# LEARNING WITH PRESERVING FOR CONTINUAL MUL TITASK LEARNING

Anonymous authors

Paper under double-blind review

# ABSTRACT

Artificial Intelligence (AI) drives advancements across fields, enabling capabilities previously unattainable. Modern intelligent systems integrate increasingly specialized tasks, such as improving tumor classification with tissue recognition or advancing driving assistance with lane detection. Typically, new tasks are addressed by training single-task models or re-training multitask models, which becomes impractical when prior data is unavailable or new data is limited. This paper introduces Continual Multitask Learning (CMTL), a novel problem category critical for future intelligent systems yet overlooked in current research. CMTL presents unique challenges beyond the scope of traditional Continual Learning (CL) and Multitask Learning (MTL). To address these challenges, we propose Learning with Preserving (LwP), a novel approach for CMTL that retains previously learned knowledge while supporting diverse tasks. LwP employs a Dynamically Weighted Distance Preservation loss function to maintain representation integrity, enabling learning across tasks without a replay buffer. We extensively evaluate LwP on three benchmark datasets across two modalities-inertial measurement units of multivariate time series data for quality of exercises assessment and image datasets. Results demonstrate that LwP outperforms existing continual learning baselines, effectively mitigates catastrophic forgetting, and highlights its robustness and generalizability in CMTL scenarios.

028 029

031 032

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

## 1 INTRODUCTION

Artificial intelligence is driving progress across numerous critical fields, enabling innovations that 033 were once beyond reach. Increasingly, more specialized and detailed tasks are being integrated into 034 existing intelligent systems, enhancing their capabilities. For example, in medical imaging, tumor 035 classification may evolve to include incremental annotation of tumor shape recognition and tissue density analysis Kaustaban et al. (2022); Freeman et al. (2021). Similarly, intelligent driving assis-037 tance systems advance from basic object detection to identifying lanes and recognizing traffic signs Shaheen et al. (2022). Traditionally, additional tasks are integrated by training new single-task models or retraining multitask models, which fall short when previous data is inaccessible or new data is 040 limited. Compiling comprehensive datasets with all labels simultaneously is often unfeasible due to data privacy, resource constraints, or the sequential nature of data collection and annotation. There-041 fore, labels arrive sequentially, requiring a suitable learning paradigm. Recent works in continual 042 and multitask learning, such as Mirzadeh et al. (2020) and Liao et al. (2022), <sup>1</sup> have addressed these 043 challenges by integrating aspects of continual and multitask learning. However, these approaches 044 often assume access to all tasks or do not generalize new tasks with previous ones. In this paper, 045 we propose *Continual Multitask Learning* (CMTL), a new problem category where input originates 046 from same dataset across tasks, but each task introduces distinct data to label spaces. This reflects 047 real-world scenarios where data drawn from a specific domain are annotated with different attributes 048 over time, requiring models to generate inferences for all the learned attributes for each input. CMTL 049 poses additional challenges compared to traditional Continual Learning (CL) and Multitask Learning (MTL). It requires models to retain knowledge from previous tasks (a CL challenge) and develop 051 shared representations beneficial to multiple tasks (an MTL goal), all while handling new tasks se-052 quentially without access to previous data. In traditional CL, especially task-incremental learning,

<sup>053</sup> 

<sup>&</sup>lt;sup>1</sup>Further discussion in Appendix A



Figure 1: Comparison among CL, MTL, and CMTL. Two key differences of CMTL compared to the other two scenarios are: (1) inputs stem from a consistent underlying distribution, with labels representing features that any input might have, much like in MTL, and (2) labels are provided sequentially, akin to CL. Models must generalize shared representations while minimizing catastrophic forgetting.

061

062

063

064

054

068 models handle a single task where new classes or labels are introduced over time-like recogniz-069 ing additional colors in image classification—within the same domain. In contrast, CMTL involves learning different tasks sequentially (e.g., color, shape, size), requiring models to adapt to new task 071 domains while preserving shared representations, as shown in Figure 1. Unlike CL, which focuses 072 on learning new tasks and mitigating catastrophic forgetting Wang et al. (2023; 2022), and MTL, 073 which learns multiple tasks simultaneously, CMTL balances both challenges in a sequential frame-074 work. This introduces complexities such as task interference and the need for models to generalize across tasks not available concurrently. 075 076

Although CMTL can be classified as a subcategory of Task-incremental Learning (Task-IL) Van 077 De Ven et al. (2022), conventional CL approaches often fail to surpass the performance of multitask 078 models or even simple single-task models under these conditions De Lange et al. (2021); Yoon et al. 079 (2019). Our experiments corroborate this, as shown in Table 1 in Section 4. We hypothesize this 080 shortfall arises because traditional CL methods treat new tasks in isolation, focusing narrowly on 081 task-specific distinctions without considering the broader feature space.

082 To address these challenges, we bridge the gap by introducing Learning with Preserving (LwP). In 083 this novel approach, we preserve previously learned knowledge in a way that remains applicable and 084 beneficial across diverse tasks that may share underlying knowledge structures. This enhancement 085 is designed to maintain both implicitly and explicitly acquired knowledge, ensuring that the learned representations are rich and generalizable enough to facilitate learning in future tasks without interference. The main **contributions** of this paper can be summarized as follows: a) We propose a 087 new scenario of continual learning, CMTL, highlighting its unique challenges and significance in 088 practical applications where labels arrive sequentially and comprehensive datasets are impractical. 089 **b**) We introduce *Learning with Preserving* (LwP), a novel framework along with a preserving loss 090 function that maintains and distills the integrity of the latent space, ensuring it is conducive to learn-091 ing across prior and future tasks without relying on a replay buffer. c) We demonstrate, through 092 extensive evaluation across two modalities — IMU sensing data (assessing the quality of exercise) 093 and image datasets ---- that our method outperforms existing baselines, including traditional CL and 094 MTL models, and showcases capabilities in CMTL scenarios.

095 096 097

098

PROBLEM FORMULATION: CONTINUAL MULTITASK LEARNING 2

099 We propose a subcategory of incremental learning scenarios that closely resemble multitask learning 100 settings. We consider a sequential learning scenario involving T tasks  $\{\mathcal{T}_t\}_{t=1}^T$ . Each task  $\mathcal{T}_t$  is 101 associated with a label space  $\mathcal{Y}_t$  and involves learning a mapping  $f_t: \mathcal{X} \to \mathcal{Y}_t$ . The input space  $\mathcal{X}$ 102 is common across all tasks, with input data  $x \in \mathcal{X}$  drawn from an identical distribution  $P_X$ . At each 103 time step t, we receive a dataset  $D_t^1 = \{(x_i, y_i^t)\}_{i=1}^{n_t}$  sampled from distribution  $\mathcal{D}_t$ , where  $x_i \sim P_X$ 104 is an input sample,  $y_i^t \in \mathcal{Y}_t$  is the corresponding label for task  $\mathcal{T}_t$ . Note that for time t, only label  $y_i^t$ 105 is available. Other labels  $y_i^j$  for  $j \neq t$  cannot be observed at time t. 106

Our goal is to find a predictor  $\varphi(\boldsymbol{x}; \theta_s, \theta_t) : \mathcal{X} \to \mathcal{Y}_1 \times \mathcal{Y}_2 \times \cdots \times \mathcal{Y}_T$  parameterized by a set of 107 shared parameters  $\theta_s$  and task-specific parameters  $\theta_t$ , such that



Figure 2: Overview of the LwP training framework on a human face dataset. While learning task  $\mathcal{T}1$  (wearing hat) with  $\mathcal{L}_{cur}$ , the model preserves prior knowledge (sunglasses) through supervised pseudolabeling  $\mathcal{L}_{old}$  and implicit knowledge retention via  $\mathcal{L}_{DWDP}$ .

$$\mathcal{L}(\theta_s, \{\theta_t\}_{t=1}^T) \coloneqq \sum_{t=1}^T \mathbb{E}_{(\boldsymbol{x}, y^t) \leftarrow \mathcal{D}_t} \left[ \ell \left( y^t, \varphi(\boldsymbol{x}, t; \theta_s, \theta_t) \right) \right], \tag{1}$$

is minimized for some loss function  $\ell(\cdot, \cdot)$ .

# 3 LEARNING WITH PRESERVING

#### 3.1 OVERVIEW

We introduce LwP, a versatile framework designed to effectively manage CMTL scenarios, as depicted in Figure 2. This framework incorporates neural network functions  $f_{\theta_s}(x)$  to create a shared representation z, along with  $g_{\theta_t}(z)$ , which represents task-specific layers for task t and utilizes z to generate predictions for the  $t^{th}$  task. This requirement is essential for the model to acquire a shared and generalizable representation space in z.

When training the current task t, we preserve and freeze the previous model to generate pseudolabels for all the previous t - 1 tasks. The current model, which is a duplicate of the previous one, includes an additional task-specific layer that will take z as input and learn to predict the current task label  $y_t$  using an appropriate supervised loss. Concurrently, the outputs for the previous tasks aim to minimize their supervised loss objectives utilizing pseudolabels from the preceding model.

Following this, we present our key novelty and apply the Dynamically Weighted Distance Preservation (DWDP) loss to preserve the knowledge that has been implicitly learned. Overall, the total objective function for the model while learning task t is defined as:

$$\mathcal{L}_{lwp} = \lambda_{c} \mathcal{L}_{cur}(y_{t}, \hat{y}_{t}) + \lambda_{o} \mathcal{L}_{old}(\tilde{y}_{o}, \hat{y}_{o}) + \lambda_{d} \mathcal{L}_{DWDP}(\boldsymbol{z}^{[t]}, \boldsymbol{z}^{[t-1]}, y_{t})$$
(2)

where

$$\begin{split} \mathcal{L}_{\text{DWDP}} &= \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} m_{ij} \left( d(\boldsymbol{z}_i^{[t-1]}, \boldsymbol{z}_j^{[t-1]}) - d(\boldsymbol{z}_i^{[t]}, \boldsymbol{z}_j^{[t]}) \right)^2, \\ m_{ij} &= \begin{cases} 1, & \text{if } y_i^{[t]} = y_j^{[t]}, \\ 0, & \text{otherwise}, \end{cases} \\ \boldsymbol{z}_i^{[t]} &= f_{\boldsymbol{a}^{[t]}}(\boldsymbol{x}_i). \end{split}$$

 $d(z_i, z_j) : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$  represents either a distance or similarity metric, with  $\lambda_c$ ,  $\lambda_o$  and  $\lambda_d$  as 159 hyperparameters.  $y_t$  and  $\hat{y}_t$  denote ground truth and model output of the current task label while 160  $\tilde{y}_o$  and  $\hat{y}_o$  denoting previous model's outputs (pseudolabels) and current model's outputs for old 161 tasks, respectively.  $\mathcal{L}_{cur}$  and  $\mathcal{L}_{old}$  represent appropriate supervised learning losses for respective 162 tasks, such as cross entropy or mean squared error. Note that the previous model  $\theta^{[t-1]}$  is frozen to



Figure 3: Development of representation space over  $\mathcal{T}$ . LwP preserves the structure as new tasks are learned.

175

176

177

178

171

produce stationary pseudolabels and z. In other words, in addition to using pseudolabel to maintain performance on old tasks, we introduce a regularization term aimed at preserving the structure of shared representations by reducing the differences in pairwise similarities (distances) between the model's previous task and the current one if the pairs have the same label for the current task.

#### 179 180 3.2 PRESERVING IMPLICIT KNOWLEDGE

In the context of CMTL, if  $\theta_t$ , the task-specific parameters, are simply linear projection layers applied to the final layer of the shared parameters  $\theta_s$ , we observe that Learning without Forgetting (LwF) Li & Hoiem (2017b) is interpreted as an approximation of the multitask learning objective that encourages the formation of more informative and generalized representation space in z.

Motivated by this observation, we show that  $\mathcal{L}_{DWDP}$  is a result of incorporating *implicitly learned knowledge* as an optimization objective. We define such knowledge as the capability of the model's representation to provide an approximate solution to some unknown problem. Therefore, in order to preserve implicitly learned knowledge, we intend to find a loss function that can preserve approximate solutions for *any* problems that can be defined in z.

In order to preserve all approximate solutions from the representation space alone, we exploit that kernel methods with the Gaussian kernel are universal approximators Hammer & Gersmann (2003).

Given two sets of representations  $Z, Z' \in \mathbb{R}^{n \times d}$ , where each row corresponds to z, our objective is to ensure that Z' maps to the same Reproducing Kernel Hilbert Space (RKHS) as Z under the Gaussian kernel. To achieve this, we derive a loss function  $\mathcal{L}_{pres}$  that encourages the alignment of the pairwise similarities encoded by the Gaussian kernel in both representation spaces.

The Gaussian kernel is defined as  $k(\mathbf{z}_i, \mathbf{z}_j) = \exp\left(-\frac{\|\mathbf{z}_i - \mathbf{z}_j\|^2}{2\sigma^2}\right)$ , where  $\mathbf{z}_i, \mathbf{z}_j \in \mathbb{R}^d$  are representations, and  $\sigma > 0$  is the bandwidth parameter controlling the kernel's sensitivity to distance. The Gaussian kernel is a positive definite function, inducing an RKHS  $\mathcal{H}$  with an implicit feature mapping  $\phi : \mathbb{R}^d \to \mathcal{H}$  such that  $k(\mathbf{z}_i, \mathbf{z}_j) = \langle \phi(\mathbf{z}_i), \phi(\mathbf{z}_j) \rangle_{\mathcal{H}}$ .

For a set of representations Z, the Gram matrix  $K(Z) \in \mathbb{R}^{n \times n}$  is constructed with entries  $K_{ij}(Z) = k(z_i, z_j)$ . Similarly, we construct K(Z') for Z'. Our goal is to align K(Z) and K(Z') such that the pairwise similarities in Z' match those in Z. This alignment ensures that Z and Z' are mapped to the same locations in the RKHS up to an isometry.

To formalize the alignment objective, we define the loss function  $\mathcal{L}_{pres}$  as the squared Frobenius norm of the difference between the two kernel matrices:

$$\mathcal{L}_{pres}(Z, Z') = \|K(Z) - K(Z')\|_{F}^{2} = \sum_{i=1}^{n} \sum_{j=1}^{n} \left(k(\boldsymbol{z}_{i}, \boldsymbol{z}_{j}) - k(\boldsymbol{z}_{i}', \boldsymbol{z}_{j}')\right)^{2}.$$
 (3)

210 211

208 209

212 Minimizing  $\mathcal{L}_{pres}$  with respect to Z' (while keeping Z fixed) encourages the kernel matrices to 213 become identical, i.e.,  $K(Z') \approx K(Z)$ . This implies that for all pairs  $(i, j), k(\mathbf{z}'_i, \mathbf{z}'_j) \approx k(\mathbf{z}_i, \mathbf{z}_j)$ .

By minimizing  $\mathcal{L}_{pres}$ , we effectively align the images of Z and Z' under the feature map  $\phi$ :

<

$$\phi(\mathbf{z}_i), \phi(\mathbf{z}_j)\rangle_{\mathcal{H}} \approx \langle \phi(\mathbf{z}'_i), \phi(\mathbf{z}'_j) \rangle_{\mathcal{H}}, \quad \forall i, j.$$
 (4)



Figure 4: The impact of  $\mathcal{L}_{DWDP}$  on a two-dimensional toy dataset, where  $y_1$  (O vs. X) indicates an XOR problem and  $y_2$  (blue vs. red) signifies a concentric circle problem. The figure shows the representation space after training on  $y_2$  without  $\mathcal{L}_{DWDP}$  (left) and with  $\mathcal{L}_{DWDP}$  (right). The latter successfully preserves the cluster structures of the former representation, which is advantageous for learning  $y_1$  in subsequent phases.

This alignment implies that there exists an isometry  $T : \mathcal{H} \to \mathcal{H}$  such that:

$$\phi(\mathbf{z}_i') = T(\phi(\mathbf{z}_i)), \quad \forall i.$$
(5)

For any function  $f \in \mathcal{H}$ , there exists a weight vector  $w \in \mathcal{H}$  such that  $f(z) = \langle w, \phi(z) \rangle_{\mathcal{H}}$ . The evaluation of f at  $z'_i$  becomes:

$$f(\boldsymbol{z}_{i}') = \langle w, \phi(\boldsymbol{z}_{i}') \rangle_{\mathcal{H}} = \langle w, T(\phi(\boldsymbol{z}_{i})) \rangle_{\mathcal{H}}.$$
(6)

Because T is an isometry, its adjoint  $T^*$  is also an isometry, and we can write:

$$f(\mathbf{z}_i) = \langle T^* w, \phi(\mathbf{z}_i) \rangle_{\mathcal{H}}.$$
(7)

Define  $w' = T^*w$  and  $f'(z) = \langle w', \phi(z) \rangle_{\mathcal{H}}$ . Then:

 $f(\boldsymbol{z}_i') = f'(\boldsymbol{z}_i), \quad \forall i.$ (8)

Thus, Z' becomes an alternative representation that is functionally equivalent to Z in terms of any operations performed within the RKHS induced by the Gaussian kernel. Now, consider a learning problem defined on Z:

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \ell(f(\boldsymbol{z}_i), y_i) + \Omega(f),$$
(9)

and the corresponding problem on Z':

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \ell(f(\boldsymbol{z}_{i}^{\prime}), y_{i}) + \Omega(f).$$

$$(10)$$

Using the relationship  $f(z'_i) = f'(z_i)$ , the loss terms satisfy  $\ell(f(z'_i), y_i) = \ell(f'(z_i), y_i)$ . Since  $||f||_{\mathcal{H}} = ||f'||_{\mathcal{H}}$ , the regularization terms are equal:  $\Omega(f) = \Omega(f')$ . Thus, the risk functionals for the problems on Z and Z' are equivalent when considering f and f':

$$\frac{1}{n}\sum_{i=1}^{n}\ell(f(\boldsymbol{z}_{i}'), y_{i}) + \Omega(f) = \frac{1}{n}\sum_{i=1}^{n}\ell(f'(\boldsymbol{z}_{i}), y_{i}) + \Omega(f').$$
(11)

Because the risk functionals are equivalent, the optimal solutions  $f^*$  obtained on Z' correspond to the optimal solutions  $f'^*$  on Z via the isometry  $T^*$ :

$$f^*(z_i') = f'^*(z_i).$$
 (12)

284

286 287 288

289

290 291

292

270

This means any model trained on Z can be transformed to a model on Z' with identical performance, and vice versa.

Through empirical observation, we have determined that maintaining the squared Euclidean distance instead of k directly leads to enhanced performance. This is likely because the non-exponentiated distance metric more effectively retains the global structure of the representation space. Refer to Appendix 4.6 for the experimental data. Additionally, in Appendix C, we show that the difference in RBF kernel values is bounded by the difference in the squared  $L^2$  norm.

Hereby we define a family of such losses that preserve some distance (or similarity) metric between pairs of representations as the following:

$$\mathcal{L}_{\text{pres}}(\boldsymbol{z}, \boldsymbol{z}') = \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} \left( d(\boldsymbol{z}_i, \boldsymbol{z}_j) - d(\boldsymbol{z}'_i, \boldsymbol{z}'_j) \right)^2,$$
(13)

(14)

where d represents either a distance or a similarity function. Note that it no longer needs to be a kernel to include a broader variety of metrics.

## 3.3 DYNAMIC WEIGHTING

 $\mathcal{L}_{\text{pres}}$  is designed to maintain the implicitly learned knowledge of the input data in the representation space. However, in scenarios where there are distinct classes or labels, this loss can conflict with other objectives, such as separating distinct classes.

To address this issue, we introduce the Dynamically Weighted Distance Preservation (DWDP) Loss,  $\mathcal{L}_{DWDP}$ . This loss function adapts the preservation loss by applying a dynamic mask  $m_{ij}$ , which controls the contribution of each pairwise comparison based on their label similarity. The intuition behind this modification is to deactivate the preservation requirement for pairs with different labels, thus preventing conflicts with the separation objectives.

 $m_{ij} = \begin{cases} 1, & \text{if } y_i^{[t]} = y_j^{[t]}, \\ 0, & \text{otherwise,} \end{cases}$ 

301 The dynamic mask  $m_{ij}$  is defined as follows:

302 303

304 305

306 307

where  $y^{[t]}$  represents the labels of the current task.

Thus, the DWDP Loss is then given by:

310 311

312 313

308

$$\mathcal{L}_{\text{DWDP}}(\boldsymbol{z}^{[t-1]}, \boldsymbol{z}^{[t]}, y^{[t]}) = \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} m_{ij} \left( d(\boldsymbol{z}_i^{[t-1]}, \boldsymbol{z}_j^{[t-1]}) - d(\boldsymbol{z}_i^{[t]}, \boldsymbol{z}_j^{[t]}) \right)^2$$
(15)

Consequently, this modification alleviates the objective conflict issue at the cost of reducing the scope for preservation to intraclass sets of the current task. Our detailed pseudo-code algorithm is presented in Appendix B.

317 318

4 EVALUATION4.1 OVERVIEW

319 320

321

We present a set of experiments designed to rigorously validate our approach using three benchmark datasets that span multiple modalities. 1) We conduct a comprehensive performance evaluation of our method, LwP, comparing it to state-of-the-art CL techniques. This assessment focuses on the average accuracy across all tasks after training is completed. 2) We analyze the extent of catastrophic
forgetting in each model, employing the Backward Transfer (BWT) metric Lopez-Paz & Ranzato
(2017a) to quantify the accuracy degradation over successive tasks. 3) We benchmark the rate of
performance improvement per train iteration not only against other baseline CL models but also
against MTL oracles, which serve as an upper bound for the problem. 4) We explore the effects of
dynamic weighting and various distance/similarity functions, d, as proposed in other studies, on the
performance metrics of LwP.

331 Additionally, in the Appendix D, we explain more detailed setup including attributes and conduct 332 further evaluations through additional experiments. A detailed comparison with MTL methods is 333 presented in Appendix D.2. Furthermore, we supply supplemental diagrams that illustrate the pro-334 gression of the accuracy over time for both the PhysiQ and FairFace datasets, which can be found in Appendix D.3. Appendix D.4 elaborates on the influence of the number of training examples on 335 the overall performance of each model. In Appendix D.5, we benchmark the rate of performance 336 improvement per train iteration not only against other baseline CL models but also against MTL 337 methods. In Appendix D.6, we analyze the impact of model size and image size on the performance 338 of all the methods. Additionally, in Appendix D.7, we explore training the first 5 tasks of the CelebA 339 dataset using an MTL scheme, followed by a CL setting for the remaining 5 tasks. 340

341

343

342 4.2 EXPERIMENT SETUP

**Datasets** We utilize three datasets from two distinct modalities, each structured for taskincremental learning. In this setting, each task is only exposed to a subset of training samples:

The CelebA dataset Liu et al. (2018), consisting of 200,000 images with 40 facial attributes. For our
 work, we focus on 10 of the most balanced attributes. The train dataset is equally subdivided for
 each task, leading to 20,000 images per task. For simplicity, input images are resized to 32x32.

The PhysiQ dataset Wang & Ma (2023), which contains approximately 4,500 samples collected using inertial measurement units (IMUs) to capture the quality of physical exercises. The data is collected on accelerometer and gyroscope modality of 50 Hz sampling rate for 31 participants with three attributes. Each task corresponds to one of these attributes with around 1,500 samples.

The Fairface dataset Karkkainen & Joo (2021), which includes 100,000 images with three attributes.
 Following the same subdivison procedure, the dataset results in containing approximately 33,333 images with a resolution of 128x128 per task. Not only the tasks differ from those of CelebA, but also the images are not resized in order to show our approach is scalable.

358

Baselines Our primary emphasis is on CL baselines since integrating many MTL methods into 359 CMTL often requires substantial modifications to accommodate the incremental characteristics of 360 CMTL. For CL, we compare against Online Bias Correction (OBC) Chrysakis & Moens (2023), 361 Dual View Consistency (DVC) Gu et al. (2022), Dark Experience Replay (DER) Buzzega et al. 362 (2020), DER++ Boschini et al. (2022), Function Distance Regularization (FDR) Benjamin et al. 363 (2019), Experience Replay (ER) Robins (1995); Ratcliff (1990), Gradient-based Sample Selection 364 (GSS) Aljundi et al. (2019b), online Elastic Weight Consolidation (oEWC) Kirkpatrick et al. (2017), Synaptic Intelligence (SI) Zenke et al. (2017), and Learning without Forgetting (LwF) Li & Hoiem 366 (2017b). In addition, we compare our approach with MTL methods, which are detailed in Appendix 367 D.2. These include the basic MTL approach of training all tasks simultaneously with different 368 predictors Caruana (1997), as well as more advanced techniques like PCGrad Yu et al. (2020), Impartial MTL (IMTL) Liu et al. (2021), and NashMTL Navon et al. (2022). We also include a 369 single task learning (STL) baseline, where each task is learned separately. For the choice of distance 370 metric d, we test common options such as Euclidean distance and cosine similarity, as well as loss 371 functions designed to preserve relational knowledge, such as those proposed in RKD Park et al. 372 (2019) and Co2L Cha et al. (2021). 373

374

375 **Model Architectures** We use an untrained ResNet-18 for CelebA and FairFace datasets. Each 376 task is predicted after a linear projection layer applied to the flattened last shared layer  $z \in \mathbb{R}^{512}$ . 377 Similarly, for PhysiQ, we use a 3-layer 1DCNN model with  $z \in \mathbb{R}^{128}$  final shared layer connected to task-specific linear layers. We evaluate additional architectures and image sizes in Appendix D.6.



Figure 5: Matrices showcasing the accuracy progression for various models for Dataset CelebA. Each column corresponds to an iteration of the task, arranged sequentially from left to right. We generate the confusion matrices normalized on all the tasks for all the models for consistency.

Method Type	Model	CelebA (10 Tasks)	PhysiQ (3 Tasks)	FairFace (3 Tasks)
STL	-	$72.230 \pm 7.297$	$87.167 \pm 10.102$	$64.435 \pm 3.660$
	LwF	$64.626 \pm 10.806$	$69.952 \pm 21.090$	$61.034 \pm 6.162$
	oEWC	$69.666 \pm 9.019$	$82.640 \pm 12.166$	$63.604 \pm 3.122$
	ER	$67.598 \pm 7.452$	$76.798 \pm 16.347$	$63.220 \pm 4.730$
	SI	$68.735 \pm 10.545$	$83.727 \pm 11.828$	$63.359 \pm 3.451$
CL	GSS	$71.680 \pm 8.468$	$85.741 \pm 10.950$	$64.230 \pm 3.918$
	FDR	$69.514 \pm 8.917$	$71.859 \pm 18.687$	$63.709 \pm 3.151$
	DER	$70.703 \pm 8.388$	$84.796 \pm 11.168$	$64.114 \pm 3.484$
	DERPP	$67.693 \pm 9.425$	$82.838 \pm 13.775$	$63.806 \pm 3.694$
	DVC	$71.441 \pm 7.640$	$85.100 \pm 10.381$	$63.848 \pm 3.193$
	OBC	$70.829 \pm 8.267$	$83.999 \pm 11.377$	$63.872 \pm 3.449$
CMTL	LwP	$\textbf{73.484} \pm \textbf{8.019}$	$\textbf{88.242} \pm \textbf{12.010}$	$\textbf{66.482} \pm \textbf{3.138}$

Table 1: Accuracy Percentage Comparison Across Models and Datasets

In this experiment, we evaluate the performance of Our LwP against several state-of-the-art CL methods. All methods, except for LwF and LwP, are provided with a buffer size of 512 for the CelebA and FairFace datasets, and 46 for the PhysiQ dataset, corresponding to approximately 2-3% of the training set for each dataset. Each model is trained five times using different random seeds. The standard training protocol consists of 20 epochs, with a batch size of 256 for image-based datasets and 32 for PhysiQ, coupled with early stopping. For PhysiQ, we only compare the average accuracy across the final task iteration due to the training instability caused by smaller dataset size. Table 1 reports the average test accuracy, along with the standard deviation over five runs for each method and dataset. Fig. 5 visualizes the progression of task accuracy in task iterations (left to right). Additional results are provided in Appendix D.3 

Table 1 highlights that LwP consistently achieves superior performance across all three benchmarks
and is the only method to exceed the Single Task Learning (STL) baseline. This suggests that other
continual learning methods likely experience significant task interference. Additional results with
MTL are provided in Appendix D. Furthermore, our approach is modality-agnostic, as evidenced
by LwP's ability to generalize across different domains. This is demonstrated by the results on the
PhysiQ dataset from the IMU sensor domain, which underscores LwP's robustness against challenges unique to non-image-based tasks.

431 The results suggest that LwP demonstrates competitive performance compared to existing continual learning methods across a range of benchmarks in CMTL settings. LwP consistently achieves higher

accuracy than other approaches, including surpassing the STL baseline, indicating its potential to reduce task interference and catastrophic forgetting in different modalities.

# 4.4 BACKWARD TRANSFER



Figure 6: Backward Transfer Diagrams for Various Datasets

The Backward Transfer Lopez-Paz & Ranzato (2017a) is a metric to evaluate the influence of learning the current task on the performance of previous tasks. A positive backward transfer value indicates that, on average, accuracies on the previous tasks have increased during the current task iteration and vice versa. It is defined as:

$$BWT = \frac{1}{T-1} \sum_{i=1}^{T-1} R_{T,i} - R_{i,i},$$
(16)

where T is the index of the current task, i is an index of previous tasks ranging from 1 to T - 1,  $R_{T,i}$  is the accuracy on task i after training up to task T, and  $R_{i,i}$  is the accuracy on task i after learning.

As illustrated in Fig. 6, we observe that LwP outperforms all baselines in terms of BWT across all benchmarks. This result is consistent with the visualization shown in Fig. 3, where LwP can maintain the accuracy of each task since its initial training.

### 4.5 REPRESENTATION SPACE VISUALIZATION VIA T-SNE





Figure 7: Representation space progression over task iteration. Colors indicate different label values.

Fig. 7 shows the 2D visualization of the representation z for each model trained on the PhysiQ dataset. It was constructed using the dimensionality reduction algorithm t-SNE van der Maaten & Hinton (2008) on the PhysiQ test dataset at the end of each task iteration. x and y axes represent the two new dimensions created by the algorithm to project the high-dimensional data onto a 2D plane. Note these dimensions do not have an intrinsic meaning and rather constructed to reflect the relative distances between data points in the high-dimensional space. As demonstrated, the z produced by LwP maintains coherent cluster formations as it progressively learns new tasks without introducing considerable distortions when compared to the baseline model. This behavior is comparable to the example provided with the toy dataset depicted in Fig. 4.

### 4.6 THE ABLATION STUDY OF EFFECTIVENESS OF THE LOSS FUNCTION

To evaluate the impact of the proposed loss function, we perform experiments by selectively disabling the dynamic weighting feature and comparing it with other loss functions that also aim to

	Table 2: Ablation	comparison	on $\mathcal{L}_{DWDP}$	implementation
--	-------------------	------------	-------------------------	----------------

Method on PhysiQ	$\mathbf{LwP}\left(L^{2}\right)$	LwP (Cosine)	LwP (RBF)	IRD (Co2L)	RKD
Dynamic Weighting	$\textbf{88.2} \pm \textbf{12.0}$	$85.4 \pm 13.1$	$84.5\pm13.7$	$86.4 \pm 11.5$	$85.1\pm13.3$
W/o Dynamic Weighting	$86.0 \pm 12.3$	$84.1\pm14.4$	$84.8\pm14.5$	$79.9 \pm 17.1$	$85.9 \pm 11.9$

preserve structures. In our evaluation, we include CO2L Cha et al. (2021), RKD Park et al. (2019), and two novel variation of baselines: cosine similarity and the RBF kernel as described in eq. 3.

The findings in Table 2 indicate that the loss function with both dynamic weighting and Euclidean distance consistently surpasses the other options. We believe that the effectiveness of Euclidean distance with dynamic weighting is due to its loss not being normalized across batches, unlike previously proposed approaches.

498 499 500

501

486

493

494 495

496

497

# 5 RELATED WORK

502 MTL enhances generalization and computational efficiency by leveraging shared representations across related tasks Caruana (1997); Sener & Koltun (2018). However, optimizing multiple objec-504 tives often presents conflicting gradients. Approaches like the Multiple Gradient Descent Algorithm 505 (MGDA) Sener & Koltun (2018) seek Pareto optimal solutions through convex combinations of 506 task-specific gradients, while Gradient Surgery (PCGrad) Yu et al. (2020) projects conflicting gra-507 dients onto the normal plane of each other to reduce interference. Navon et al. (2022) 508 modeled gradient combination as a cooperative bargaining game to ensure fairness among tasks. 509 Loss balancing is also crucial, with methods like IMTL Liu et al. (2021) incorporating both gradient and loss balancing mechanisms. CL enables sequential task learning without catastrophic forgetting 510 Ratcliff (1990); Robins (1995). Techniques like MER Riemer et al. (2018) focus on maximizing 511 knowledge transfer while minimizing interference, while HAL Chaudhry et al. (2021) anchors past 512 knowledge to prevent representation drift. Bridging MTL and CL, continual multitask learning aims 513 to manage performance across sequential and concurrent tasks Wu et al. (2023), using methods like 514 MC-SGD Mirzadeh et al. (2020) to enhance CL by leveraging linear mode connectivity. Task-free 515 CL Aljundi et al. (2019a) eliminates task boundaries. More detailed discussions are available in the 516 Appendix A.

517 518 519

520

# 6 CONCLUSION

We explored the limitations of existing continual learning methods in CMTL. Our findings show 521 that conventional approaches often underperform compared to single-task models, largely due to 522 their focus on preserving explicit information while neglecting broadly useful, implicit features. To 523 address this, we introduced *Learning with Preserving* (LwP) with a dynamically weighted distance 524 preservation function. This approach maintains the structure of the representation space, preserv-525 ing implicit knowledge without needing replay buffers, making it especially valuable in privacy-526 sensitive domains like healthcare. Our experiments across various datasets demonstrated that LwP 527 surpasses state-of-the-art baselines and outperforms single-task models, consistently retaining ac-528 curacy and mitigating catastrophic forgetting. The results emphasize the importance of preserving 529 implicit knowledge and the effectiveness of our loss function. Future work could explore LwP's 530 application with non-stationary dataset or unlabeled data (i.e., investigation on KL divergence vs. 531 LwP performance), and its integration with pre-trained foundation models.

532 533 534

537

# References

Tameem Adel, Han Zhao, and Richard E Turner. Continual learning with adaptive weights (claw).
 *arXiv preprint arXiv:1911.09514*, 2019.

Rahaf Aljundi, Klaas Kelchtermans, and Tinne Tuytelaars. Task-free continual learning. In *Proceed- ings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11254–11263, 2019a.

- 540 Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection 541 for online continual learning. Advances in neural information processing systems, 32, 2019b. 542
- 543 Ari S. Benjamin, David Rolnick, and Konrad Kording. Measuring and regularizing networks in function space, June 2019. URL http://arxiv.org/abs/1805.08289. arXiv:1805.08289 544 [cs, stat].
- 546 Matteo Boschini, Lorenzo Bonicelli, Pietro Buzzega, Angelo Porrello, and Simone Calderara. Class-547 incremental continual learning into the extended der-verse. IEEE Transactions on Pattern Analy-548 sis and Machine Intelligence, 2022. 549
- 550 Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and SIMONE CALDER-551 ARA. Dark Experience for General Continual Learning: a Strong, Simple Baseline. In Advances in Neural Information Processing Systems, volume 33, pp. 15920–15930. Curran 552 Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper\_files/ 553 paper/2020/hash/b704ea2c39778f07c617f6b7ce480e9e-Abstract.html. 554
- 555 Kaidi Cao, Jiaxuan You, and Jure Leskovec. Relational multi-task learning: Modeling relations 556 between data and tasks. In International Conference on Learning Representations, 2022. URL https://openreview.net/forum?id=8Py-W8lSUgy. 558
- 559 Rich Caruana. Multitask learning. Machine learning, 28:41-75, 1997.

563

564

565

566

569

570

571

572

575

576

577

578

583

- Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co\$^2\$L: Contrastive Continual Learning, June 2021. 561 URL http://arxiv.org/abs/2106.14413. arXiv:2106.14413 [cs]. 562
  - 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 567 lifelong learning with a-gem. arXiv preprint arXiv:1812.00420, 2018b. 568
  - Arslan Chaudhry, Albert Gordo, Puneet Dokania, Philip Torr, and David Lopez-Paz. Using hindsight to anchor past knowledge in continual learning. In Proceedings of the AAAI conference on artificial intelligence, volume 35, pp. 6993–7001, 2021.
- 573 Aristotelis Chrysakis and Marie-Francine Moens. Online bias correction for task-free continual 574 learning. ICLR 2023 at OpenReview, 2023.
  - Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385, 2021.
- 579 Lea Duncker, Laura Driscoll, Krishna V Shenoy, Maneesh Sahani, and David Sussillo. Organizing 580 recurrent network dynamics by task-computation to enable continual learning. Advances in neural 581 information processing systems, 33:14387–14397, 2020. 582
- Beverly Freeman, Naama Hammel, Sonia Phene, Abigail Huang, Rebecca Ackermann, Olga Kanzheleva, Miles Hutson, Caitlin Taggart, Quang Duong, and Rory Sayres. Iterative quality 584 control strategies for expert medical image labeling. Proceedings of the AAAI Conference on Hu-585 man Computation and Crowdsourcing, 9(1):60–71, Oct. 2021. doi: 10.1609/hcomp.v9i1.18940. 586 URL https://ojs.aaai.org/index.php/HCOMP/article/view/18940.
- 588 Joshua P Gardner, Ian Simon, Ethan Manilow, Curtis Hawthorne, and Jesse Engel. MT3: Multi-task multitrack music transcription. In International Conference on Learning Representations, 2022. 590 URL https://openreview.net/forum?id=iMSjopcOn0p.
- Yanan Gu, Xu Yang, Kun Wei, and Cheng Deng. Not just selection, but exploration: Online class-592 incremental continual learning via dual view consistency. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 7442–7451, June 2022.

594	Barbara Hammer and Kai Gersmann. A Note on the Universal Approximation Capability of
595	Support Vector Machines. Neural Processing Letters, 17(1):43-53, 2003. ISSN 13704621.
596	doi: 10.1023/A:1022936519097. URL http://link.springer.com/10.1023/A:
597	1022936519097.
598	

- Geoffrey Hinton. Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531, 2015. 600
- 601 Adrián Javaloy and Isabel Valera. Rotograd: Gradient homogenization in multitask learning. arXiv preprint arXiv:2103.02631, 2021. 602
- 603 Kimmo Karkkainen and Jungseock Joo. Fairface: Face attribute dataset for balanced race, gender, 604 and age for bias measurement and mitigation. In Proceedings of the IEEE/CVF Winter Conference 605 on Applications of Computer Vision, pp. 1548–1558, 2021. 606
- Veena Kaustaban, Qinle Ba, Ipshita Bhattacharya, Nahil Sobh, Satarupa Mukherjee, Jim Martin, 607 Mohammad Saleh Miri, Christoph Guetter, and Amal Chaturvedi. Characterizing Continual 608 Learning Scenarios for Tumor Classification in Histopathology Images. In Yuankai Huo, Bryan A. 609 Millis, Yuyin Zhou, Xiangxue Wang, Adam P. Harrison, and Ziyue Xu (eds.), Medical Optical 610 Imaging and Virtual Microscopy Image Analysis, pp. 177-187, Cham, 2022. Springer Nature 611 Switzerland. ISBN 978-3-031-16961-8. 612
- 613 Donggyun Kim, Seongwoong Cho, Wonkwang Lee, and Seunghoon Hong. Multi-task processes. In International Conference on Learning Representations, 2022. URL https://openreview. 614 net/forum?id=9otKVlgrpZG. 615
- 616 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL 617 https://arxiv.org/abs/1412.6980. 618
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, An-619 drei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis 620 Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic 621 forgetting in neural networks. Proceedings of the National Academy of Sciences, 114(13): 622 3521-3526, March 2017. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1611835114. URL 623 https://pnas.org/doi/full/10.1073/pnas.1611835114. 624
- 625 Yuanpeng Li, Liang Zhao, Kenneth Church, and Mohamed Elhoseiny. Compositional language continual learning. In International Conference on Learning Representations, 2019. 626
- 627 Zhizhong Li and Derek Hoiem. Learning without forgetting. IEEE transactions on pattern analysis 628 and machine intelligence, 40(12):2935–2947, 2017a. 629
- Zhizhong Li and Derek Hoiem. Learning without Forgetting, February 2017b. URL http:// 630 arxiv.org/abs/1606.09282. arXiv:1606.09282 [cs, stat]. 631
- 632 Weibin Liao, Haoyi Xiong, Qingzhong Wang, Yan Mo, Xuhong Li, Yi Liu, Zeyu Chen, Siyu Huang, 633 and Dejing Dou. Muscle: Multi-task self-supervised continual learning to pre-train deep models 634 for x-ray images of multiple body parts. In International Conference on Medical Image Comput-635 ing and Computer-Assisted Intervention, pp. 151–161. Springer, 2022.

- Liyang Liu, Yi Li, Zhanghui Kuang, Jing-Hao Xue, Yimin Chen, Wenming Yang, Qingmin Liao, and Wayne Zhang. Towards impartial multi-task learning. In International Conference on Learning 638 *Representations*, 2021. URL https://openreview.net/forum?id=IMPnRXEWpvr. 639
- 640 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Large-scale celebfaces attributes (celeba) 641 dataset. Retrieved August, 15(2018):11, 2018.
- 642 David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. 643 In Proceedings of the 31st International Conference on Neural Information Processing Sys-644 tems, NIPS'17, pp. 6470-6479, Red Hook, NY, USA, 2017a. Curran Associates Inc. ISBN 645 9781510860964. 646
- David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. 647 Advances in neural information processing systems, 30, 2017b.

648	Arun Mallya and Svetlana Lazebnik, Packnet, Adding multiple tasks to a single network by iterative
649	pruning In Proceedings of the IEEE conference on Computer Vision and Pattern Recognition
650	pp. 7765–7773, 2018.
651	rr

- Seyed Iman Mirzadeh, Mehrdad Farajtabar, Dilan Gorur, Razvan Pascanu, and Hassan Ghasemzadeh. Linear mode connectivity in multitask and continual learning. *arXiv preprint arXiv:2010.04495*, 2020.
- Aviv Navon, Aviv Shamsian, Idan Achituve, Haggai Maron, Kenji Kawaguchi, Gal Chechik, and
   Ethan Fetaya. Multi-task learning as a bargaining game. *arXiv preprint arXiv:2202.01017*, 2022.
- Jeremy Nixon, Michael W. Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran. Measuring calibration in deep learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.
  - Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational Knowledge Distillation, May 2019. URL http://arxiv.org/abs/1904.05068. arXiv:1904.05068 [cs].
- Roger Ratcliff. Connectionist models of recognition memory: constraints imposed by learning and
   forgetting functions. *Psychological review*, 97(2):285, 1990.
- 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 interfer ence. arXiv preprint arXiv:1810.11910, 2018.
- Anthony Robins. Catastrophic forgetting, rehearsal and pseudorehearsal. *Connection Science*, 7(2): 123–146, 1995.
- Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. Advances in neural information processing systems, 31, 2018.
- Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *International conference on machine learning*, pp. 4548–4557. PMLR, 2018.
- Khadija Shaheen, Muhammad Abdullah Hanif, Osman Hasan, and Muhammad Shafique. Continual Learning for Real-World Autonomous Systems: Algorithms, Challenges and Frameworks. *Journal of Intelligent & Robotic Systems*, 105(1):9, May 2022. ISSN 0921-0296, 1573-0409. doi: 10.1007/s10846-022-01603-6. URL https://link.springer.com/10.1007/s10846-022-01603-6.
- Gido M. Van De Ven, Tinne Tuytelaars, and Andreas S. Tolias. Three types of incremental learning. *Nature Machine Intelligence*, 4(12):1185–1197, December 2022. ISSN 2522-5839. doi: 10.1038/s42256-022-00568-3. URL https://www.nature.com/articles/s42256-022-00568-3.
- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of Machine Learning Research, 9(86):2579-2605, 2008. URL http://jmlr.org/papers/v9/ vandermaaten08a.html.
- Johannes Von Oswald, Christian Henning, Benjamin F Grewe, and João Sacramento. Continual learning with hypernetworks. *arXiv preprint arXiv:1906.00695*, 2019.
- Hanchen David Wang and Meiyi Ma. PhysiQ: Off-site Quality Assessment of Exercise in Physical Therapy. Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies, 6(4):208:1-208:25, January 2023. doi: 10.1145/3570349. URL https://doi.org/10. 1145/3570349.

702 703 704	Zhenyi Wang, Li Shen, Le Fang, Qiuling Suo, Donglin Zhan, Tiehang Duan, and Mingchen Gao. Meta-learning with less forgetting on large-scale non-stationary task distributions. In <i>European Conference on Computer Vision</i> , pp. 221–238. Springer, 2022.
705 706 707	Zhenyi Wang, Enneng Yang, Li Shen, and Heng Huang. A comprehensive survey of forgetting in deep learning beyond continual learning. <i>arXiv preprint arXiv:2307.09218</i> , 2023.
707 708 709	Zhenyi Wang, Yan Li, Li Shen, and Heng Huang. A unified and general framework for continual learning arXiv preprint arXiv:2403 13249, 2024
710 711 712	Hongxin Wei, Renchunzi Xie, Hao Cheng, Lei Feng, Bo An, and Yixuan Li. Mitigating neural network overconfidence with logit normalization, 2022. URL https://arxiv.org/abs/
712 713 714	2205.09310. Zibao Wu, Huy Tran, Hamed Pirsiayash, and Sobeil Kolouri. Is multi-task learning an upper bound
715 716	for continual learning? In ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1–5. IEEE, 2023.
717 718 719 720 721	Enneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guibing Guo, Xingwei Wang, and Dacheng Tao. Adamerging: Adaptive model merging for multi-task learning. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=nZP6NgD3QY.
722 723	Jaehong Yoon, Saehoon Kim, Eunho Yang, and Sung Ju Hwang. Scalable and order-robust continual learning with additive parameter decomposition. <i>arXiv preprint arXiv:1902.09432</i> , 2019.
724 725 726	Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. <i>Advances in Neural Information Processing Systems</i> , 33:5824–5836, 2020.
727 728 729	Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In <i>International conference on machine learning</i> , pp. 3987–3995. PMLR, 2017.
731 732	
733 734	
735 736 737	
738 739	
740 741	
742 743	
744 745 746	
747 748	
749 750	
751 752	
753 754 755	
100	

# 756 A RELATED WORKS

# 758 A.1 MULTITASK LEARNING

Multitask learning (MTL) has been extensively explored for its ability to leverage shared representations across multiple related tasks, thereby enhancing generalization and computational efficiency Caruana (1997); Sener & Koltun (2018); Javaloy & Valera (2021); Gardner et al. (2022); Kim et al. (2022); Cao et al. (2022); Yang et al. (2024); Liu et al. (2021). In MTL, models are trained on multiple tasks simultaneously, with the assumption that learning tasks together allows the model to capture commonalities and differences among tasks, leading to better performance than training each task separately.

767 A central challenge in MTL is the optimization of multiple objectives, which often present conflict-768 ing gradients that can impede the convergence and performance of the model. To address this, var-769 ious gradient balancing approaches have been proposed. Sener and Koltun Sener & Koltun (2018) introduced the Multiple Gradient Descent Algorithm (MGDA), which seeks Pareto optimal solu-770 tions by finding a convex combination of task-specific gradients. Building on this, Yu et al. Yu et al. 771 (2020) proposed Gradient Surgery (PCGrad), which directly modifies conflicting gradients by pro-772 jecting them onto the normal plane of each other to reduce negative interference. More recently, 773 Navon et al. Navon et al. (2022) approached MTL from a game-theoretic perspective, modeling the 774 gradient combination step as a cooperative bargaining game and employing the Nash Bargaining 775 Solution to ensure proportional fairness among tasks. 776

In addition to gradient balancing, loss balancing is crucial for stable and unbiased learning in MTL.
Liu et al. Liu et al. (2021) introduced IMTL, which incorporates both gradient and loss balancing
mechanisms. Their method, IMTL-G, ensures unbiased updates to task-shared parameters by finding
the geometric angle bisector of task gradients, while IMTL-L automatically learns loss weighting
parameters to harmonize the scales of different task losses.

While these MTL methods have advanced the ability to learn multiple tasks simultaneously, they typically assume that all task data is available at training time and can be processed jointly. This assumption does not hold in scenarios where tasks and their associated data arrive sequentially, as in our defined problem, Continual Multitask Learning (CMTL). In such cases, models must learn new tasks without access to all previous data, and ideally, they should leverage new tasks to improve performance on prior tasks.

Our work differs from traditional MTL approaches by addressing the sequential arrival of tasks and data, where tasks are learned iteratively rather than simultaneously. Unlike MTL methods that focus on balancing gradients and losses across tasks trained together, our approach must handle the challenge of incorporating new tasks without retraining on previous tasks' data. Furthermore, we introduce mechanisms to utilize new task data to enhance the model's generalizability on earlier tasks, which is not considered in standard MTL frameworks.

793

# A.2 CONTINUAL LEARNING (CL)

CL aims to enable models to learn sequentially from a stream of tasks without forgetting previously acquired knowledge, addressing the challenge of catastrophic forgetting Ratcliff (1990); Robins (1995). Various methods have been developed to tackle this problem, broadly categorized into *rehearsal-based methods, knowledge distillation,* and *regularization-based techniques.*

800

801 **Rehearsal-based methods** Early works such as Ratcliff (1990); Robins (1995) introduced Expe-802 rience Replay (ER), where old data samples are mixed with current ones during training. Building 803 upon this concept, Robins Robins (1995) explored pseudorehearsal techniques. More recent meth-804 ods like Meta-Experience Replay (MER) Riemer et al. (2018) reformulate ER within a meta-learning 805 framework, aiming to enhance knowledge transfer between past and present tasks while reducing in-806 terference. Gradient-based Sample Selection (GSS) Aljundi et al. (2019b) modifies ER by selecting 807 optimal examples for storage in the memory buffer, improving retention of past knowledge. Another method, Hindsight Anchor Learning (HAL) Chaudhry et al. (2021), augments ER with an additional 808 goal to prevent forgetting key data points. Gradient Episodic Memory (GEM) Lopez-Paz & Ranzato (2017b) and its more efficient variant Averaged-GEM (A-GEM) Chaudhry et al. (2018b) use previ810 ous training data to impose optimization constraints on the current update, ensuring better retention 811 of learned information. Additionally, Yoon et al. (2019) introduced Additive Parameter 812 Decomposition (APD), an architectural approach that represents the parameters for each task as a 813 sum of task-shared and task-adaptive parameters. APD ensures scalability and order-robustness by 814 preventing catastrophic forgetting and addressing order-sensitivity through parameter decomposition. Lastly, Aljundi et al. (2019a) introduce task-free continual learning, eliminating 815 the need for task boundaries and enabling more flexible adaptation to new tasks without explicit task 816 identifiers. 817

These methods, while effective in certain scenarios, rely heavily on storing and replaying data from previous tasks, which may not be feasible due to privacy concerns or memory constraints. In contrast, our approach does not require storing raw data from previous tasks. Instead, we utilize pseudolabels generated by the frozen previous model and introduce a novel regularization term to preserve the structure of shared representations. This enables the model to retain and improve upon prior knowledge without explicit rehearsal.

824

825

**Knowledge Distillation** Methods leveraging Knowledge Distillation Hinton (2015) address the 826 issue of forgetting by using a previous iteration of the model as a teacher. Learning Without For-827 getting (LwF) Li & Hoiem (2017a) generates a softened version of the model's current outputs 828 on new data at the onset of each task, minimizing output drift throughout training. iCaRL Re-829 buffi et al. (2017) combines distillation with replay techniques, using a memory buffer to train a 830 nearest-mean-of-exemplars classifier while applying a self-distillation loss to preserve learned rep-831 resentations across tasks. Moreover, Li et al. Li et al. (2019) proposed a continual learning method 832 tailored for sequence-to-sequence tasks, leveraging compositionality to enable knowledge transfer 833 and prevent catastrophic forgetting. Their approach extends traditional label prediction continual 834 learning methods to handle more complex tasks like machine translation and instruction learning.

While these methods use knowledge distillation to maintain performance on old tasks, they typically focus on preserving output logits or feature representations without considering the underlying relational structure between data points. Our method extends this idea by not only preserving the output predictions via pseudolabels but also maintaining the pairwise relationships in the representation space through our Dynamically Weighted Distance Preservation (DWDP) loss. This helps in better retaining the learned structure and prevents the model from drifting away from previously acquired knowledge.

842 843

844 **Regularization-based techniques** These methods modify the loss function to include a penalty that restricts changes to the model's parameters. Examples include Elastic Weight Consolidation 845 (EWC) Duncker et al. (2020), its online variant (oEWC) Kirkpatrick et al. (2017), Synaptic Intel-846 ligence (SI) Zenke et al. (2017), and Riemannian Walk (RW) Chaudhry et al. (2018a). In contrast, 847 architectural methods such as Progressive Neural Networks (PNN) Rusu et al. (2016) incrementally 848 expand the model by adding new networks for each task, which leads to increased memory usage. 849 To address this, methods like PackNet Mallya & Lazebnik (2018) and Hard Attention to the Task 850 (HAT) Serra et al. (2018) reuse the same architecture for multiple tasks, dynamically allocating 851 resources to prevent performance degradation. Recent advances include a generalized framework 852 with additional loss functions proposed by Wang et al. Wang et al. (2024). Another promising 853 architectural method is task-conditioned hypernetworks Von Oswald et al. (2019), which generate 854 weights for the target network based on task identity. These hypernetworks do not need to recall all input-output relationships for previously seen tasks, as they instead rehearse task-specific weight 855 realizations. Moreover, Adel et al. (2019) introduced Continual Learning with Adaptive 856 Weights (CLAW), which employs a probabilistic modeling approach to adaptively identify which 857 parts of the network should be shared across tasks in a data-driven manner. This method balances 858 between modeling each task separately to prevent catastrophic forgetting and sharing components 859 to allow transfer learning and reduce model size. 860

Our approach differs from these methods as we do not rely on parameter regularization or expand ing architectures. Instead, we focus on preserving the learned representations and their relational
 structure between tasks through the DWDP loss, which provides a more scalable solution without
 incurring additional memory overhead.

# A.3 CONTINUAL AND MULTITASK LEARNING

Our work distinguishes itself from existing approaches in CMTL by introducing a new problem domain where new tasks and their associated datasets arrive incrementally. In this setting, the model is not only required to adapt to new tasks but also to utilize these new datasets to enhance its performance on previous tasks. Specifically, when new data for additional tasks becomes available, it is used to further train the existing model. This training process enables the model to reinforce and improve its understanding of prior tasks, effectively allowing it to remember and perform better on both past and current tasks.

Building upon the extensive research in multitask learning (MTL) Caruana (1997); Sener & Koltun (2018); Yu et al. (2020); Navon et al. (2022); Li & Hoiem (2017a) and continual learning (CL) Ratcliff (1990); Robins (1995); Riemer et al. (2018); Aljundi et al. (2019b); Chaudhry et al. (2021); Li & Hoiem (2017a), the emerging field of continual multitask learning seeks to bridge the two paradigms to effectively manage performance across sequential and concurrent tasks Wu et al. (2023).

One of the most related works to ours, Mirzadeh et al. Mirzadeh et al. (2020), focus on the linear mode connectivity between solutions obtained through sequential and simultaneous training. While they demonstrate that a linear path of low error exists for more than twenty tasks and introduce algorithms like Mode Connectivity SGD (MC-SGD) to enhance continual learning, their work does not address the use of *new tasks* to improve performance on previous ones, particularly using a similar setup to traditional continual learning, which means their works fit more on the realm of CL.

Similarly, Liao et al. Liao et al. (2022) propose MUSCLE, a multitask self-supervised continual learning framework designed to pre-train deep models on diverse X-ray datasets. This work, similar to ours, operates in the domain of medical imaging to process classification and segmentation in different body areas. However, their work differs from ours because their focus is on pre-training the model on different tasks for better generalization, which they refer to as "multitask continual learning." We specifically differentiate our CMTL approach from theirs in that our tasks are seen *iteratively*; we do not have access to all tasks at the same time, and the tasks themselves could be orthogonal to previously seen tasks.

In summary, our approach introduces a novel aspect to CMTL by leveraging new tasks and their data not only to learn the new tasks but also to generalize on prior tasks, all within an iterative framework where tasks arrive sequentially and are potentially unrelated. This sets our work apart from existing CMTL methods, which typically do not utilize new tasks to enhance previous ones in this manner.

#### 918 B LEARNING WITH PRESERVING ALGORITHM OVERVIEW 919

In this section, we present the pseudocode for our algorithm presented in Section 3.

921 922 Algorithm 1 Learning with Preserving (LwP) 923 1: Input: Sequence of tasks  $\{D_t\}_{t=1}^T$ , hyperparameters  $\lambda_c$ ,  $\lambda_o$ ,  $\lambda_d$ 924 2: **Output:** Final model parameters  $\theta^{[T]}$ 925 3: Initialize initial model parameters  $\theta^{[0]}$ 926 4: **for** t = 1 to *T* **do** 927 Initialize current model parameters:  $\theta^{[t]} \leftarrow \theta^{[t-1]}$ 5: 928 Add new task-specific layer  $g_{\theta_t}$  for task t to  $\theta^{[t]}$ 6: 929 Freeze previous model parameters  $\theta^{[t-1]}$ 7: 930 for each minibatch  $\{(\boldsymbol{x}_i, y_i^t)\}_{i=1}^N$  from  $D_t$  do 8: 931 Compute shared representations:  $\boldsymbol{z}_{i}^{[t]} = f_{\theta^{[t]}}(\boldsymbol{x}_{i})$ 9: 932 Compute output for current task:  $\hat{y}_i^t = g_{\theta^{[t]}}(\boldsymbol{z}_i^{[t]})$ 933 10: 934 Compute representations from frozen model:  $z_i^{[t-1]} = f_{a^{[t-1]}}(x_i)$ 11: 935 12: for o = 1 to t - 1 do 936 13: Compute outputs for previous task o: 937 Current model output:  $\hat{y}_i^o = g_{\theta^{[t]}}(\boldsymbol{z}_i^{[t]})$ 14: 938 Pseudolabel from frozen model:  $\tilde{y}_i^o = g_{\theta_{\alpha}^{[t-1]}}(\boldsymbol{z}_i^{[t-1]})$ 939 15: 940 16: end for 941 Compute loss for new task:  $\mathcal{L}_{cur} \leftarrow \mathcal{L}_{cur}(y_i^t, \hat{y}_i^t)$ 17: Compute loss for old tasks:  $\mathcal{L}_{old} \leftarrow \sum_{o=1}^{t-1} \mathcal{L}_{old}(\tilde{y}_i^o, \hat{y}_i^o)$ 942 18: 943 19: Compute dynamic mask  $m_{ij}$ : 944  $m_{ij} = \begin{cases} 1, & \text{if } y_i^t = y_j^t, \\ 0, & \text{otherwise} \end{cases}$ 945 946 947 Compute DWDP loss: 20: 948 949  $\mathcal{L}_{\text{DWDP}} \leftarrow \frac{1}{N^2} \sum_{i=1}^{N} \sum_{i=1}^{N} m_{ij} \left( d(\boldsymbol{z}_i^{[t-1]}, \boldsymbol{z}_j^{[t-1]}) - d(\boldsymbol{z}_i^{[t]}, \boldsymbol{z}_j^{[t]}) \right)^2$ 950 951 952 21: Compute total loss: 953  $\mathcal{L}_{lwp} \leftarrow \lambda_{c} \mathcal{L}_{cur} + \lambda_{o} \mathcal{L}_{old} + \lambda_{d} \mathcal{L}_{DWDP}$ 954 955 Update parameters  $\theta^{[t]}$  by minimizing  $\mathcal{L}_{lwp}$ 22: 956 end for 23: 957 24: end for 958 959 960 961 962 963 964 965 966 967 968 969 970 971

#### 972 C JUSTIFICATION ON USING EUCLIDEAN DISTANCE 3.2

Here, we show that preserving the squared Euclidean distances between the data points in Z and Z' is sufficient to achieve the same effect.

**Squared Euclidean Distance Preservation** We define the squared Euclidean distance between two points  $z_i$  and  $z_j$  as:

$$D_{ij}(Z) = ||z_i - z_j||^2.$$
(17)

Similarly, we compute  $D_{ij}(Z')$  for Z'.

Our goal is to minimize the difference between the squared distances in Z and Z', which we formalize with the following loss function:

$$\mathcal{L}_{\text{dist}}(Z, Z') = \sum_{i=1}^{n} \sum_{j=1}^{n} \left( \|z_i - z_j\|^2 - \|z'_i - z'_j\|^2 \right)^2.$$
(18)

Minimizing  $\mathcal{L}_{dist}$  with respect to Z' encourages the squared distances between all pairs of points in Z' to match those in Z:

$$||z'_i - z'_j||^2 \approx ||z_i - z_j||^2, \quad \forall i, j.$$
 (19)

Since the exponential function is Lipschitz continuous on compact subsets, small changes in the squared distance result in small changes in the kernel value. Specifically, if the squared distances are preserved within a small error  $\epsilon > 0$ :

$$\left| \|z_i - z_j\|^2 - \|z_i' - z_j'\|^2 \right| < \epsilon,$$
(20)

then the difference in the kernel values can be bounded:

$$\left|k(z_i, z_j) - k(z'_i, z'_j)\right| = \left|\exp\left(-\frac{\|z_i - z_j\|^2}{2\sigma^2}\right) - \exp\left(-\frac{\|z'_i - z'_j\|^2}{2\sigma^2}\right)\right|$$
(21)

$$\leq \frac{1}{2\sigma^2} \exp\left(-\frac{\min(\|z_i - z_j\|^2, \|z_i' - z_j'\|^2)}{2\sigma^2}\right) \left|\|z_i - z_j\|^2 - \|z_i' - z_j'\|^2\right|$$

$$\leq \frac{1}{2\sigma^2}\epsilon \tag{23}$$

$$\leq L_k \epsilon$$
 (24)

where  $L_k$  is a Lipschitz constant dependent on  $\sigma$ .

Therefore, preserving the squared Euclidean distances between Z and Z' implies that the Gaussian kernel matrices K(Z) and K(Z') are approximately equal:

$$k(z_i, z_j) \approx k(z'_i, z'_j), \quad \forall i, j.$$
(25)

# 1013 D ADDITIONAL DETAILS ON EXPERIMENTAL RESULTS

### 1015 D.1 HYPERPARAMETERS

In the following section, we provide an extensive description of the hyperparameters utilized during the training phase. Across all datasets and models, the Adam optimizer Kingma & Ba (2017) was employed universally. For the CelebA and FairFace datasets, a consistent learning rate of 0.0001 was maintained, coupled with a batch size configuration of 256. In contrast, for the PhysiQ dataset, a higher learning rate of 0.01 was utilized alongside a smaller batch size of 32. Furthermore, we adhered to fixed model-specific hyperparameters for all datasets and models to ensure uniformity and consistency, including the LwP parameters. In the case of LwP, the parameters set as follows:  $\lambda_n$  as a value of 1,  $\lambda_o$  as a value of 1, and  $\lambda_d$  with a value of 0.01. Additionally, the 10 tasks used for CelebA are wearing lipsticks, smiling, mouth slightly open, high cheekbones, attractive, heavy makeup, male, young, wavy hair, and straight hair. PhysiQ dataset includes three attributes assessing exercise quality: stability, range of motion, and exercise variation.

Details of all models and their hyperparameter selection have been documented in the codebase. For
 in-depth understanding and additional information, please consult our code repository available at
 [ANONYMOUS LINK].

# 1030 D.2 COMPARISON WITH MTL METHODS

1031 1032 1033

1048

1049

1029

Table 3: Comparison of Accuracy across Different Models and Datasets

Method Type	Model	CelebA	PhysiQ	FairFace
STL	-	$72.230\pm7.297$	$87.167 \pm 10.102$	$64.435 \pm 3.660$
MTL	MTL PCGrad IMTL NashMTL	$\begin{array}{c} \textbf{76.526} \pm \textbf{7.616} \\ 75.506 \pm 8.146 \\ 76.280 \pm 7.248 \\ 75.506 \pm 8.146 \end{array}$	$\begin{array}{c} \textbf{93.536} \pm \textbf{5.739} \\ 91.910 \pm 8.491 \\ 92.661 \pm 6.617 \\ 91.518 \pm 7.118 \end{array}$	$71.418 \pm 4.169 70.061 \pm 4.892 71.399 \pm 3.887 71.607 \pm 3.577$
CMTL	LwP	$73.484\pm8.019$	$88.242 \pm 12.010$	$68.545 \pm 4.454$

All MTL approaches utilize the same model architecture as LwP. Despite being supplied with all labels for every input data point, the amount of training samples for MTL models matches that seen by CL models per task iteration. Aligning with earlier studies, MTL approaches frequently represent the upper bound for all CL models. An interesting discovery is that all MTL models deliver nearly identical performance on the benchmark.

## D.3 ACCURACY PROGRESSION FOR EACH TASK ITERATION



10

1066

Figure 8: Confusion matrices for different models on the PhysiQ dataset

Figures 8 and 9 illustrate that the application of LwP reduces the issue of catastrophic forgetting in the PhysiQ and FairFace datasets as well. This effect is particularly pronounced when applied to datasets with a large number of samples, such as Fairface and CelebA, in comparison to smaller datasets such as PhysiQ. These observations imply that LwP is a scalable and effective solution to mitigate catastrophic forgetting in continual multitask learning models. We further investigate the effect of the number of training samples on performance in D.4.

1073

1074 D.4 INFLUENCE OF TRAINING SAMPLE

We include experiment results on the influence of number of training samples to the performance, as shown in Fig. 10a. It shows that our approach outperforms others from 1000 labels and onward, when trained and tested on the PhysiQ dataset.

1079 Fig. 10b illustrates the Expected Calibration Error (ECE) Nixon et al. (2019) for each model in relation to the number of training samples. The ECE quantifies how much confidence a model deviates



We show that LwP demonstrates a faster improvement per iteration given the same batch size as other CL and MTL models. Here, we plot the evolution of average accuracy across all tasks seen on the test set over training iteration for the top performing CL and MTL baselines along with LwP.
We used the CelebA dataset with 5 task splits. To make a fair comparison with CL models, MTL models were trained on the amount of train data that CL models saw in each iteration with access to all 5 tasks. The accuracies of MTL models are calculated up to what CL models have learned so far.

Fig. 11a shows how quickly each CL and MTL model learns the first task. This can be understood as the speed at which the models acquire knowledge when they have no prior information to "recall".
As shown, LwP learns consistently faster per iteration compared to CL and MTL baselines. As MTL models simultaneously learn multiple labels, their convergence per iteration is generally slower compared to CL models in this setting.

1145 Conversely, Fig. 11b illustrates a case where MTL models are trained from the beginning with la-1146 bels available for all t tasks, whereas CL models, having been pretrained on t-1 tasks, must now 1147 incrementally learn the  $t^{th}$  task while maintaining performance on old tasks. This configuration is 1148 crucial in real-world scenarios where the cost of labeling data typically exceeds that of data collec-1149 tion, prompting the decision to gather more partially labeled data rather than re-labeling existing 1150 data. Analogous to the prior scenario, the progression of test accuracy over iterations demonstrates 1151 that LwP consistently exceeds other CL models and exhibits performance that is competitive with MTL models, which are considered the upper bound for continual learning. This highlights the com-1152 parative benefit of LwP when users face the choice between relabeling existing data and obtaining 1153 new data with different labels. 1154

### 1155

### 1156

### D.6 EFFECT OF MODEL PARAMETERS AND IMAGE SIZES ON TRAINING PERFORMANCE

1157 1158 1159

Table 4: Accuracy Percentage Comparison Across Models on CelebA Dataset

Method Type	Model	<b>ResNet50</b> (32 × 32)	<b>ResNet101</b> (32 × 32)	<b>ResNet50</b> ( $224 \times 224$ )
	LwF	59.277 ± 11.920	58.279 ± 11.202	$60.012 \pm 14.448$
	oEWC	$66.975 \pm 10.110$	$67.159 \pm 10.506$	$68.511 \pm 13.352$
	ER	$65.335 \pm 9.298$	$65.646 \pm 8.784$	$65.973 \pm 14.729$
	SI	$66.698 \pm 10.030$	$67.456 \pm 9.880$	$67.747 \pm 13.754$
CI	GSS	$65.926 \pm 13.120$	$65.587 \pm 13.142$	$69.817 \pm 18.771$
CL	FDR	$61.753 \pm 11.943$	$61.720 \pm 12.017$	$65.225 \pm 15.545$
	DER	$62.105 \pm 12.114$	$63.797 \pm 10.774$	$69.859 \pm 12.690$
	DERPP	$62.814 \pm 11.071$	$62.957 \pm 11.577$	$68.102 \pm 13.557$
	DVC	$67.084 \pm 10.380$	$65.340 \pm 11.427$	$70.921 \pm 13.823$
	OBC	$64.220 \pm 11.237$	$66.058 \pm 10.370$	$69.319 \pm 13.607$
CMTL	LwP	$\textbf{67.388} \pm \textbf{11.125}$	$\textbf{69.432} \pm \textbf{10.416}$	$\textbf{85.064} \pm \textbf{5.388}$
	Method Type CL CMTL	Method TypeModelLwF oEWCER SICLGSS FDR DER DERPP DVC OBCCMTLLwP	$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$

1171 1172

Table 4 illustrates that the LwP method scales effectively with increased input resolution and model 1173 size. We find that preserving the Gaussian kernel, as shown in eq. 3, results in improved performance 1174 on larger scales, especially with respect to input resolution. In the ResNet50 benchmark utilizing a 1175 224x224 image size, LwP notably surpasses other baselines by achieving an 85% accuracy, which is 1176 about 15% percentage points greater than the runner-up. This suggests that, as the input allows the 1177 model to create more insightful representations, LwP becomes increasingly advantageous because it 1178 can maintain these representations. We also note that the bigger models with the same input size are 1179 not performing as well as the one with resnet18. This is due to the fact that the inputs do not have enough information to capture generalized patterns, resulting in overfitting. 1180

1181

# 1182 D.7 TRAINING FROM MTL TO CL

1184 We initially train the model on the first five tasks using a MTL setting, employing ResNet18 as 1185 the encoder with input images of size  $64 \times 64 \times 3$ . After completing the MTL phase, we extract 1186 the encoder and freeze its weights. This frozen encoder is then used to train classifiers for the first 1187 five tasks in a continual learning CL setting with various models. Subsequently, we train the entire 1186 models for the last five tasks under the same CL framework, utilizing the same frozen encoder on the



1211 tasks provides a robust foundation for learning new tasks in a continual fashion. Our method effectively mitigates catastrophic forgetting by preserving essential features learned during the MTL
1213 phase while adapting to new tasks, given the continual tasks are shorter now. This balance between stability and plasticity still allows LwP to maintain high accuracy in the continual learning tasks.

Table 5: Accuracy Percentage Comparison Across Models on CelebA Dataset, Trained on MTL onfirst 5 tasks then CL on last 5 tasks

1218			
1219	Method Type	Model	<b>ResNet18</b> (64 × 64)
1220		LwF	$74.057 \pm 11.364$
1221		oEWC	$82.250 \pm 6.362$
1222		ER	$77.245 \pm 8.434$
1993		SI	$82.194 \pm 6.460$
1004	CI	GSS	$80.563 \pm 8.239$
1224	CL	FDR	$81.271 \pm 7.738$
1225		DER	$81.010 \pm 8.674$
1226		DERPP	$78.177 \pm 9.532$
1227		DVC	$81.387 \pm 7.821$
1228		OBC	$80.516 \pm 8.446$
1229	CMTL	LwP	$\textbf{83.652} \pm \textbf{7.069}$
1230			

Moreover, the lower standard deviation in LwP's performance indicates consistent results across different runs, highlighting the reliability of our approach. The results confirm that combining MTL pre-training with our proposed CL strategy enhances the model's ability to generalize and adapt to new tasks without compromising performance on previously learned tasks.

Similarly in Figure 12, we averaged the results across task iterations to evaluate performance over time. Our method, LwP, demonstrates minimal accuracy loss when training on new tasks, high-lighting its performance against forgetting. The standard deviation—represented as 20% of the total for visualization purposes—remains low, indicating consistent performance. Although there is a slight increase in standard deviation during later tasks, suggesting a potential drop in accuracy due to forgetting, LwP still preserves knowledge at a superior level compared to other baselines. Even with the first five tasks trained in a multitask setting, our method maintains the best overall accuracy, outperforming other models in preserving learned information.