# PART: PARALLEL LEARNING TOWARDS ROBUST AND TRANSPARENT AI

**Mahsa Paknezhad**
**Hamsawardhini Rengarajan**
**Chenghao Yuan**
Bioinformatics Institute, A*STAR,
30 Biopolis Street, 07-01, Matrix
Singapore, 138671
mahsap@bii.a-star.edu.sg

**Savitha Ramasamy**
**Manas Gupta**
**Sujanya Suresh**
I2R, A*STAR,
1 Fusionopolis Way, 21-01 Connexis
Singapore 138632

**Hwee Kuan Lee**
Bioinformatics Institute, A*STAR,
30 Biopolis Street, 07-01, Matrix
Singapore, 138671
National University of Singapore
Singapore, 119077
Singapore Eye Research Institute
20 College Road
Singapore, 169856
leehk@bii.a-star.edu.sg

February 24, 2022

## ABSTRACT

This paper takes a parallel learning approach for robust and transparent AI. A deep neural network is trained in parallel on multiple tasks, where each task is trained only on a subset of the network resources. Each subset consists of network segments, that can be combined and shared across specific tasks. Tasks can share resources with other tasks, while having independent task-related network resources. Therefore, the trained network can share similar representations across various tasks, while also enabling independent task-related representations. The above allows for some crucial outcomes. (1) The parallel nature of our approach negates the issue of catastrophic forgetting. (2) The sharing of segments uses network resources more efficiently. (3) We show that the network does indeed use learned knowledge from some tasks in other tasks, through shared representations. (4) Through examination of individual task-related and shared representations, the model offers transparency in the network and in the relationships across tasks in a multi-task setting. Evaluation of the proposed approach against complex competing approaches such as Continual Learning, Neural Architecture Search, and Multi-task learning shows that it is capable of learning robust representations. This is the first effort to train a DL model on multiple tasks in parallel. Our code is available at https://github.com/MahsaPaknezhad/PaRT.

## 1 Introduction

Training deep learning (DL) models on multiple heterogeneous tasks is one of the main steps in the direction of offering robust and generalizable AI solutions. There exists an extensive body of research on multi-task learning [43, 44, 45]. However, the effectiveness of these methods in learning heterogeneous tasks with different data distributions is under question. Research in continual learning (CL) [46, 47] is recently gaining a lot of attention as it enables learning a sequence of tasks and leveraging on shared representations across those tasks while avoiding catastrophic forgetting of the previously learnt tasks. However, certain issues still remain in CL algorithms such as: 1) existing CL solutions have not been able to fully address catastrophic forgetting yet. 2) some CL solutions require more network resources to perform architectural adaptation, and 3) it is difficult to draw conclusions regarding representational similarities across different tasks with the current CL solutions. In this paper, we take a parallel approach to learning heterogeneous tasks using a single network. Parallel learning enables sharing representations across tasks, while also having task-related representations per individual tasks. Once a new task with a different input data distribution is introduced, the proposed algorithm can learn the new task alongside the existing tasks while sharing representations with the existing tasks. Parallel learning also draws a window of transparency through which we can see representational similarities
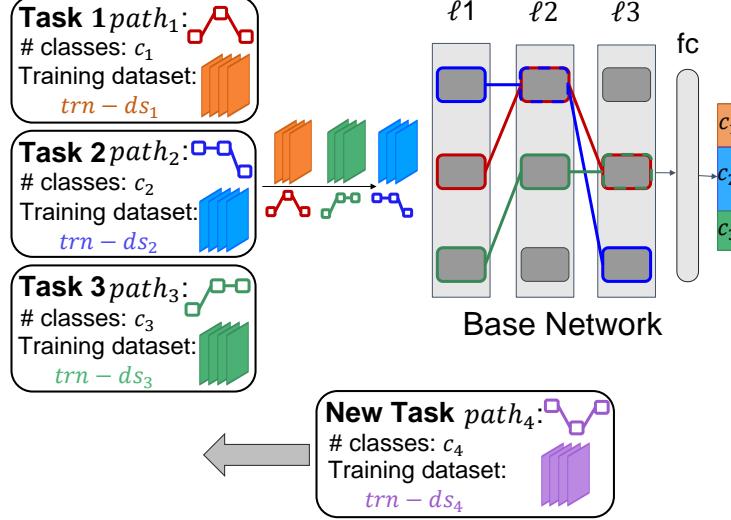
Figure 1: Figure shows a diagram of how parallel learning is carried out. Parallel learning encompasses a base neural network on which heterogeneous tasks are trained in parallel. Certain modules in the base network are assigned to each task and learning that task is done using those modules. Training on multiple tasks is performed by randomly sampling batches from the training dataset of those tasks. Once a new task is introduced, it will be added to the bucket of tasks on which the base network is trained.

as well as the amount of resource sharing in the network between tasks. CL is aimed at learning a sequence of tasks, while parallel learning is capable of learning multiple tasks in parallel. Furthermore, it can leverage on existing representations to learn new tasks swiftly. Parallel learning will be useful for learning multiple tasks from multiple sources of data within a single institution such as banks and hospitals. Parallel learning will equip these entities with a single DL model that learns from multiple datasets to produce multiple heterogeneous outcomes of interest with almost the same performance compared to a DL model trained on each dataset and outcome of interest. Our contributions are summarized below:

- We develop an algorithm for training a DL model on multiple heterogeneous tasks with diverse input data distributions in parallel.

- We demonstrate that parallel learning is a simple but effective approach that has improved performance over multi-task learning, sequential CL algorithms, single task learning algorithms and the state-of-the-art algorithms in other domains.

- Using our algorithm, we are able to see task-related and shared representations across tasks in a multi-task setting. We demonstrate that trained models using parallel learning learn more similar representations across tasks.

## 2 Related Work

**Multi-task learning, Transfer learning and Meta-learning**: Multi-task learning [19, 20, 21, 22, 23] and transfer learning [24, 25] algorithms aim to reduce the required amount of data and the need for different models by using a single model to learn different tasks. Multilinear Relationship Network (MRN) [19] is one of the recently proposed multi-task learning algorithms. MRN jointly learns transferable features and multilinear relationships of tasks and features. It is important to note that the parallel learning is not a multi-task learning algorithm. In multi-task learning, a DL model is trained to produce multiple outcomes of interest on an input dataset. In parallel learning, however, different outcomes of interest are defined on different input datasets. Also, our proposed DL model is different from transfer learning. Parallel learning shares a portion of the neurons across multiple tasks and trains them on those tasks simultaneously whereas in transfer learning, neurons are initialized with values learned on an earlier task and are trained on a new task. Another group of algorithms called Meta-learning [26, 27, 28] aim to reduce the required amount of data to learn a new task using trained models on other tasks. Although meta-learning algorithms have been proposed to learn tasks with heterogeneous attribute spaces [37], these methods are yet to be extended to heterogeneous
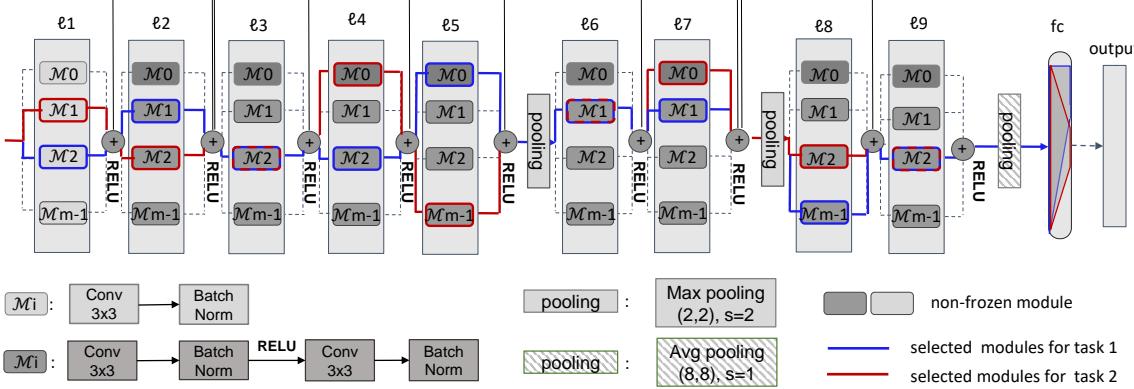
Figure 2: Figure shows a simple architecture for the base network in our parallel learning algorithm. One row of modules is equal to a ResNet-18 architecture in this figure. For each task, a set of $N$ modules is randomly chosen in each layer and are assigned to the task. Training on the task is performed only on the selected modules for that task.

tasks. As will be shown in section 4, parallel learning can be used to learn heterogeneous tasks with diverse input data distributions.

**Neural architecture search (NAS)**: Studies in the NAS domain [8, 11, 10, 12, 14, 15, 16, 17, 18] target automatic and fast design of DL models for the task in hand. pDarts [6] is an optimized version of the Darts algorithm [13] which is one of the most benchmarked algorithms in the NAS domain. In pDarts, a network is formed by stacking multiple cells together. The architecture within a cell is derived using a gradient-based approach. ENAS [9] is another algorithm which trains a controller using policy gradient to search an optimal network architecture. A multi-agent approach to NAS where agents control parts of the network and coordinate to reach the optimal architectures is proposed as MANAS [7]. However, NAS algorithms generate a separate DL model for each task. Consequently, while they reduce the time needed to develop a model, they are not considered to be memory-efficient and they do not take advantage of learned representations from other tasks.

**Continual learning (CL)**: A large amount of work has been done in CL [1, 2, 39, 4, 5] to train a DL model on new tasks while retaining its performance on preceding tasks. Training on tasks is done sequentially in many such algorithms. Once training on a task finishes, the contributing neurons to that task are frozen or are allowed to change slightly to avoid catastrophic forgetting. In Elastic Weight Consolidation (EWC) [2], for instance, a regularization term pulls back parameter values to their old learned values on previous tasks. Hung et al. [5] train multiple tasks sequentially by expanding the neural network if necessary. They add a model compression step before training the DL model on a new task. A closer group of algorithms to our proposed method are CL algorithms with generative replay algorithms such as the work by Shin et al. [39], van de Ven et al. [30] and other methods [40, 41, 42] where a generator is trained to synthesize samples for previous task, or a small representative set of data from previously learned tasks called a coreset [49, 48, 49] is kept. The samples are mixed with the training data of the new task and are used during training. Random Path Selection (RPS) [4], is a good example which deploys coresets. The main reason for training a generator or keeping coresets is to make the algorithm more memory-efficient. We, however, target entities such as hospitals at which clinical data will always be kept. As a result, memory-efficiency is not a major concern. What is important, instead, is to train a DL model which can learn to do multiple important tasks such as cancer grading and disease prediction with high accuracy and limited amount of data, which is the main objective of the proposed approach in this paper.

## 3 Method

In this section, we introduce PaRT: a parallel learning algorithm for robust and transparent AI. This is the first effort to train a DL model on multiple tasks in parallel. We draw inspirations for our algorithm from existing efforts in the CL domain. PaRT encompasses an algorithm for parallel learning of multiple tasks using a base network architecture. A simple diagram of our algorithm is show in Fig. 1. In the following, we explain our proposed algorithm.

Fig. 2 shows the architecture we deployed as our base network. The base network consists of multiple layers, with a total of $M$ modules in each layer. One row in the base network can be equivalent to any standard network architecture such as ResNet-18 or ResNet-50 architectures. Let us assume there are $k$ tasks to learn using the proposed parallel

---

**Algorithm 1** Parallel Learning

---

    **input**: $k$ : The number of tasks
    $(trn\text{-}ds_i, val\text{-}ds_i)$: The training and validation datasets for the $i^{th}$ task
    $c_i$: The number of classes for the $i^{th}$ classification task
    $C_\theta$: The base network
    $path_i$: A randomly selected set of modules for the $i^{th}$ task
    **output**: $C_\theta$ trained on the $k$ tasks

1:  **procedure** PARALLELLEARNING
2:     $\{trn\text{-}ds_1^{new}, trn\text{-}ds_2^{new}, ..., trn\text{-}ds_k^{new}\} \leftarrow$ Update the $trn\text{-}ds_i$ for $i = 1, ..., k$ to have equal size by oversampling
3:     **for each** training epoch
4:         **while** there exists untrained batches
5:             $bucket \leftarrow \{trn\text{-}ds_1^{new}, ..., trn\text{-}ds_k^{new}\}$
6:             Randomly sample a dataset $trn\text{-}ds_i^{new}$ from the $bucket$ without replacement
7:             Sample a few batches $batch\text{-}set$ from $trn\text{-}ds_i^{new}$
8:             Set $start = \sum_{n=1}^{i-1} c_n$ and $end = \sum_{n=1}^{i} c_n$
9:             **for** $batch_j$ from $batch\text{-}set$
10:               Train $C_\theta$ on $batch_j$ over $path_i$ using output neurons $[start, end]$
11: **end procedure**

---

learning algorithm and the base network. At the beginning, the training datasets for all tasks are oversampled to have the same size. For each task, $N$ modules will be randomly selected from the existing $M$ modules in each layer independent of the other layers and are assigned to the task for training. The parameters $N$, $N$ are hyperparameters. Some modules may be selected by multiple tasks. As a result, such modules will be shared and trained on those tasks. Let us define the randomly selected set of modules at all the layers in the base network for the $i^{th}$ task as $path_i$. The training process on the base network is as follows: In each training epoch, a few batches called $batch\text{-}set$s are taken from all k training datasets. We then trained these $batch\text{-}set$s in a random order. Training on batches from $task i$ is performed on the modules specified in $path_i$ in the base network. We then repeated the process until all samples have been taken from the training sets and used for training, and one epoch finishes. The size of the $batch\text{-}set$ is also a hyperparameter. The outputs of modules selected in the same layer for task i are summed together, and the sum is fed to the modules selected by task i in the next layer(s). No module is frozen during training, and tasks are trained on their corresponding modules. The number of neurons in the output layer is the sum of the number of classes for all tasks. The parallel nature of training in the proposed algorithm negates catastrophic forgetting. We use $ADAM$ as our optimization and $CrossEntropy$ for our loss function. Our parallel learning algorithm for training the base network on these $k$ tasks is shown in Algorithm 1. Our code will be available online.

Once training the base network finishes, validation for the $i^{th}$ task is performed by passing samples from the validation dataset of the task, $val\text{-}ds_i$, through $path_i$ and comparing the predicted labels with ground truth labels. We compare our algorithm to sequential learning, a basic version of Continual Learning. In sequential learning, the model is trained on tasks one after another. Similar to parallel learning, a set of modules defined by $path_i$ are assigned to each task $i$. Once training for task $i$ finishes, the modules in $path_i$ will be frozen and will no longer be trainable for subsequent tasks even if those frozen modules are assigned to a new task. We also benchmark our algorithm against single task learning which is defined as training $path_i$ in the base network from scratch on only task $i$. Single task learning represents the achievable performance for each task on our base network. In the next section, we will apply our parallel learning algorithm to learning multiple homogeneous as well as heterogeneous tasks. We will see certain changes have to be made to our base network architecture when Parallel Learning is applied on heterogeneous tasks.

We use the representation similarity metric proposed by Kornblith et al. [38] to compare the learned task-related and shared representations by models trained using PaRT with the learned representations by models trained on each task using single task learning. This metric is called centered kernel alignment (CKA) and is closely related to canonical correlation analysis (CCA). CKA is a stronger metric compared to CCA and other similarity metrics as it allows measuring meaningful similarities between representations of higher dimension than the number of data points.

## 4 Experimental Evaluation and Analysis

We designed three experiments with increasing heterogeneity in tasks. In the following, we will explain each experiment and analyze the results for the experiments.
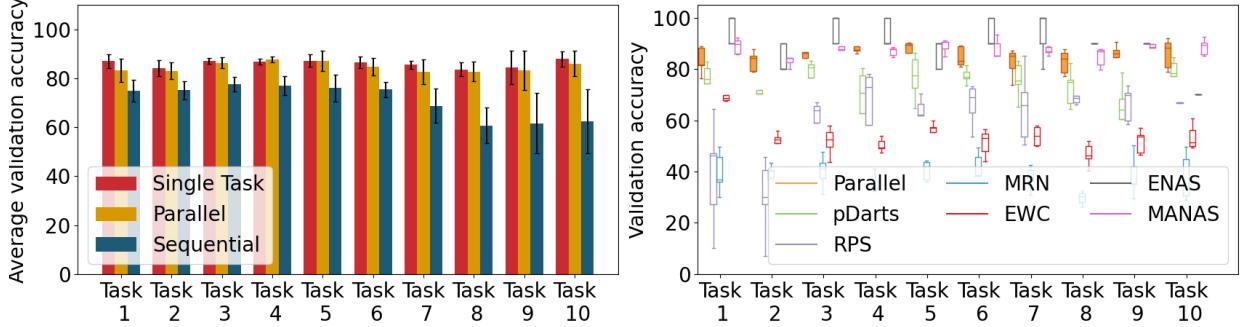
Figure 3: Mean validation accuracy of DL models trained on 10 tasks defined on the CIFAR100 [36] dataset. Each task consists of classifying 10 randomly selected classes in the CIFAR100 dataset into 10 categories. (Left) compares the performance of PaRT with single task learning and sequential learning methods. (Right) shows performance of PaRT and those of the competing algorithms. Five DL models were trained for each algorithm.

| Learning Method | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 | Task 6 | Task 7 | Task 8 | Task 9 | Task 10 | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 86.96 | 84.12 | 87.04 | 86.88 | 87.24 | 86.38 | 85.58 | 83.68 | 84.52 | 87.85 | 86.03 |
| | ± 2.89 | ± 3.23 | ± 1.11 | ± 1.23 | ± 2.64 | ± 2.37 | ± 1.61 | ± 2.95 | ± 6.85 | ± 3.13 | |
| Single Task | Tasks | Tasks | Tasks | Tasks | Tasks | Tasks | Tasks | Tasks | Tasks | Tasks | |
| | $All\backslash 1$ | $All\backslash 2$ | $All\backslash 3$ | $All\backslash 4$ | $All\backslash 5$ | $All\backslash 6$ | $All\backslash 7$ | $All\backslash 8$ | $All\backslash 9$ | $All\backslash 10$ | |
| | 10.11 | 9.58 | 10.03 | 9.36 | 10.36 | 10.06 | 10.01 | 9.39 | 9.5 | 9.89 | 9.83 |
| | ±1.58 | ±1.53 | ±1.31 | ±0.82 | ±0.83 | ±1.31 | ±1.29 | ±1.29 | ±1.66 | ±3.49 | |
| Learning Method | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 | Task 6 | Task 7 | Task 8 | Task 9 | Task 10 | |
| Sequential | 74.96 | 75.06 | 77.48 | 76.96 | 76.02 | 75.40 | 68.80 | 60.76 | 61.66 | 62.42 | 70.95 |
| | ± 4.54 | ± 3.66 | ± 3.01 | ± 3.86 | ± 5.49 | ± 3.04 | ± 7.10 | ± 7.32 | ± 12.41 | ± 13.23 | |
| Parallel | 80.82 | 75.96 | 79.08 | 73.06 | 75.64 | 78.44 | 79.06 | 80.30 | 84.46 | 89.66 | 79.64 |
| | ± 5.16 | ± 6.34 | ± 3.26 | ± 6.84 | ± 9.38 | ± 5.25 | ± 2.99 | ± 5.27 | ± 3.84 | ± 3.51 | |

Table 1: Table compares validation accuracy of models trained on 10 different tasks defined on the CIFAR100 dataset for single task learning, sequential learning and PaRT. For DL models trained using single task learning, the measured validation accuracy on the learned task and unlearned tasks are provided. Five different DL models were trained for each task and learning method and the validation accuracy of the five models were averaged.

## 4.1 Training on Homogeneous Tasks

In this section, we applied PaRT on 10 homogeneous tasks defined on the CIFAR100 [36] dataset. Each task consists of classifying images from 10 randomly selected classes from the CIFAR100 dataset into 10 categories. ResNet-18 network was used to generate the base network. The total number of modules in each layer in the base network was 12 ($M$=12). Four modules per layer were randomly selected and assigned to each task ($N$=4). Due to random assignment of modules to different tasks, certain modules were shared between different tasks. The module sharing profile for this experiment is shown in Fig. 4 (Left). As an example, the plot shows that around 20 modules were shared among four tasks for PaRT experiments.

In five runs using five different seeds, tasks were defined as explained above and certain modules were assigned to each task randomly. The selected modules were trained on their corresponding task following three different learning methods: 1) single task learning, 2) sequential learning and 3) PaRT. For all learning methods, training was performed until convergence. Fig. 3 (Left) compares the validation accuracy of models trained using these three learning methods. One can see that the average validation accuracy of DL models trained with single task learning method is high for the learned task but very low for the other nine unlearned tasks. Validation accuracy of the DL models trained using PaRT is closer to the achievable performance for each task compared to using sequential learning. It can also be seen that the validation accuracy of the second half of the tasks for sequential learning gradually decreases. We believe the reason is that the base network runs out of unfrozen modules to train on these tasks, and hence insufficient representations are learned for these tasks.

Fig. 3 (Right) compares accuracy of DL models trained on the defined 10 tasks using PaRT with competing algorithms in CL such as RPS [4] and EWC [2], NAS such as pDarts [6], ENAS [9] and MANAS [7], and multi-task learning such as MRN [19]. The experiment setup and hyperparameter values for our method and competing algorithms are provided in Appendix B and the mean and standard deviation of the validation accuracy for each competing algorithm is provided in Appendix C of the Supplementary Materials. It is shown that parallel learning using the simple base network architecture shown in Fig. 2 and with random assignment of modules can outperform many existing complex algorithms from different DL domains. Only a NAS algorithm, MANAS, constantly outperformed PaRT in this experiment. Table 5 shows the mean and standard deviation of the measured validation accuracy for the trained models for PaRT, sequential and single task learning methods.
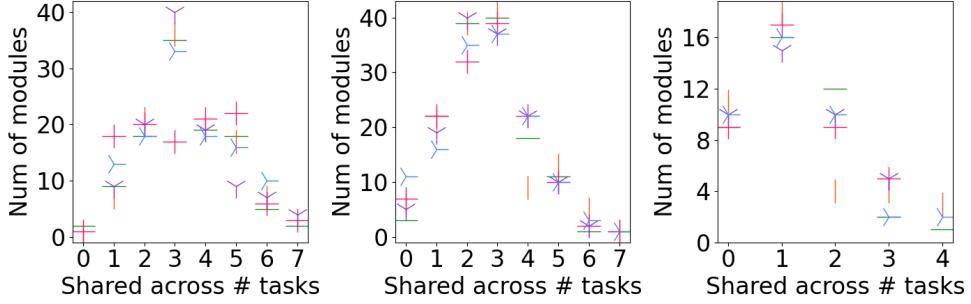


Figure 4: Figure shows the module sharing profile when using PaRT for the CIFAR100 experiments (Left), the CIFAR10 and CIFAR100 experiments (Centre) and FiveTasks experiments (Right).

## 4.2 Heterogenous Tasks

In this section, we applied PaRT on more heterogeneous tasks. Two experiments were conducted. The experiments are presented in increasing order of complexity.

### 4.2.1 CIFAR10 and CIFAR100 Experiment

In this experiment, five tasks were defined on the CIFAR10 [36] dataset, and five tasks were defined on the CIFAR100 [36] dataset. The tasks defined on the CIFAR10 dataset are about classifying images from two randomly selected classes in the CIFAR10 dataset into two categories; the tasks defined on the CIFAR100 dataset are about classifying images from 20 randomly selected classes in the CIFAR100 dataset into 20 categories. ResNet-18 network was used to generate the base network. The total number of modules in each layer in the base network was 15 ($M$=15). Four modules per layer were randomly selected and assigned to each task ($N$=4). The module sharing profile for this experiment is shown in Fig. 4 (Centre).
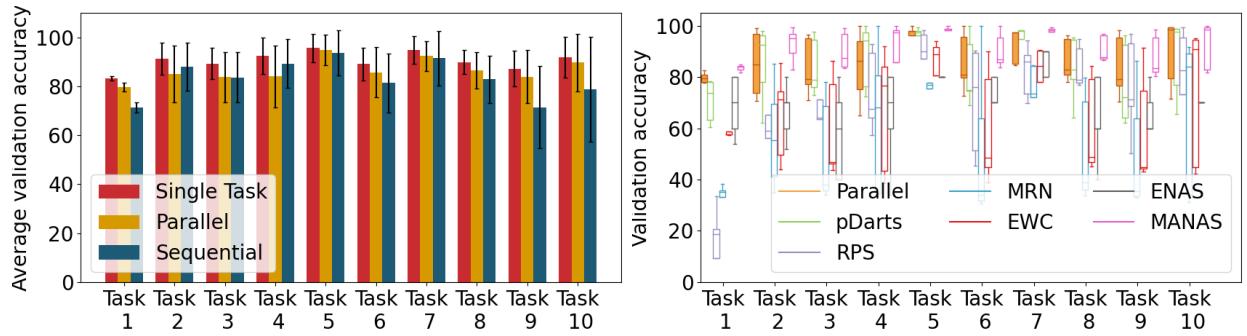


Figure 5: Mean validation accuracy of DL models trained on five tasks defined on the CIFAR10 dataset and five tasks defined on the CIFAR100 [36] dataset. The tasks defined on the CIFAR10 dataset classify two randomly selected classes from the CIFAR10 dataset into two categories. The tasks defined on the CIFAR100 dataset classify 20 randomly selected classes from the CIFAR100 dataset into 20 categories. (Left) compares the performance of PaRT with single task learning and sequential learning. (Right) compares the performance of PaRT with those of the competing algorithms. Five models were trained for each learning method.

6

| Learning Method | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 | Task 6 | Task 7 | Task 8 | Task 9 | Task 10 | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Single Task | 83.29 | 91.30 | 89.24 | 92.45 | 95.64 | 89.07 | 94.90 | 89.82 | 87.22 | 91.85 | 90.48 |
| | ± 0.96 | ± 6.48 | ± 6.45 | ± 7.57 | ± 5.75 | ± 6.62 | ± 5.63 | ± 4.91 | ± 7.39 | ± 8.34 | |
| | Tasks $All\backslash 1$ | Tasks $All\backslash 2$ | Tasks $All\backslash 3$ | Tasks $All\backslash 4$ | Tasks $All\backslash 5$ | Tasks $All\backslash 6$ | Tasks $All\backslash 7$ | Tasks $All\backslash 8$ | Tasks $All\backslash 9$ | Tasks $All\backslash 10$ | |
| | 28.64 | 27.55 | 26.44 | 27.31 | 24.74 | 25.16 | 26.95 | 25.70 | 24.87 | 25.16 | 26.25 |
| | ± 8.11 | ± 8.87 | ± 7.57 | ± 8.27 | ± 8.55 | ± 10.97 | ± 2.19 | ± 7.71 | ± 7.13 | ± 7.35 | |
| Learning Method | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 | Task 6 | Task 7 | Task 8 | Task 9 | Task 10 | |
| Sequential | 71.28 | 88.05 | 83.62 | 89.32 | 93.54 | 81.33 | 91.51 | 82.84 | 71.47 | 78.71 | 83.17 |
| | ± 2.11 | ± 9.81 | ± 10.24 | ± 10.10 | ± 9.21 | ± 11.97 | ± 11.11 | ± 9.67 | ± 16.80 | ± 21.47 | |
| Parallel | 79.75 | 85.00 | 83.77 | 84.10 | 94.88 | 85.76 | 92.33 | 86.61 | 83.94 | 89.63 | 86.58 |
| | ± 1.74 | ± 11.58 | ± 10.22 | ± 12.62 | ± 6.17 | ± 10.37 | ± 6.15 | ± 7.45 | ± 10.89 | ± 11.85 | |

Table 2: Table compares validation accuracy of DL models trained on 10 different tasks defined on the CIFAR10 and CIFAR100 datasets for single task learning, sequential learning and PaRT. For DL models trained using single task learning, the measured validation accuracy on the learned task and unlearned tasks are provided. Five different models were trained for each task and learning method and the validation accuracy of the five models were averaged.

Five DL models were trained for each learning method. Fig. 5 (Left) shows the results for this experiment using PaRT, single task learning, and sequential learning. Table 2 shows the mean and standard deviation of the measured validation accuracy for the trained models for PaRT, sequential and single task learning methods. Similar to the previous experiment, the mean validation accuracy for DL models trained using single task learning on the learned tasks is higher that PaRT and sequential learning. However, their performance on the other nine tasks is found to be low. In addition, even though we increased the total number of modules in each layer ($M$) to 15, the DL models trained using the Sequential Learning method still do not have as good performance on the defined 10 tasks as the DL models trained using PaRT. Similar behaviour is seen when comparing the validation accuracy of the DL models trained on the 10 tasks using PaRT with competing algorithms (Fig. 5 (Right)). The experiment setup and hyperparameter values for our method and each competing algorithm is provided in Appendix B and the mean and standard deviation of the validation accuracy for each competing algorithm is provided in Appendix C of the Supplementary Materials.

### 4.2.2 FiveTasks Experiment

In this experiment, we ran experiments on more heterogeneous classification tasks. We borrowed this experimental design from the work by Hung et al. [5]. Five classification tasks were defined on five different image datasets. Table 3 shows a summary of each dataset. To handle the heterogeneity in the input data distributions, we modified the base network by assigning a specific batch normalization module to each task in all the batch normalization layers in the base network. Fig. 6 shows how the batch normalization layers were modified to have multiple batch normalization modules. Once a new task is added for training, a new batch normalization module will be added to each batch normalization layer in the base network.

In this experiment, ResNet-50 network was used to generate the base network. A detailed representation of the base network architecture generated using ResNet-50 is shown in Appendix A of the supplementary Materials. The total number of modules in each layer was set to 8 ($M$=8). The base network was initialized with parameters of a ResNet-50 model trained on the ImageNet dataset following Hung et al. [5]. Two modules in each layer were randomly selected and assigned to each task ($N$=2). The module sharing profile for this experiment is shown in Fig. 4 (Right).

| Task | Dataset | #Train | #Eval | #Classes |
|---|---|---|---|---|
| Task 1 | CUBs [33] | 5,994 | 5,794 | 200 |
| Task 2 | Stanford Cars [34] | 8,144 | 8,041 | 196 |
| Task 3 | Flowers [35] | 2,040 | 6,149 | 102 |
| Task 4 | WikiArt [32] | 42,129 | 10,628 | 195 |
| Task 5 | Sketch [31] | 16,000 | 4,000 | 250 |

Table 3: Table shows brief information about the five datasets that were deployed to define five heterogeneous classification tasks in section 4.2.2.
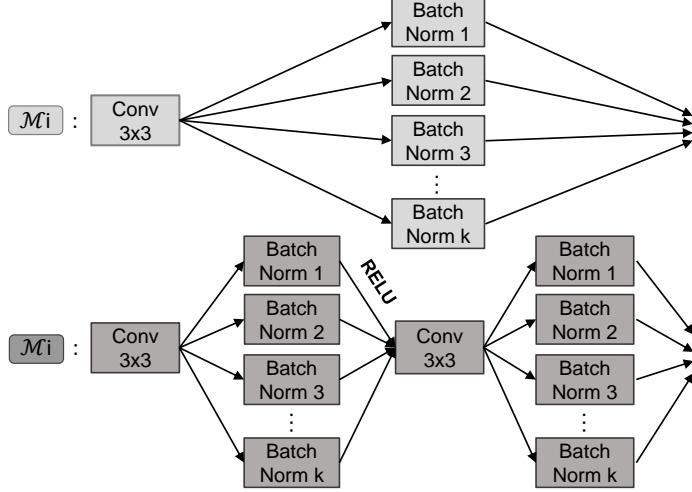
Figure 6: To handle the heterogeneity of the tasks defined in our FiveTasks experiments, each task was assigned a separate batch normalization instance in each batch normalization layer in the base network. In case a new task is introduced, a new batch normalization instance will be added to each batch normalization layer in the base network for that task.

| Method | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 | Avg |
|---|---|---|---|---|---|---|
| | 82.75 ± 0.40 | 92.08 ± 0.26 | 95.32 ± 0.42 | 74.01 ± 0.88 | 79.35 ± 0.23 | 84.70 |
| Single | Tasks $All\backslash 1$ | Tasks $All\backslash 2$ | Tasks $All\backslash 3$ | Tasks $All\backslash 4$ | Tasks $All\backslash 5$ | |
| | 0.57 ± 0.36 | 0.60 ± 0.43 | 0.43 ± 0.16 | 0.57 ± 0.50 | 0.49 ± 0.18 | 0.53 |
| | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 | |
| Parallel | 80.63 ± 0.27 | 92.08 ± 0.29 | 94.38 ± 0.15 | 74.44 ± 1.02 | 78.32 ± 0.65 | 83.97 |
| CPG [5] | 83.59 | 92.80 | 96.62 | 77.15 | 80.33 | 86.10 |

Table 4: Table compares validation accuracy of DL models trained for the FiveTasks experiment using single task learning and PaRT. For DL models trained using single task learning, the measured validation accuracy on the learned task and unlearned tasks are provided. The results are compared with those of the proposed algorithm by Hung et al.[5] named CPG. Five different models were trained for each task for single task learning and PaRT and the validation accuracy of all the five models were averaged.

Since PaRT outperformed sequential learning in the previous experiments, we compared PaRT solely with single task learning in this section. Fig. 7 shows the average validation accuracy for each task over five DL models trained with different seeds using single task learning and PaRT. It is shown that the average validation accuracy of DL models trained using PaRT is only slightly lower than that of DL models trained with single task learning. Table 4 shows the mean and standard deviation of validation accuracy for each task for models trained using single task learning and PaRT. For Task 4, the mean validation accuracy of PaRT is even higher than that of DL models trained using single task learning. We have also shown the validation accuracy achieved by Hung et al. [5]. CPG [5] which is a more complex CL algorithm outperforms PaRT with an average of 2.13% in validation accuracy. Considering the simplicity of our base network and the current strategy for module assignment, we believe a better base network and a more efficient module assignment strategy may train models that outperform the performance of models trained using CPG [5].

# 5 Discussion

In sections 4.1 and 4.2.1, we showed that DL models trained using PaRT can utilize network resources more efficiently compared to models trained with the sequential learning algorithm. We also showed that DL models trained with PaRT on more heterogeneous tasks in section 4.2.2 can perform as well as models trained with single task learning. In order to explain how PaRT enables these improvements in network resource efficiency and task performance compared to sequential learning, we compared the task-related and shared representations learned by the DL modules for each task.
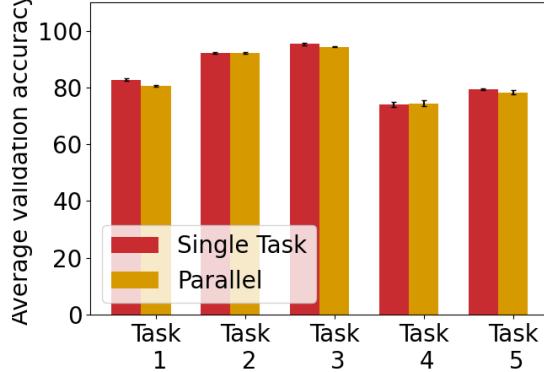
Figure 7: Validation accuracy of tasks in the FiveTasks experiments using single task learning and PaRT. The validation accuracy is measured on five DL models trained for each learning method.

## 5.1 Learned Representations

We designed a separate experiment to study the effect of module-sharing on the learned representations by the base network using PaRT. For this experiment, we considered a base network using the ResNet-50 architecture with four modules in each layer ($M$=4). Similar to the FiveTasks experiment, we initialized the network with parameters of a ResNet-50 model trained on the ImageNet dataset. The base network was trained on two tasks only. Two modules in each layer were selected and assigned to each task ($N$=2), however, selection of the modules was not done randomly. Ten different setups were considered with each setup defining a specific way of sharing modules between the two tasks. For instance, in 'no layer' setup, no modules were shared between the two tasks. This setup is equivalent to single task learning of each of the two tasks. In 'layer 1' setup, two modules in layer 1 were shared between the two tasks. and in 'layer 12345' setup, two modules at each layer were shared between the two tasks. In experiment (A), the base network was trained using PaRT on two tasks which were classification of Stanford Cars and CUBs datasets. In experiment (B), the base network was trained using PaRT on two tasks which were classification of Flowers and Sketches datasets. Three DL models with different random seeds were trained using each setup and the similarity of their learned representations for the two tasks were compared with each other in both the experiments.

We deployed the representation similarity metric proposed by Kornblith et al. [38] called CKA. For this comparison, 1000 samples from the validation dataset of each task with equal number of samples per class were passed through each DL model and the similarity of the representations for the two tasks at each layer were compared with DL models trained with ten different setups. Fig. 8 shows the average similarity of learned representations by the modules for the two tasks at each layer for different module-sharing setups for experiments (A) and (B). Multiple important points can be observed in these plots including:

- Learned representations at the middle layers are more similar across the two tasks compared to learned representations at the first and last layers.

- Compared to single task learning, sharing modules at layer 3 results in more similar representations across the two tasks at the surrounding layers but not at the first and last layers.

- Compared to single task learning, sharing modules at layers 1, 2, and 3 results in more similar representations across the two task at layer 4 but not at the last layer.

- Compared to other module-sharing setups, sharing modules at all layers leads to less similar representations at different layers across the two tasks.

We have also shown the selected path for each task and the average pair-wise CKA similarity between different modules at the same layer for DL models trained using 'no layer', 'layer 1', 'layer 3', 'layer 123', and 'layer 5' in Figures 9, 10, 11 and 12, respectively. The average CKA similarity metric is measured across 3 DL models trained using each setup with different random seeds. We show the CKA comparison results using an RBF kernel with $\sigma$=0.5 bandwidth which controls the extent to which similarity of small distances is emphasized over large distances. A value close to one, shown by a bight yellow color, means more similar representations were learned for the two tasks by the corresponding module. The shared modules between the two tasks in DL models trained using PaRT are highlighted with a green bounding box. This experiment provides a high level of transparency over what features a DL model learns at different layers. We believe the low CKA similarity seen across different tasks at the first and final layers
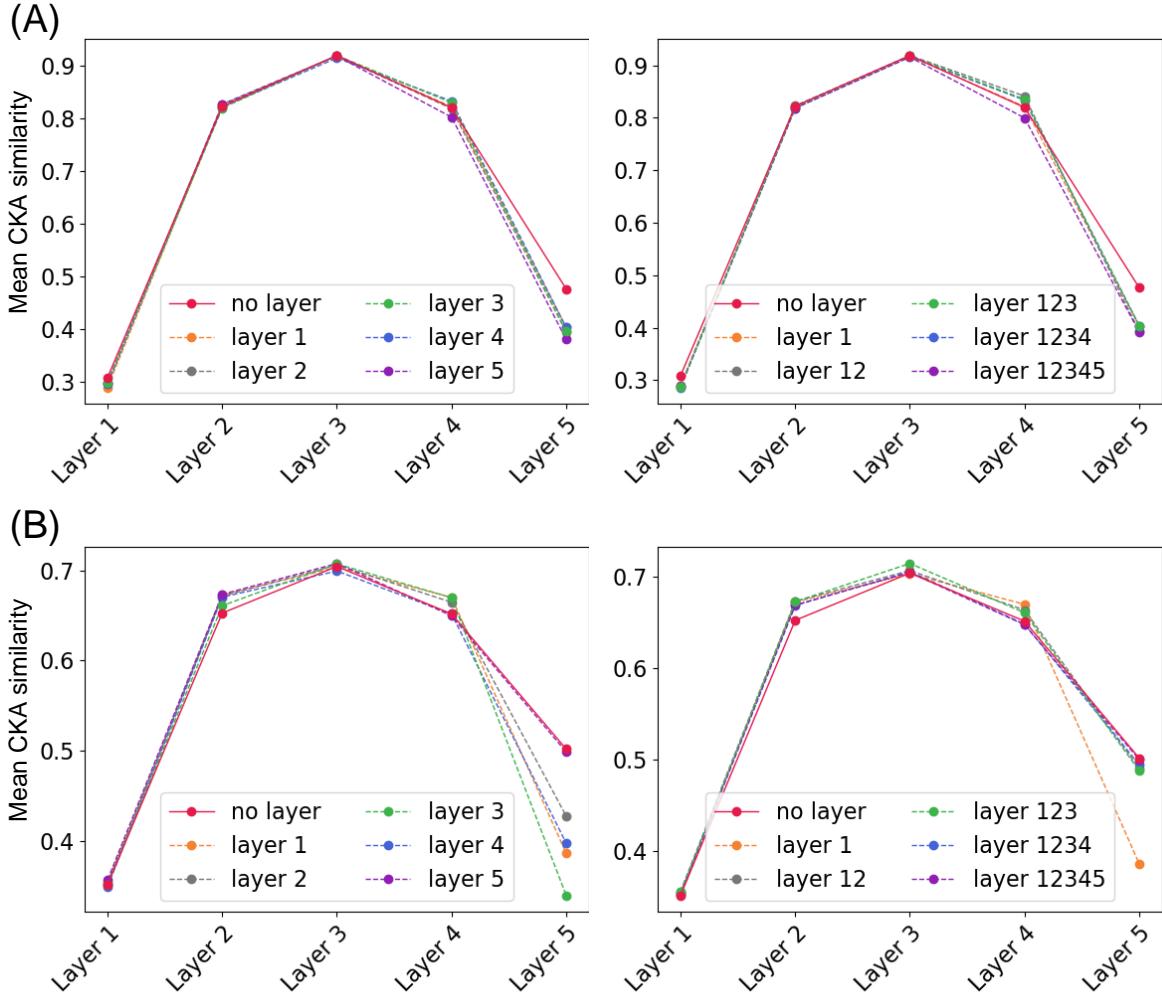
9

Figure 8: (A) Figure compares learned representations at each layer by DL models trained with PaRT on classification of the Stanford Cars and CUBs datasets with different amounts of module-sharing between the two tasks. The label 'no layer' means a DL model was trained on the two tasks while no modules were shared between the two tasks which is equivalent to single task learning. Label 'layer 1' means a DL model was trained on the two tasks while modules in layer 1 were shared between the two tasks and label 'layer 12345' means a DL model was trained on the two tasks while all modules at all layers were shared between the two tasks. (B) Figure compares learned representations by DL models trained with PaRT on classification of the Flowers and Sketches datasets with different amounts of module-sharing between the two tasks. The labels convey the same meaning as in (A).

implies that the basic features learned by DL models for different tasks and the final decision making is quite different across tasks whereas the abstract features extracted in the middle layers of the network are more similar across different tasks.

## 5.2 Limitations

In a normal training scenario, a single ResNet-18 or ResNet-50 model usually suffices to learn one task. However, PaRT currently assigns $N > 1$ modules in each layer to a task. In an ideal scenario, PaRT should be able to learn multiple tasks using almost the same amount of network resources that would be required while training a DL model on each task. Reducing utilization of network resources is a part of our future plan for extending PaRT. PaRT should also be tested on more diverse tasks such as classification and segmentation tasks. Moreover, while the random selection of modules in the paths for different tasks already gives decent results, this process of selecting modules can be optimized by incorporating Reinforcement Learning which takes the average accuracy or loss as the reward for the agent to learn.
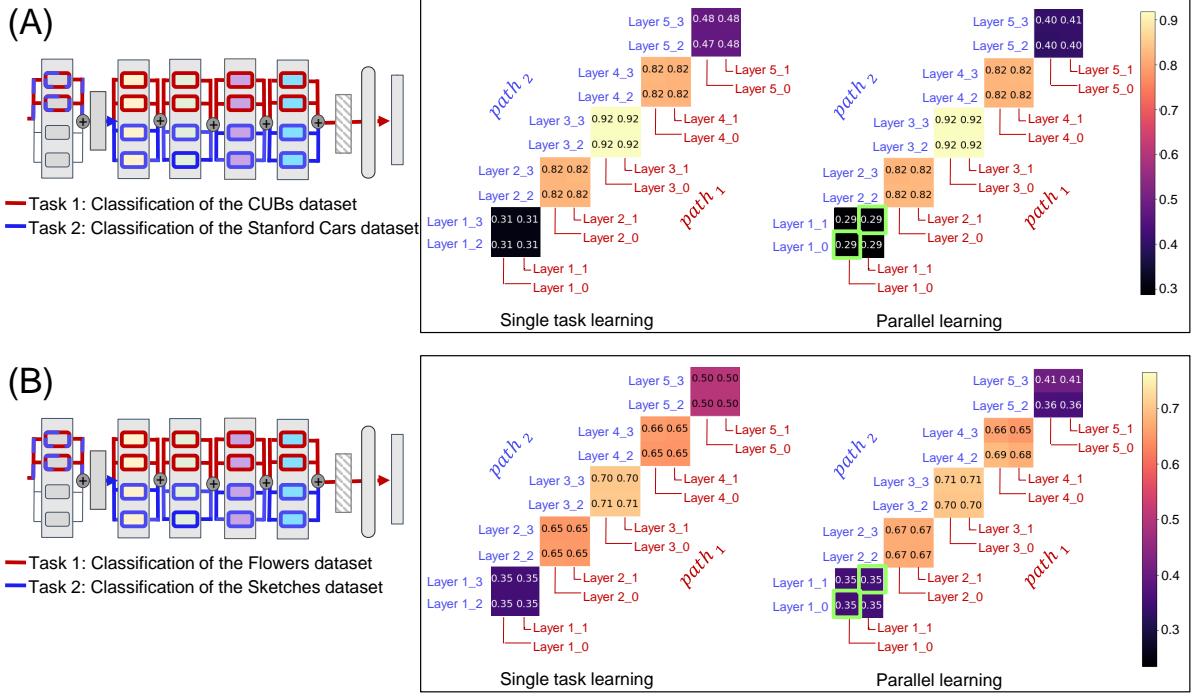
Figure 9: Figure compares learned representations by a DL model trained with PaRT while no modules were shared between the two tasks,which is equivalent to single task learning, with learned representations by a DL model trained with PaRT on those two tasks while sharing two modules at the first layer of the network. (A) provides this comparison for the Stanford Cars and CUBs datasets. (B) shows this comparison for the Flowers and WikiArt datasets. CKA with RBF kernel ($\sigma$=0.5) [38] was used as the representation similarity metric. Higher values show more similar representations. The tasks and their selected modules are shown at the left. The shared modules are also shown by a green bounding box on the heatmaps.

## 6    Conclusion

In this paper, we take a parallel learning approach to train DL models on multiple tasks simultaneously. Parallel learning achieves a better performance than many sequential continual learning methods and state-of-the-art methods that are proposed in other domains. This is the first effort that harnesses shared representations to improve individual task performance and increases efficient use of network resources. It also increases transparency in the model, therefore allowing us to draw explanations out of the model's predictions. Parallel learning has widespread applications in the real-world. One area where it can be harnessed is the healthcare industry. Our model enables a single DL model to have the capability of identifying multiple diseases such as lung cancer, and heart failure, hence largely lowering the cost of incorporating DL in hospitals.

## References

[1] Lesort, Timothée, Continual Learning: Tackling Catastrophic Forgetting in Deep Neural Networks with Replay Processes, arXiv preprint arXiv:2007.00487, 2020 3

[2] Kirkpatrick, James and Pascanu, Razvan and Rabinowitz, Neil and Veness, Joel and Desjardins, Guillaume and Rusu, Andrei A and Milan, Kieran and Quan, John and Ramalho, Tiago and Grabska-Barwinska, Agnieszka and others, Overcoming catastrophic forgetting in neural networks, Proceedings of the national academy of sciences, 114, 13, 3521–3526, 2017 3, 6, 17, 18, 19, 20

[3] Shin, Hanul and Lee, Jung Kwon and Kim, Jaehong and Kim, Jiwon, Continual learning with deep generative replay, arXiv preprint arXiv:1705.08690, 2017 3

[4] Rajasegaran, Jathushan and Hayat, Munawar and Khan, Salman H and Khan, Fahad Shahbaz and Shao, Ling, Random path selection for continual learning, Advances in Neural Information Processing Systems 32 (NeurIPS 2019) 3, 6, 17, 18, 19, 20
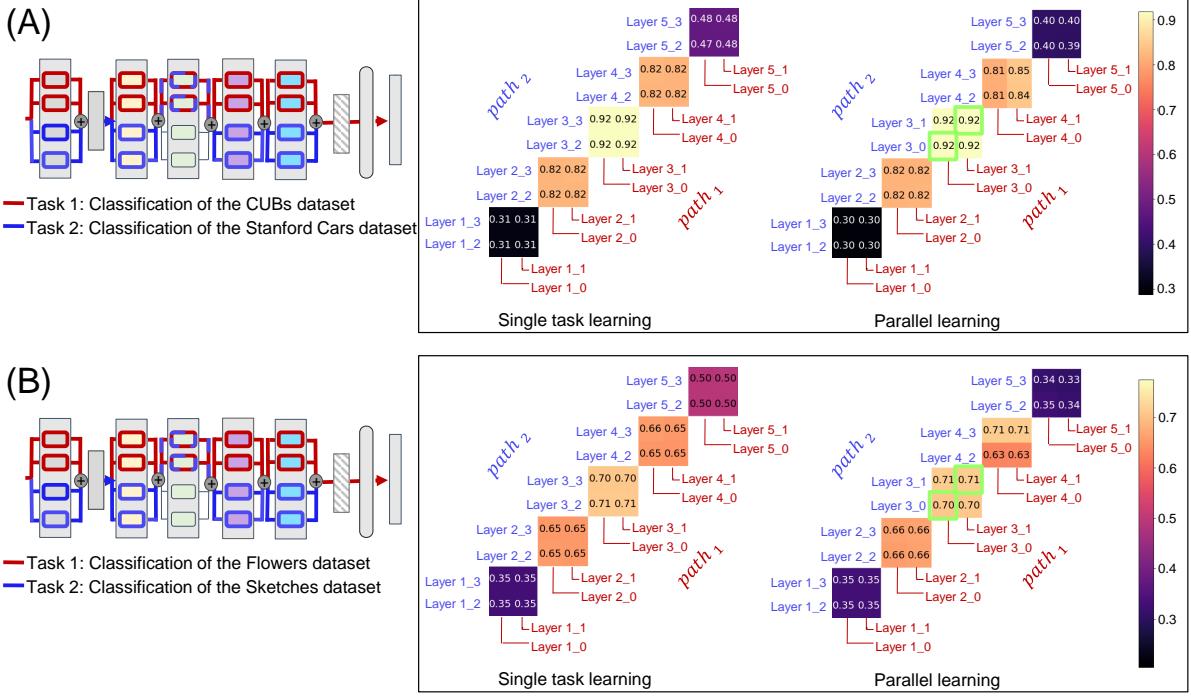
Figure 10: Figure compares learned representations by a DL model trained with PaRT while no modules were shared between the two tasks, which is equivalent to single task learning, with learned representations by a DL model trained with PaRT on those two tasks while sharing two modules at the third layer of the network. (A) provides this comparison for the Stanford Cars and CUBs datasets. (B) shows this comparison for the Flowers and WikiArt datasets. CKA with RBF kernel ($\sigma$=0.5) [38] was used as the representation similarity metric. Higher values show more similar representations. The tasks and their selected modules are shown at the left. The shared modules are also shown by a green bounding box on the heatmaps.

[5] Hung, Steven CY and Tu, Cheng-Hao and Wu, Cheng-En and Chen, Chien-Hung and Chan, Yi-Ming and Chen, Chu-Song, Compacting, picking and growing for unforgetting continual learning, arXiv preprint arXiv:1910.06562, 2019 3, 7, 8, 18

[6] Chen, Xin and Xie, Lingxi and Wu, Jun and Tian, Qi, Progressive differentiable architecture search: Bridging the depth gap between search and evaluation, Proceedings of the IEEE/CVF International Conference on Computer Vision, 1294–1303, 2019 3, 6, 17, 18, 19, 20

[7] Carlucci, Fabio Maria and Esperança, Pedro M and Singh, Marco and Gabillon, Victor and Yang, Antoine and Xu, Hang and Chen, Zewei and Wang, Jun, MANAS: multi-agent neural architecture search, arXiv preprint arXiv:1909.01051, 2019 3, 6, 17, 18, 19, 20

[8] Isensee, Fabian and Petersen, Jens and Klein, Andre and Zimmerer, David and Jaeger, Paul F and Kohl, Simon and Wasserthal, Jakob and Koehler, Gregor and Norajitra, Tobias and Wirkert, Sebastian and others, nnu-net: Self-adapting framework for u-net-based medical image segmentation, arXiv preprint arXiv:1809.10486, 2018 3

[9] Pham, Hieu and Guan, Melody and Zoph, Barret and Le, Quoc and Dean, Jeff, Efficient neural architecture search via parameters sharing, International Conference on Machine Learning, 4095–4104, 2018 3, 6, 17, 18, 19, 20

[10] Zoph, Barret and Le, Quoc V, Neural architecture search with reinforcement learning, , arXiv preprint arXiv:1611.01578, 2016 3

[11] Zoph, Barret and Vasudevan, Vijay and Shlens, Jonathon and Le, Quoc V, Learning transferable architectures for scalable image recognition, Proceedings of the IEEE conference on computer vision and pattern recognition, 8697–8710, 2018 3

[12] Kirsch, Louis and Kunze, Julius and Barber, David, Modular networks: Learning to decompose neural computation,, arXiv preprint arXiv:1811.05249, 2018 3

[13] Liu, Hanxiao and Simonyan, Karen and Yang, Yiming, Darts: Differentiable architecture search, arXiv preprint arXiv:1806.09055, 2018 3

[14] Gou, Yuanbiao and Li, Boyun and Liu, Zitao and Yang, Songfan and Peng, Xi, CLEARER: Multi-scale neural architecture search for image restoration, Advances in Neural Information Processing Systems, 33, 2020 3
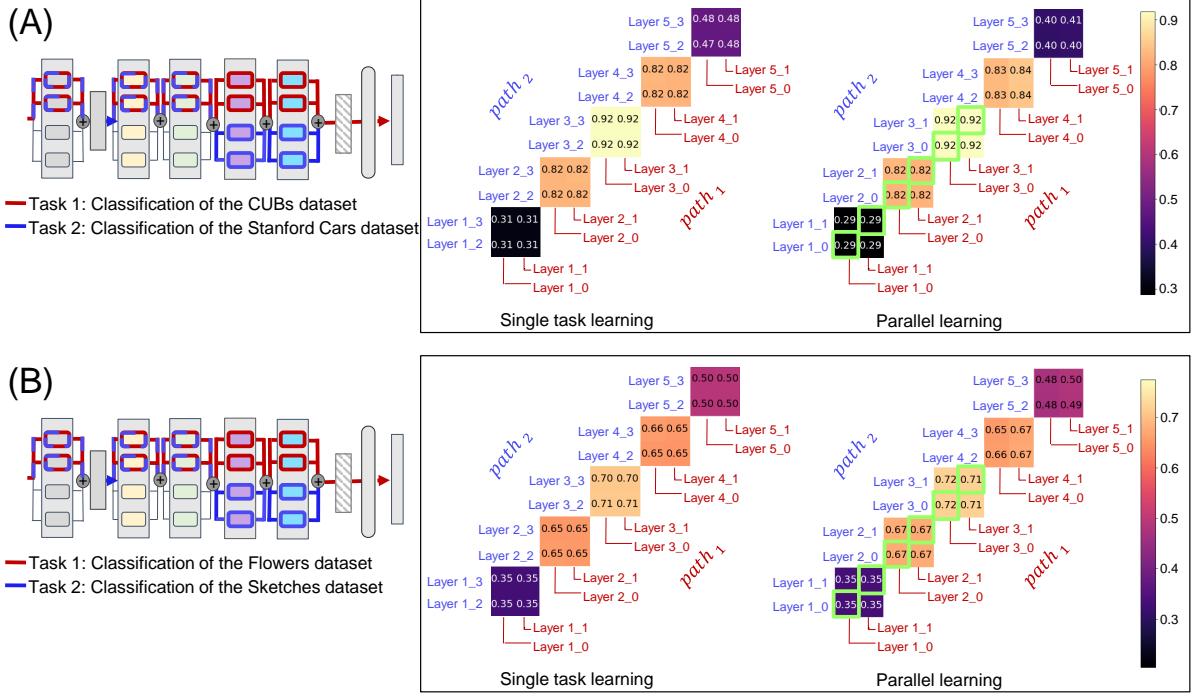
Figure 11: Figure compares learned representations by a DL model trained with PaRT while no modules were shared between the two tasks, which is equivalent to single task learning, with learned representations by a DL model trained with PaRT on those two tasks while sharing six modules at the first, second and third layers of the network. (A) provides this comparison for the Stanford Cars and CUBs datasets. (B) shows this comparison for the Flowers and WikiArt datasets. CKA with RBF kernel ($\sigma$=0.5) [38] was used as the representation similarity metric. Higher values show more similar representations. The tasks and their selected modules are shown at the left. The shared modules are also shown by a green bounding box on the heatmaps.

[15] Negrinho, Renato and Gordon, Geoff, Deeparchitect: Automatically designing and training deep architectures, arXiv preprint arXiv:1704.08792, 2017 3

[16] Liu, Chenxi and Zoph, Barret and Neumann, Maxim and Shlens, Jonathon and Hua, Wei and Li, Li-Jia and Fei-Fei, Li and Yuille, Alan and Huang, Jonathan and Murphy, Kevin, Progressive neural architecture search, Proceedings of the European conference on computer vision (ECCV), 19–34, 2018 3

[17] Yao, Lewei and Xu, Hang and Zhang, Wei and Liang, Xiaodan and Li, Zhenguo, SM-NAS: structural-to-modular neural architecture search for object detection, Proceedings of the AAAI Conference on Artificial Intelligence, 34, 07, 12661–12668, 2020 3

[18] Real, Esteban and Moore, Sherry and Selle, Andrew and Saxena, Saurabh and Suematsu, Yutaka Leon and Tan, Jie and Le, Quoc V and Kurakin, Alexey, Large-scale evolution of image classifiers, International Conference on Machine Learning, 2902–2911, 2017 3

[19] Long, Mingsheng and Cao, Zhangjie and Wang Jianmin and Philip S., Yu, Learning multiple tasks with multilinear relationship networks, arXiv preprint arXiv:1506.02117, 2, 1, 2015 2, 6, 16, 18, 19, 20

[20] Zhou, Yue and Chen, Houjin and Li, Yanfeng and Liu, Qin and Xu, Xuanang and Wang, Shu and Yap, Pew-Thian and Shen, Dinggang, Multi-task learning for segmentation and classification of tumors in 3D automated breast ultrasound images, Medical Image Analysis, 70, 101918, 2021 2

[21] Ruder, Sebastian, An overview of multi-task learning in deep neural networks, arXiv preprint arXiv:1706.05098, 2017 2

[22] Vandenhende, Simon and Georgoulis, Stamatios and Van Gool, Luc, Mti-net: Multi-scale task interaction networks for multi-task learning, European Conference on Computer Vision, 527–543, 2020 2

[23] Misra, Ishan and Shrivastava, Abhinav and Gupta, Abhinav and Hebert, Martial, Cross-stitch networks for multi-task learning, Proceedings of the IEEE conference on computer vision and pattern recognition, 3994–4003, 2016 2

[24] Tan, Chuanqi and Sun, Fuchun and Kong, Tao and Zhang, Wenchang and Yang, Chao and Liu, Chunfang, A survey on deep transfer learning, International conference on artificial neural networks, 270–279, 2018 2
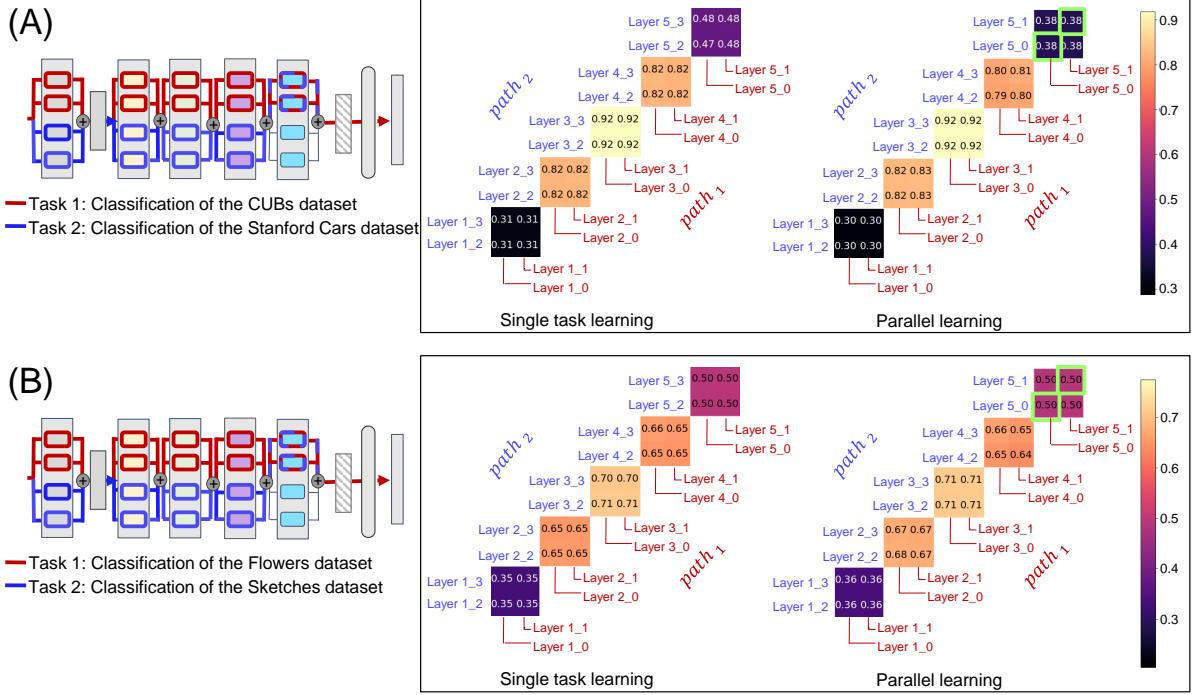
Figure 12: Figure compares learned representations by a DL model trained with PaRT while no modules were shared between the two tasks, which is equivalent to single task learning, with learned representations by a DL model trained with PaRT on those two tasks while sharing two modules at the fifth layer of the network. (A) provides this comparison for the Stanford Cars and CUBs datasets. (B) shows this comparison for the Flowers and WikiArt datasets. CKA with RBF kernel ($\sigma$=0.5) [38] was used as the representation similarity metric. Higher values show more similar representations. The tasks and their selected modules are shown at the left. The shared modules are also shown by a green bounding box on the heatmaps.

[25] Zhuang, Fuzhen and Qi, Zhiyuan and Duan, Keyu and Xi, Dongbo and Zhu, Yongchun and Zhu, Hengshu and Xiong, Hui and He, Qing, A comprehensive survey on transfer learning, Proceedings of the IEEE, 109, 1, 43–76, 2020 2

[26] Liu, Lang and Fard, Mahdi Milani and Zhao, Sen, Distribution Embedding Network for Meta-Learning with Variable-Length Input, OpenReview, 2020 2

[27] Hospedales, Timothy and Antoniou, Antreas and Micaelli, Paul and Storkey, Amos, Meta-learning in neural networks: A survey, arXiv preprint arXiv:2004.05439, 2020 2

[28] Rusu, Andrei A and Rao, Dushyant and Sygnowski, Jakub and Vinyals, Oriol and Pascanu, Razvan and Osindero, Simon and Hadsell, Raia, Meta-learning with latent embedding optimization, arXiv preprint arXiv:1807.05960, 2018 2

[29] Shin, Hanul and Lee, Jung Kwon and Kim, Jaehong and Kim, Jiwon, Continual learning with deep generative replay, arXiv preprint arXiv:1705.08690, 2017 3

[30] van de Ven, Gido M and Siegelmann, Hava T and Tolias, Andreas S, Brain-inspired replay for continual learning with artificial neural networks, Nature communications, 11, 1, 1–14, 2020 3

[31] Eitz, Mathias and Hays, James and Alexa, Marc, How Do Humans Sketch Objects?, ACM Trans. Graph. (Proc. SIGGRAPH), 2012, 31, 4, 44:1–44:10 7

[32] Saleh, Babak and Elgammal, Ahmed, Large-scale classification of fine-art paintings: Learning the right metric on the right feature, arXiv preprint arXiv:1505.00855, 2015 7

[33] Wah, Catherine and Branson, Steve and Welinder, Peter and Perona, Pietro and Belongie, Serge, The caltech-ucsd birds-200-2011 dataset, 2011 7

[34] Krause, Jonathan and Stark, Michael and Deng, Jia and Fei-Fei, Li, 3d object representations for fine-grained categorization, Proceedings of the IEEE international conference on computer vision workshops, 554–561, 2013 7

[35] Nilsback, Maria-Elena and Zisserman, Andrew, Automated flower classification over a large number of classes, 2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing, 722–729, 2008 7

[36] Krizhevsky, Alex and Hinton, Geoffrey and others, Learning multiple layers of features from tiny images, Citeseer, 2009 5, 6

[37] Iwata, Tomoharu and Kumagai, Atsutoshi, Meta-learning from Tasks with Heterogeneous Attribute Spaces, Advances in Neural Information Processing Systems, 33, 2020 2

[38] Kornblith, Simon and Norouzi, Mohammad and Lee, Honglak and Hinton, Geoffrey, Similarity of neural network representations revisited, International Conference on Machine Learning, 3519–3529, 2019 4, 9, 11, 12, 13, 14

[39] Shin, Hanul and Lee, Jung Kwon and Kim, Jaehong and Kim, Jiwon, Continual learning with deep generative replay, arXiv preprint arXiv:1705.08690, 2017 3

[40] Van de Ven, Gido M and Tolias, Andreas S, Generative replay with feedback connections as a general strategy for continual learning, arXiv preprint arXiv:1809.10635, 2018 3

[41] Lesort, Timothée and Caselles-Dupré, Hugo and Garcia-Ortiz, Michael and Stoian, Andrei and Filliat, David, Generative models from the perspective of continual learning, 2019 International Joint Conference on Neural Networks (IJCNN), 1–8, 2019 3

[42] Aljundi, Rahaf and Caccia, Lucas and Belilovsky, Eugene and Caccia, Massimo and Lin, Min and Charlin, Laurent and Tuytelaars, Tinne, Online continual learning with maximally interfered retrieval, arXiv preprint arXiv:1908.04742, 2019 3

[43] Zhang, Yu and Yang, Qiang, A survey on multi-task learning, IEEE Transactions on Knowledge and Data Engineering, 2021 1

[44] Zhang, Yu and Yang, Qiang, An overview of multi-task learning, National Science Review, 5, 1, 30–43, 2018 1

[45] Thung, Kim-Han and Wee, Chong-Yaw, A brief review on multi-task learning, Multimedia Tools and Applications, 77, 22, 29705–29725, 2018 1

[46] Parisi, German I and Kemker, Ronald and Part, Jose L and Kanan, Christopher and Wermter, Stefan, Continual lifelong learning with neural networks: A review, Neural Networks, 113, 54–71, 2019 1

[47] Hadsell, Raia and Rao, Dushyant and Rusu, Andrei A and Pascanu, Razvan, Embracing change: Continual learning in deep neural networks, Trends in cognitive sciences, 2020 1

[48] Rebuffi, Sylvestre-Alvise and Kolesnikov, Alexander and Sperl, Georg and Lampert, Christoph H, icarl: Incremental classifier and representation learning, Proceedings of the IEEE conference on Computer Vision and Pattern Recognition 2001–2010, 2017 3

[49] Lopez-Paz, David and Ranzato, Marc'Aurelio, Gradient episodic memory for continual learning, Advances in neural information processing systems, 30, 6467–6476, 2017 3

[50] Chaudhry, Arslan and Dokania, Puneet K and Ajanthan, Thalaiyasingam and Torr, Philip HS, Riemannian walk for incremental learning: Understanding forgetting and intransigence, Proceedings of the European Conference on Computer Vision (ECCV), 532–547, 2018 17

# 7    APPENDIX A: Our Base Network Generated Using The ResNet-50 Network Architecture

Fig 13 shows the architecture of the base network generated using the ResNet-50 network architecture. This base network architecture was deployed in the FiveTasks experiments presented in section 4.2.2 in the original paper. Similar to the base network architecture generated using ResNet-18, the base network consists of multiple layers, with a total of $M$ modules in each layer. The outputs of modules in the same layer are summed together, and the sum is fed to all the modules in the next layer. One row in the base network is equivalent to the ResNet-50 network architecture. The architecture of each module is shown at the bottom of the figure.

# 8    APPENDIX B: Experiment Setups and Hyperparameter Values

In this section, we provide the experiment setup for each learning method for all three experiments conducted in the original paper and the hyperparameters we eventually decided on after fine-tuning them.

## 8.1    CIFAR100 Experiments

This section provides the experiment setup for each compared learning method in section 4.1 in the original paper.

### 8.1.1    Image Pre-processing Pipeline

The same image pre-processing pipeline was applied to the images before running the learning methods mentioned below. We normalized images in the CIFAR100 training dataset using mean values of $\{0.507, 0.486, 0.44\}$ and standard deviation values of $\{0.267, 0.256, 0.276\}$ for the $RGB$ channels, respectively. We also applied image augmentation methods including $RandomPerspective$ with probability of 0.75 and distortion scale of 0.2,
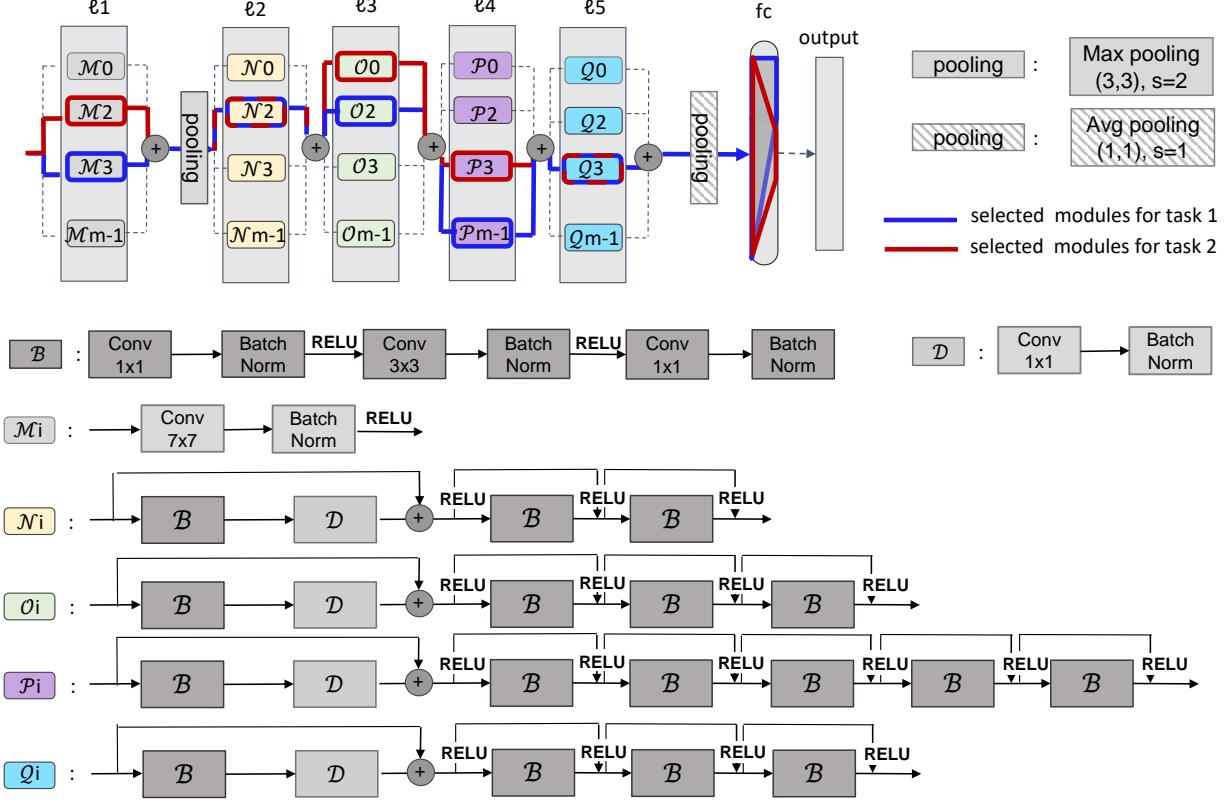
Figure 13: Figure shows the architecture of the base network that was deployed for the FiveTasks experiments in section 4.2.2 in the original paper. One row in the base network is equivalent to the ResNet-50 network architecture.

$RandomHorizontalFlip$ and $RandomCrop$ with output size of 32 pixels and 4 pixels padding. The preprocessing pipeline for the images in the CIFAR100 test dataset included only image normalization using the same mean and standard deviation values.

### 8.1.2  Single Task Learning

In this experiment, we used the base network generated using the ResNet-18 network architecture. At each layer, $M$=12 modules were specified. For each task, $N$=4 modules were randomly selected and were assigned to the task. $ADAM$ optimizer with $CrossEntropy$ loss function was deployed. The initial learning rate for the optimizer was set to 0.001. The learning rate was halved at epochs 20, 30, and 40. Models were trained until convergence.

### 8.1.3  Sequential Learning

For this experiment, the same hyperparameter values and experiment setups were followed as for the single task learning experiments. During each run, the classes assigned to each task and the order of tasks to learn were set randomly. The model was trained on each task sequentially. The selected modules for each task were frozen after the training for the task ended.

### 8.1.4  Parallel Learning

The hyperparameter values were set to the same values used in the single task learning experiments and the $batch$-$set$ size was set to 10.

### 8.1.5  MRN [19]

We used the official code provided by the authors. The original code uses VGG-16 as the underlying network architecture, and the last fully connected layer in VGG-16 is used as the head for each task. We replaced this network with

16

ResNet-18, and the last fully connected layer in ResNet-18 as the head for each task. The initial learning rate was set to 0.05, the gamma value was set to 0.2 with 0.75 power and 3000 stepsize.

### 8.1.6  pDarts [6]

We used the official code provided by the authors. The initial learning rate is set to 0.025 with momentum of 0.9 and weight decay of $3 \times 10^{-4}$. The number of search epochs is set to 5 and training epochs to 60. We specified 16 initial search channels and 36 initial train channels. The total number of search layers was set to 5, and the total number of train layers was set to 20. We did not use cutout. Drop path probability and training portion were set to 0.3. The initial $DropOut$ probability was set to 0.0, 0.3, 0.6 for stage 1, 2, and 3, respectively. An $ADAM$ optimizer with a learning rate of 0.0006 and weight decay of 0.001 and momentum of (0.5,0.999) were used for architecture encoding. Finally, gradient clipping was set to 5.

### 8.1.7  ENAS [9]

For ENAS experiments, we used the Pytorch implementation of this algorithm on GitHub: ENAS-Pytorch. The parameters were trained with a momentum value of 0.9 where the learning rate followed the cosine schedule with $l_{max} = 0.05$, $l_{min} = 0.0005$, $T_0 = 10$, $T_{mul} = 2$, and a weight decay of $10^{-4}$. The policy parameters were trained using $ADAM$ optimizer with a learning rate of 0.0035, where a TanH constant of 1.10, TanH reduce rate of 2.5, and a temperature of 5.0 were applied to the controller's logits. The controller entropy was added to the reward, weighted by 0.0001. The LSTM size was set to 64, with 1 layer, 0 keep probability.

### 8.1.8  MANAS [7]

We received the MANAS code from the authors. For MANAS experiments we set the number of layers to 20, reward baseline to 5, number of iterations to 390, the starting temperature to 1000, the ending temperature to 200, and the number of initial channels to 16.

### 8.1.9  EWC [2]

We used the implementation of EWC [2] that is applied to the CIFAR100 dataset by Chaudhry et al. [50]. We trained the ResNet-18 network architecture with multiple heads similar to our approach and used the same hyperparameters specified in the paper for the CIFAR100 dataset that is optimizer=$ADAM$, lambda=1000, arch=$ResNet\text{-}B$ and method=$RWalk$.

### 8.1.10  RPS [4]

For RPS experiments, we used our base network generated using ResNet-18 network architecture. It is important to mention that our base network is different from the network deployed in this paper. In this paper, the authors have added one extra module at layers 4, 6, and 8 in the network. These extra modules are never frozen and the output of the modules for all tasks pass through these extra modules. We set he number of modules at each layer to $M$=12 and the number of randomly selected modules for each task to $N$=4. We used the same hyperparameter values used in this paper, which include an $ADAM$ optimizer with $t_e = 2$ and an initial learning rate of $10^{-3}$ that is halved at epochs 20, 40 and 60. The scaling factor and the path selection frequency ($J$) values for the controller are set to 2.5 and 1, respectively.

## 8.2  CIFAR10 and CIFAR100 Experiments

This section provides the experimental setup for each compared learning method in section 4.2.1 in the original paper.

### 8.2.1  Image Pre-processing Pipeline

The same image pre-processing pipeline was applied to the images before running the learning methods mentioned below. The same processing steps mentioned in section 8.1 for the CIFAR100 dataset were deployed. For the CIFAR10 dataset, we normalized images using mean values of $\{0.4914, 0.4822, 0.4465\}$ and standard deviation values of $\{0.2023, 0.1994, 0.2010\}$ for the $RGB$ channels, respectively. The image augmentation methods used for the CIFAR100 training dataset were applied to the CIFAR10 training dataset.

### 8.2.2 Single Task Learning

In this experiment, we used the base network generated using the ResNet-18 network architecture. At each layer, $M$=15 modules were specified. For each task, $N$=4 modules were randomly selected and were assigned to the task. $ADAM$ optimizer with $CrossEntropy$ loss function was deployed. The initial learning rate for the optimizer was set to 0.001. The learning rate was halved at epochs 20, 30, and 40. Models were trained until convergence.

### 8.2.3 Sequential Learning

For this experiment, the same hyperparameter values and experiment setups were followed as for the single task learning experiments. During each run, the classes assigned to each task and the order of tasks to learn were set randomly. The model was trained on each task sequentially. The selected modules for each task were frozen after training for the task ended.

### 8.2.4 Parallel Learning

The hyperparameter values were set to the same values used in the single task learning experiments and the $batch\text{-}set$ size was set to 10.

### 8.2.5 MRN [19]

We used the same hyperparameter values that were specified for the CIFAR100 experiments.

### 8.2.6 pDarts [6]

We used the same hyperparameter values that were specified for the CIFAR100 experiments.

### 8.2.7 ENAS [9]

We used the same hyperparameter values that were specified for the CIFAR100 experiments.

### 8.2.8 MANAS [7]

We used the same hyperparameter values that were specified for the CIFAR100 experiments.

### 8.2.9 EWC [2]

We used the same hyperparameter values that were specified for the CIFAR100 experiments.

### 8.2.10 RPS [4]

For this experiment, we set the number of modules at each layer to $M$=15. The rest of hyperparameter values are the same as the values specified for the CIFAR100 experiments.

## 8.3 FiveTasks Experiments

This section provides the experimental setup for each compared learning method in section 4.2.2 in the original paper.

### 8.3.1 Image Pre-processing Pipeline

The same image pre-processing pipeline was applied to the images before running the learning methods mentioned below. We applied the same image pre-processing steps described by Hung et al. [5]. The processing steps are as follows: For all five image datasets, the images in the training dataset were normalized using mean values of $\{0.485, 0.456, 0.406\}$ and standard deviation values of $\{0.229, 0.224, 0.225\}$. Image augmentation methods including $RandomHorizontalFlip$, $RandomResizedCrop$ with crop size of 224 were applied on images in the training dataset. Images in the validation dataset were normalized using the same mean and standard deviation and were resized to image size of 224 by 224 pixels.

### 8.3.2 Single Task Learning

In this experiment, we used the base network generated using the ResNet-50 network architecture. At each layer, $M$=8 modules were specified. For each task, $N$=2 modules were randomly selected and were assigned to the task. The modules were initialised with a ResNet-50 model trained on the ImageNet dataset. $ADAM$ optimizer with $CrossEntropy$ loss function was deployed. The initial learning rate for the optimizer was set to 0.0001. The learning rate was halved at epochs 20, 30, and 40. Models were trained until convergence.

### 8.3.3 Parallel Learning

The hyperparameter values were set to the same values used in the single task learning experiments and the $batch\text{-}set$ size was set to 50.

## 9 APPENDIX C: Validation Accuracy for All the Compared Algorithms

We provided the mean and standard deviation of validation accuracy for DL models trained using single task learning, sequential learning and PaRT for CIFAR100 experiments in sections 4.1 and CIFAR10 and CIFAR100 experiments in section 4.2 in table 1 and table 2 in the original paper. In this section, we provide the mean and standard deviation of validation accuracy for DL models trained with algorithms that PaRT was compared to, namely, MRN [19], pDarts [6], ENAS [9], MANAS [7], EWC [2] and RPS [4]. The mean and standard deviation of validation accuracy were measured for five DL models trained using each learning method.

### 9.1 CIFAR100 Experiments

Table 5 shows the mean and standard deviation of DL models trained with RN [19], pDarts [6], ENAS [9], MANAS [7], EWC [2] and RPS [4] for the tasks defined on the CIFAR100 dataset in section 4.1 in the original paper. For each algorithm, five DL models were trained, and the mean and standard deviation of their validation accuracy were measured.

| Learning Method | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 | Task 6 | Task 7 | Task 8 | Task 9 | Task 10 | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|
| MRN [19] | 39.64 ± 7.04 | 38.06 ± 4.98 | 40.12 ± 5.56 | 35.24 ± 3.18 | 40.58 ± 3.51 | 38.98 ± 9.17 | 37.58 ± 4.07 | 29.76 ± 2.33 | 38.50 ± 7.05 | 37.60 ± 8.22 | 38.53 |
| pDarts [6] | 75.56 ± 6.36 | 70.20 ± 6.59 | 73.64 ± 13.84 | 64.98 ± 16.87 | 76.96 ± 7.77 | 77.32 ± 2.64 | 75.80 ± 6.29 | 72.88 ± 6.49 | 63.44 ± 10.77 | 79.92 ± 2.94 | 71.10 |
| ENAS [9] | 96.00 ± 4.90 | 84.00 ± 4.90 | 92.00 ± 11.66 | 94.00 ± 4.90 | 86.00 ± 4.90 | 94.00 ± 4.90 | 94.00 ± 8.00 | 90.00 ± 6.32 | 90.00 ± 0.00 | 70.00 ± 6.32 | 87.00 |
| MANAS [7] | 88.96 ± 2.60 | 84.38 ± 3.65 | 88.32 ± 2.18 | 86.84 ± 1.50 | 89.02 ± 2.35 | 87.68 ± 3.38 | 87.18 ± 1.38 | 84.64 ± 3.31 | 85.96 ± 6.21 | 88.58 ± 2.83 | 85.94 |
| EWC [2] | 69.54 ± 2.39 | 52.14 ± 2.56 | 51.78 ± 4.90 | 50.36 ± 2.25 | 55.78 ± 3.37 | 51.24 ± 4.64 | 53.88 ± 3.45 | 46.72 ± 4.03 | 51.70 ± 4.05 | 53.52 ± 4.37 | 53.67 |
| RPS [4] | 39.02 ± 18.67 | 30.18 ± 13.30 | 60.12 ± 7.92 | 68.84 ± 8.80 | 61.30 ± 8.26 | 66.24 ± 7.30 | 65.18 ± 12.46 | 70.54 ± 5.45 | 66.46 ± 6.12 | 68.30 ± 3.70 | 62.11 |

Table 5: Table compares the validation accuracy of DL models trained with different learning methods on 10 tasks defined on the CIFAR100 dataset. Definition of the tasks are provided in section 4.1 in the original paper. For each learning method, five DL models were trained, and the mean and standard deviation of validation accuracy for all five models are shown.

### 9.2 CIFAR10 and CIFAR100 Experiments

Table 5 shows the mean and standard deviation of DL models trained with RN [19], pDarts [6], ENAS [9], MANAS [7], EWC [2] and RPS [4] for the tasks defined on the CIFAR10 and CIFAR100 datasets in section 4.2 in the original paper. For each algorithms, five DL models were trained, and the mean and standard deviation of their validation accuracy were measured.

| Learning Method | Task 1 | Task 2 | Task 3 | Task 4 | Task 5 | Task 6 | Task 7 | Task 8 | Task 9 | Task 10 | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|
| MRN [19] | 34.37 | 57.20 | 50.85 | 63.56 | 70.40 | 52.02 | 70.05 | 51.73 | 50.34 | 65.94 | 59.00 |
| | ± 3.07 | ± 18.40 | ± 18.66 | ± 25.73 | ± 18.98 | ± 26.99 | ± 17.89 | ± 19.28 | ± 21.37 | ± 27.39 | ± 25.92 |
| pDarts [6] | 70.80 | 85.36 | 83.99 | 88.08 | 93.12 | 82.30 | 91.53 | 83.22 | 77.45 | 87.51 | 83.77 |
| | ± 7.46 | ± 13.53 | ± 10.03 | ± 11.38 | ± 9.19 | ± 11.95 | ± 11.44 | ± 11.44 | ± 14.34 | ± 13.82 | ± 15.60 |
| ENAS [9] | 68.80 | 66.40 | 58.00 | 64.80 | 73.20 | 71.60 | 77.60 | 64.00 | 62.00 | 66.40 | 70.00 |
| | ± 10.48 | ± 9.58 | ± 16.00 | ± 17.09 | ± 19.00 | ± 20.06 | ± 20.18 | ± 14.97 | ± 13.27 | ± 12.80 | ± 15.49 |
| MANAS [7] | 83.62 | 92.67 | 90.11 | 93.45 | 96.09 | 90.34 | 96.09 | 90.84 | 87.88 | 92.47 | 91.24 |
| | ± 1.10 | ± 5.87 | ± 6.46 | ± 6.34 | ± 5.53 | ± 6.24 | ± 4.85 | ± 4.58 | ± 7.46 | ± 8.28 | ± 7.77 |
| EWC [2] | 58.81 | 64.90 | 59.98 | 66.86 | 80.59 | 60.26 | 76.35 | 60.79 | 59.65 | 73.52 | 68.86 |
| | ± 1.95 | ± 15.56 | ± 17.94 | ± 21.75 | ± 16.92 | ± 20.37 | ± 19.31 | ± 17.17 | ± 19.77 | ± 24.45 | ± 21.26 |
| RPS [4] | 18.24 | 61.67 | 67.06 | 74.25 | 83.18 | 70.51 | 82.35 | 83.87 | 76.04 | 76.80 | 71.64 |
| | ± 8.99 | ± 14.19 | ± 16.31 | ± 14.14 | ± 19.94 | ± 18.84 | ± 9.27 | ± 7.47 | ± 17.10 | ± 23.62 | ± 28.38 |

Table 6: Table compares the validation accuracy of DL models trained with different learning methods on five tasks defined on the CIFAR10 dataset and five tasks defined on the CIFAR100 dataset. Definition of, the tasks are provided in section 4.2.1 in the original paper. For each learning method, five DL models were trained and the mean and standard deviation of validation accuracy for all five models are shown.