (parameter & task embeddings, ARI) tuples. Finally, the meta-models are trained using the SGD optimizer till convergence (an MSE error of 1e-3), and perform a single pass of the entire parameter in one batch.

## 5 RESULTS AND DISCUSSION

The performance comparison of our approach against baseline approaches on the four standard continual learning benchmarks (Split-MNIST, Permuted-MNIST, Split CIFAR-10, Split CIFAR-100) is shown in Table 1. The performance for PermutedMNIST is not reported for iCaRL (Rebuffi et al. (2017)) as the approach does not support Domain-Incremental methodology.

Our approach is observed to achieve near consistent gains over the baselines across the datasets, with the performance gains being higher for CIFAR - the more complex dataset amongst the benchmarks. Our approach achieves a higher rate of accuracy on majority of the tasks, with a consistently lower forgetting metric. ICaRL (Rebuffi et al. (2017)) achieves lower forgetting on most datasets compared to the baselines, which we believe is attributable to it's auxiliary knowledge distillation loss (Hinton et al. (2015)) that constraints any significant change in logits from previous tasks. Nevertheless, our approach outperforms ICaRL (Rebuffi et al. (2017)) even on the forgetting measure on Split CIFAR-10. We do not consider A-GEM (Chaudhry et al. (2018b)) due to it's marginal improvements over GEM (Lopez-Paz & Ranzato (2017)). We observe that the removal of experience replay causes a significant decrease in the overall accuracy, as well as an increase in the degree of forgetting across all datasets. This is to be expected given that the base models used are simple, and therefore need to be replayed stored examples.

#### 5.1 ABLATION: TRIVIAL HEURISTIC TO PICK CANDIDATE MODELS

In our ablation study, we evaluate the benefit from our ARI estimator. Table 2 shows the performance of our approach with the proposed ARI estimator along with the performance of the same method with a random selection heuristic. In this random selection heuristic, the parameters from the hypernetwork are selected during training and inference using a random heuristic. We observe a clear and large reduction in the average accuracy as well as an increase in the degree of forgetting observed. Without the proposed ARI estimator based heuristic, the approach collapses as it fails to select the appropriate parameters for the test tasks. Random selection also prevents retaining the optimal weights in the hypernetwork (memory).

## 5.2 MEMORY COMPLEXITY

Across all benchmarks, we utilize the same or smaller base models as compared to established baselines (Lopez-Paz & Ranzato (2017)). Only the parameters of the hypernetwork themselves are retained in memory, and the weights of the trained base models are discarded after training. Thus, the total parameter size stored in memory remains constant and equal to the parameters of the hypernetwork (which are also 1.4x (in case of MNIST) and equal or less than the parameters of a single base model in other benchmarks).

# 6 RELATED LITERATURE

## 6.1 NEUROSCIENCE

Vital to the process of learning and 'memorization' in humans is the continual familiarity-based modification of input instances within the Hippocampus (Kumaran et al. (2016); Caramazza & Shelton (1998); Wilson & McNaughton (1994); Alvarez & Squire (1994)). Input patterns that are similar to a familiar/stored pattern are modified to either be more similar (pattern completion) or differentiated (pattern separation) to those stored patterns (Kumaran et al. (2016)). New memory is allocated if the input instance is sufficiently distinct from stored instances. In case there is high overlap between input instance and one of the stored instances, the input instance is modified to be closer to the matched stored instance.

## 6.2 CONTINUAL LEARNING

Over the last few years, several directions of work have attempted to address these issues of catastrophic forgetting and beneficial knowledge transfer in a continual learning setting. Methods retain past knowledge either by replaying stored (Rebuffi et al. (2017); Lopez-Paz & Ranzato (2017)) or generated samples (Shin et al. (2017); van de Ven & Tolias (2018)), regularizing task-specific weights (Kurle et al. (2019); Titsias et al. (2019); Chaudhry et al. (2018a)), or scaling parameters to account for new tasks (Mallya & Lazebnik (2018); Serra et al. (2018); Diethe et al. (2019)). In attempting to explicitly prevent the learning dynamics that cause the loss of task-specific knowledge, approaches have inevitably focused on modelling the transfer-interference trade-off between gradients of different tasks (Saha et al. (2021); Deng et al. (2021); Lopez-Paz & Ranzato (2017)). In Ramesh & Chaudhari (2021), authors introduce the idea of building a growing zoo of small capacity multi-tasking models, where synergistic tasks share models, enabling transfer of knowledge between them.

## 6.3 META LEARNING APPROACHES FOR CL

Online-aware Meta Learning (OML) Finn et al. (2019) introduced the application of Meta Learning approaches to the lifelong learning setting. A meta-objective was used to learn the task distribution in an offline manner, which could then be leveraged for efficient online continual learning. More recent Meta-learning approaches to continual learning such as MER (Riemer et al. (2018)) and La-MAML (Gupta et al. (2020b)) leverage gradient alignments to enforce compatibility of tasks within a finite capacity. MER (Riemer et al. (2018)) enforces gradient alignment between observed and future tasks using replay while La-MAML (Gupta et al. (2020a)) incrementally modulates parameter-specific learning rates based on gradient alignment across tasks to reduce forgetting. Recent works also investigate orthogonal setups in which a learning agent uses all previously seen data to adapt quickly to an incoming stream of data, thereby ignoring the problem of catastrophic forgetting.

# 7 CONCLUSION

In this work, we motivate this focus on quality of adaptation to improve knowledge transfer in the lifelong learning setting. We propose to represent task-model relationships as the expected adaptation rate of a (model, task) pair. In order to leverage (model, task) relationships, we replace a single network used for lifelong learning with an equivalent set of small-capacity networks such that the overall model capacity is conserved. We show that the proposed approach can transfer knowledge to new tasks without any increase in overall model capacity, while naturally mitigating catastrophic forgetting.

#### REFERENCES

- Pablo Alvarez and Larry R Squire. Memory consolidation and the medial temporal lobe: a simple network model. *Proceedings of the national academy of sciences*, 91(15):7041–7045, 1994.
- Michael L Anderson. Neural reuse: A fundamental organizational principle of the brain. *Behavioral* and brain sciences, 33(4):245–266, 2010.
- Alfonso Caramazza and Jennifer R Shelton. Domain-specific knowledge systems in the brain: The animate-inanimate distinction. *Journal of cognitive neuroscience*, 10(1):1–34, 1998.
- Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 532–547, 2018a.
- Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018b.
- Marc N Coutanche and Sharon L Thompson-Schill. Rapid consolidation of new knowledge in adulthood via fast mapping. *Trends in cognitive sciences*, 19(9):486–488, 2015.
- Danruo Deng, Guangyong Chen, Jianye Hao, Qiong Wang, and Pheng-Ann Heng. Flattening sharpness for dynamic gradient projection memory benefits continual learning. Advances in Neural Information Processing Systems, 34:18710–18721, 2021.
- Tom Diethe, Tom Borchert, Eno Thereska, Borja Balle, and Neil Lawrence. Continual learning in practice. *arXiv preprint arXiv:1903.05202*, 2019.
- Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning. In *International Conference on Machine Learning*, pp. 1920–1930. PMLR, 2019.
- Sebastian Flennerhag, Pablo G Moreno, Neil D Lawrence, and Andreas Damianou. Transferring knowledge across learning processes. arXiv preprint arXiv:1812.01054, 2018.
- Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3 (4):128–135, 1999.
- Gunshi Gupta, Karmesh Yadav, and Liam Paull. Look-ahead meta learning for continual learning. Advances in Neural Information Processing Systems, 33:11588–11598, 2020a.
- Gunshi Gupta, Karmesh Yadav, and Liam Paull. La-maml: Look-ahead meta learning for continual learning. *arXiv preprint arXiv:2007.13904*, 2020b.
- David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. arXiv preprint arXiv:1609.09106, 2016.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv* preprint arXiv:1503.02531, 2015.
- Wonyong Jeong, Hayeon Lee, Geon Park, Eunyoung Hyung, Jinheon Baek, and Sung Ju Hwang. Task-adaptive neural network search with meta-contrastive learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021a.
- Wonyong Jeong, Hayeon Lee, Gun Park, Eunyoung Hyung, Jinheon Baek, and Sung Ju Hwang. Task-adaptive neural network search with meta-contrastive learning, 2021b.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114 (13):3521–3526, 2017.

Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

- Dharshan Kumaran, Demis Hassabis, and James L McClelland. What learning systems do intelligent agents need? complementary learning systems theory updated. *Trends in cognitive sciences*, 20(7): 512–534, 2016.
- Richard Kurle, Botond Cseke, Alexej Klushyn, Patrick van der Smagt, and Stephan Günnemann. Continual learning with bayesian neural networks for non-stationary data. In *International Conference* on Learning Representations, 2019.
- Yann LeCun. The mnist database of handwritten digits. http://yann.lecun.com/exdb/mnist/, 1998.
- David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. In *Advances in neural information processing systems*, pp. 6467–6476, 2017.
- Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7765–7773, 2018.
- Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165. Elsevier, 1989.
- Rahul Ramesh and Pratik Chaudhari. Model zoo: A growing brain that learns continually. In *International Conference on Learning Representations*, 2021.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.
- Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv* preprint arXiv:1810.11910, 2018.
- Victoria JH Ritvo, Nicholas B Turk-Browne, and Kenneth A Norman. Nonmonotonic plasticity: how memory retrieval drives learning. *Trends in cognitive sciences*, 23(9):726–742, 2019.
- Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection memory for continual learning. arXiv preprint arXiv:2103.09762, 2021.
- Jonathan Schwarz, Jelena Luketina, Wojciech M Czarnecki, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. arXiv preprint arXiv:1805.06370, 2018.
- Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. *arXiv preprint arXiv:1801.01423*, 2018.
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. In *Advances in Neural Information Processing Systems*, pp. 2990–2999, 2017.
- Michalis K Titsias, Jonathan Schwarz, Alexander G de G Matthews, Razvan Pascanu, and Yee Whye Teh. Functional regularisation for continual learning with gaussian processes. *arXiv preprint arXiv:1901.11356*, 2019.
- Tyler M Tomita, Morgan D Barense, and Christopher J Honey. The similarity structure of real-world memories. *bioRxiv*, 2021.
- Gido M van de Ven and Andreas S Tolias. Generative replay with feedback connections as a general strategy for continual learning. *arXiv preprint arXiv:1809.10635*, 2018.
- Johannes Von Oswald, Christian Henning, João Sacramento, and Benjamin F Grewe. Continual learning with hypernetworks. *arXiv preprint arXiv:1906.00695*, 2019.
- Matthew A Wilson and Bruce L McNaughton. Reactivation of hippocampal ensemble memories during sleep. *Science*, 265(5172):676–679, 1994.

David H Wolpert and William G Macready. No free lunch theorems for optimization. *IEEE* transactions on evolutionary computation, 1(1):67–82, 1997.

Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. *Proceedings of machine learning research*, 70:3987, 2017.

# A APPENDIX

You may include other additional sections here.