

Appendix

This document contains additional experimental results and analysis.

A. ANCL implementation with TA

We follow the provided implementation of ANCL [22] based on the FACIL [32] framework. We only apply the adaptation to the main teacher network fitted on the previous data, and, in cases where we use warmup, we only apply it to the new head of the main network. The auxiliary network does not undergo the teacher adaptation process, as it is already fit for the current task, and we do not see a reason to apply the warmup to this network. We keep the same training schedule and hyperparameters as described in Section 4.1.

B. Warmup phase

As mentioned in Section 4.1, we experimented with applying the warmup phase, a technique that recently appeared in several works [7, 27, 49]. With this method, before training the full student model on the new data, we first finetune the classification head added for this new task in isolation, while the rest of the network is kept frozen. Compared to the standard practice of finetuning the full student model on the new task, starting with a warmup phase reduces the initial cross-entropy loss at the start of the new task, and avoids overwriting the knowledge from the previous tasks with large gradient updates caused by random classification head initialization.

For training the new head during the warmup phase, we use SGD optimizer and OneCycle scheduler [48] with cosine annealing, maximum learning rate of 0.1 and the number of epochs of increasing learning rate set as 40. We train the new head for 200 epochs with early stopping and freeze updates of batch normalization statistics in the model during the training.

In Appendix D, we evaluate different CL benchmarks with TA and warmup separately. We generally find that the warmup works complementary to our method when we use smaller networks such as ResNet32 and train on smaller datasets. We hypothesize that for the settings with a larger dataset and a larger network better initialization of the new head does not translate to better performance in CIL because the overall magnitude of updates required to train the model in such settings is greater and leads to overwriting the initialization over the course of the training.

C. Ablation studies of Teacher Adaptation

C.1. Teacher Adaptation with different batch sizes.

Larger batch sizes should lead to better estimation of the global statistics in the batch normalization layer, which, in theory, should lead to reduced effectiveness of TA. To investigate the impact of the batch size on our method, we train the model on CIFAR100 split into 10 tasks with different batch sizes, tuning the learning rate and λ for each batch size. We present the results in Table 6. Regardless of the batch size, TA constantly improves upon standard training, and the difference between the two approaches does not scale with batch size.

Batch size	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
32	28.94±0.29	44.24±0.27	32.57±0.92	23.89±1.19
+TA	32.58±0.80	45.69±0.35	25.44±0.94	20.92±0.53
64	28.92±0.57	43.89±0.27	30.48±1.11	21.70±0.39
+TA	32.15±1.24	44.71±0.57	23.43±1.97	19.16±0.67
128	28.27±0.44	42.52±0.76	29.59±0.92	22.26±0.31
+TA	31.92±0.86	44.09±0.97	22.65±1.32	19.41±0.60
256	25.13±0.46	40.44±0.11	40.06±0.38	31.39±0.98
+TA	30.56±0.68	43.28±0.23	28.24±0.64	23.94±0.59
512	23.01±0.89	37.44±0.34	40.57±1.20	32.03±0.69
+TA	27.52±1.03	39.86±0.45	29.31±0.86	25.03±0.55

Table 6. Teacher adaptation with different batch sizes. Our method improves upon the standard training regardless of the batch size.

C.2. Different CNN architectures with Teacher Adaptation.

To test the robustness of our method, in addition to experiments with ResNets, we test it with other standard convolutional neural networks. We evaluate our method with MobileNetV2 [43], MobileNetV3-small [19] and VGG11 [46]. We use CIFAR100 and ImageNet100 split into 10 equally sized tasks and train the networks as described in Section 4.1. We show the results of those experiments in the Table 7. The final results obtained with tested networks are lower than with our standard baselines, as do not tune their hyperparameters extensively, but we still notice that using TA constantly improves the results, regardless of the architecture we use.

C.3. Batch normalization statistics investigation.

We investigate the batch normalization statistics in the teacher and student model throughout CIL on CIFAR100 split into 10 equally sized tasks. To measure the divergence between the normalization statistics in the models, we compute the average Kullback-Leibler divergence (KLD) be-

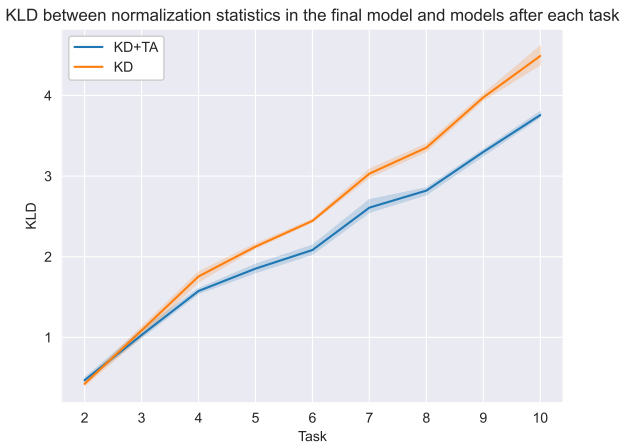
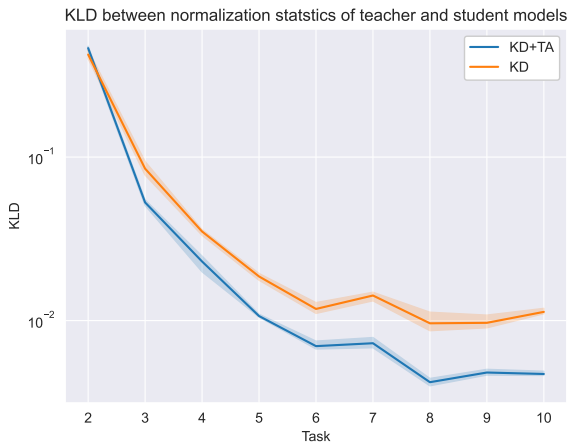


Figure 4. TA impact on batch normalization statistics. Left: KLD between the normalization statistics of the teacher and the student at the end of the training of all the tasks. Right: KLD between the normalization statistics of the model at the end of each task and the model trained on the first task.

CIFAR100				
Network	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
MobilenetV2	12.47±0.69	25.71±1.05	49.86±0.61	41.23±1.25
+TA	20.62±1.01	31.91±0.57	26.66±1.34	22.31±0.56
MobilenetV3-s	16.59±0.09	29.86±0.06	38.89±0.42	30.99±0.12
+TA	21.22±0.91	32.24±0.60	20.98±0.13	16.84±1.19
ResNet32	28.27±0.44	42.52±0.76	29.59±0.92	22.26±0.31
+TA	31.92±0.86	44.09±0.97	22.65±1.32	19.41±0.60
VGG11	17.10±0.12	34.34±0.39	62.36±0.62	52.99±0.97
+TA	24.87±0.35	42.50±0.15	43.96±0.62	32.92±0.51
ImageNet100				
Network	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
ResNet18	40.99±0.45	54.82±0.56	31.99±0.87	25.90±0.65
+TA	42.69±0.53	55.84±0.51	25.26±0.23	20.43±0.33
VGG11	24.35±0.32	44.72±0.50	61.62±0.10	46.73±0.31
+TA	30.67±0.31	48.81±0.29	50.04±0.76	37.60±0.56

Table 7. Teacher Adaptation with different network architectures on CIFAR100 and ImageNet100 split into 10 equally sized tasks.

tween the distributions of normalization statistics in every batch normalization layer. We show the results of this analysis in the Figure 4. Specifically, we measure KLD between the final teacher and student model at the end of each task, as well as the difference between the final student model trained for each task and the model initially learned after the first task. Applying TA leads to reduced KLD between the statistics in both teacher and student, and also student and initial model. This proves that our method leads to more stable representations throughout the training.

C.4. Hyperparameters for alternative adaptation methods ablation

The best results reported in Table 5 were obtained with teacher learning rate of $1e-7$. For variants with teacher per-training, the best results were obtained with training for 5 epochs. We keep λ fixed at 10 for all methods reported in the Table. While all variants that allow adaptation of batch norm statistics improve upon the baseline, the values of best hyperparameters are generally small, suggesting that the main source of improvement in all the methods comes from changes in normalization. This is further supported by the fact that all improvements vanish when normalization statistics are fixed.

C.5. Task recency bias with TA

We also conduct additional analysis of our method of Teacher Adaptation (TA) to understand the mechanism with which it improves upon the standard knowledge distillation. At Figure 5, we analyze task confusion matrices of standard knowledge distillation (LwF) and its extension with TA. We find that applying TA results in a model that is better at distinguishing between the tasks, and generally exhibits lower recency bias.

We hypothesize that the lower KD loss that we observe when using TA results in smaller updates to the model, so the difference between the magnitudes of logits learned for different tasks is smaller. Therefore, Teacher Adaptation helps to alleviate the recency bias in class-incremental learning.

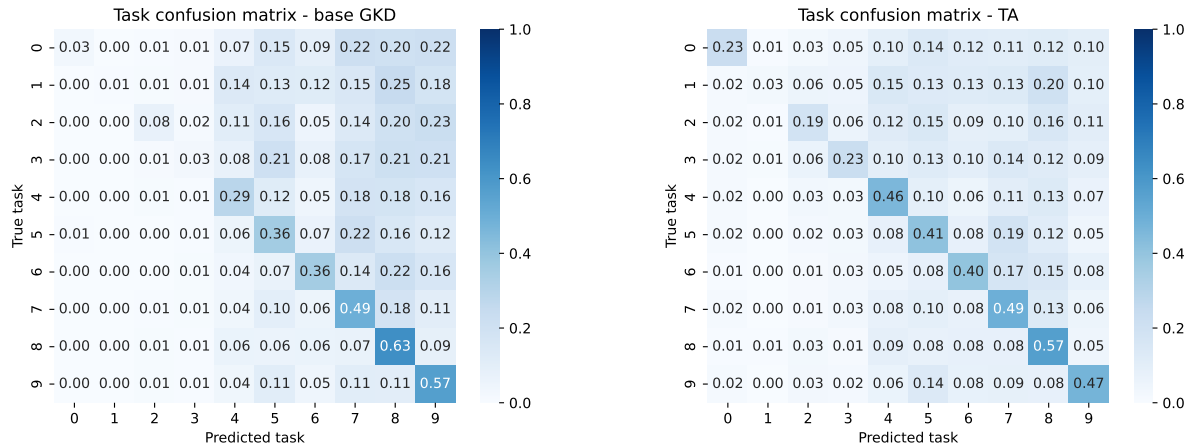


Figure 5. Task confusion matrix after learning all ten tasks on CIFAR100/10 for (left) base GKD and (right) GKD+TA. We see that TA leads to a model that is better at distinguishing between the first tasks and exhibits lower recency bias.

D. Additional experiments for different CL benchmarks

We provide the results of additional experiments conducted with our method and warmup. To shorten the notation, we denote the total number of tasks, including the first pretraining task (if present), with T , and the number of classes in the first task by S . For example, for CIFAR100 **T10S10** is a setting composed of 10 tasks with 10 classes each, while **T11S50** is a setting where the first task contains 50 classes and the next 10 tasks contain 5 classes each.

D.1. Standard benchmarks

In addition to Section 4.2, we conduct more experiments on CIFAR100, TinyImageNet200 and ImageNet100, adding two settings with a smaller number of tasks. We report final accuracy and forgetting in addition to incremental accuracy and forgetting. Additionally, we include the results for the warmup phase (WU) applied in isolation, without the Teacher Adaptation.

We report the results for CIFAR100, TinyImageNet200 and ImageNet100 in Table 8, Table 9, and Table 10 respectively.

D.2. Fine-grained classification using full datasets

We conduct experiments using the same datasets as in the Section 4.3.1, but without sampling classes, so each of the 6 datasets is treated as a single task. Additionally, we add the results for experiments with a reversed order of tasks to ensure that the ordering does not affect our results. We report final and incremental accuracy and forgetting in Table 11. We observe that Teacher Adaptation (TA) improves the final accuracy, but applying warmup (WU) in this setting results in significant drops in performance regardless of the order

of tasks.

D.3. Fine-grained classification with reverse order of datasets

We show results for the experiments conducted in Section 4.3.1, extended with additional metrics. Additionally, we add the results for the results with reversed order of tasks, as in Appendix D.2. Similar to the results in the main paper, the impact of our method is more visible in settings with a larger number of tasks.

D.4. DomainNet

We present the full experimental results for DomainNet in Table 3. In addition to incremental accuracy and forgetting, we also add final accuracy and forgetting to the metrics in the experiments. Additionally, we analyze the random initialization and impact of starting from pretrained checkpoint on ImageNet. For task-wise knowledge distillation (TKD) we apply λ equal to 1, for global knowledge distillation we notice that $\lambda = 10$ works better. We set the learning rate to 0.01 and use the scheduler with LR decay after 60, 120 and 180 epochs.

	T5S20				T6S50			
	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
GKD	37.63±0.52	48.80±0.36	23.13±2.28	19.30±2.37	40.74±0.72	51.93±1.24	25.87±0.75	18.90±0.29
+TA	40.84±0.23	50.08±0.26	16.76±1.68	16.66±1.25	43.18±1.66	52.22±1.28	18.73±1.61	14.09±0.66
+WU	38.82±0.06	49.69±0.26	22.08±0.61	19.34±1.24	43.19±0.95	53.42±0.50	24.54±2.17	18.94±2.69
+TA+WU	41.50±0.60	50.70±0.27	16.85±0.74	17.16±1.24	44.66±0.54	53.40±0.73	18.73±1.73	14.53±1.56
MKD	34.41±0.16	47.42±0.32	23.74±1.63	23.51±1.75	39.33±0.72	50.28±1.12	23.54±0.62	18.19±0.31
+TA	37.03±0.82	48.63±0.17	14.64±0.22	17.63±0.73	39.45±1.29	49.42±1.38	13.97±0.89	11.02±0.28
+WU	35.37±0.19	47.96±0.15	24.63±0.81	24.35±0.44	40.89±1.52	51.46±0.68	24.45±2.59	19.53±2.86
+TA+WU	38.26±0.24	49.58±0.06	15.18±0.41	17.96±0.65	41.42±0.78	51.04±1.17	15.00±1.07	12.03±1.32
TKD	38.33±0.70	49.56±0.48	24.78±2.58	25.04±2.55	41.19±0.42	52.07±1.35	17.31±1.20	15.02±0.39
+TA	41.12±0.35	50.87±0.15	18.12±1.55	21.09±1.39	41.36±0.89	51.88±0.80	12.84±1.25	11.42±0.64
+WU	40.03±0.66	50.50±0.17	22.79±0.73	24.19±1.61	43.16±0.26	53.51±0.39	18.40±1.73	16.03±2.77
+TA+WU	42.07±0.67	51.47±0.33	17.77±1.37	21.21±1.34	43.08±0.35	53.12±0.60	13.64±0.89	12.13±1.37
ANCL	37.13±0.33	48.60±0.27	32.41±0.26	31.71±0.81	43.41±0.62	53.40±0.49	23.54±1.29	18.08±0.40
+TA	41.04±0.22	50.36±0.43	24.07±0.81	26.18±0.51	44.64±1.65	53.33±0.66	18.43±3.12	14.44±1.29
+WU	38.94±0.57	49.26±0.61	29.84±0.88	30.98±1.98	44.74±0.30	54.25±0.29	22.42±1.75	17.49±1.56
+TA+WU	42.08±0.53	50.59±0.21	22.55±1.15	26.09±1.11	45.64±0.87	54.05±0.42	18.12±1.58	14.28±1.06
T10S10								
GKD	28.27±0.44	42.52±0.76	29.59±0.92	22.26±0.31	30.79±1.62	41.69±1.18	26.84±2.12	18.09±0.88
+TA	31.92±0.86	44.09±0.97	22.65±1.32	19.41±0.60	33.20±0.76	44.05±1.12	18.90±0.19	12.97±0.43
+WU	31.11±0.58	43.95±0.65	25.26±0.93	21.95±0.74	33.17±0.86	44.09±1.26	23.85±1.29	17.25±0.74
+TA+WU	33.76±0.78	45.25±1.02	20.87±0.13	19.87±0.34	34.72±0.77	46.27±1.09	19.13±1.77	13.98±0.98
MKD	25.41±1.07	39.36±0.70	44.04±0.98	42.74±0.52	31.31±1.22	41.04±0.93	22.56±0.94	15.37±0.33
+TA	29.97±1.44	43.55±0.96	32.37±2.17	32.61±0.80	30.90±0.69	41.67±1.35	14.76±0.25	10.51±0.36
+WU	27.32±1.28	40.78±0.95	37.89±0.22	39.03±0.53	31.63±1.16	42.45±0.78	22.57±1.57	16.98±1.23
+TA+WU	31.94±0.23	44.85±0.80	26.83±1.01	30.08±0.57	32.79±0.55	44.19±1.17	16.44±0.63	12.56±0.53
TKD	30.05±0.81	43.74±0.84	24.53±0.23	23.65±0.79	28.38±1.46	40.44±1.40	15.68±0.84	12.20±0.46
+TA	31.80±0.67	45.29±1.02	18.59±0.90	19.42±0.85	28.50±0.39	41.68±1.03	11.58±0.38	9.29±0.75
+WU	31.23±0.94	44.59±0.72	23.21±0.35	23.73±0.32	30.88±1.16	42.45±1.31	17.94±0.39	14.30±0.99
+TA+WU	33.14±1.03	46.21±0.86	17.55±0.66	20.45±0.57	30.60±0.37	44.22±1.08	14.65±0.49	11.79±0.57
ANCL	29.77±1.10	43.15±0.49	36.37±1.48	32.78±1.52	29.77±1.10	43.15±0.49	36.37±1.48	32.78±1.52
+TA	33.47±0.17	45.67±0.14	28.04±1.61	26.71±1.71	33.47±0.17	45.67±0.14	28.04±1.61	26.71±1.71
+WU	31.41±0.54	44.28±0.08	34.70±0.25	32.06±1.07	31.41±0.54	44.28±0.08	34.70±0.25	32.06±1.07
+TA+WU	34.17±0.30	46.73±0.20	28.67±0.63	26.86±1.17	34.17±0.30	46.73±0.20	28.67±0.63	26.86±1.17
T20S5								
GKD	15.59±0.32	31.89±0.45	43.28±0.56	34.68±1.87	10.10±0.71	17.64±0.93	15.29±0.27	9.67±0.26
+TA	19.55±0.24	35.99±0.79	30.38±2.08	23.32±1.79	11.99±0.66	19.37±1.73	9.05±0.63	8.31±0.68
+WU	19.22±0.80	34.67±0.55	37.05±1.03	31.14±1.41	10.41±1.01	21.87±0.16	17.92±0.48	11.02±0.65
+TA+WU	21.52±0.62	37.11±0.64	30.34±0.78	24.87±1.04	18.11±0.52	26.15±0.94	10.11±0.97	8.73±0.85
MKD	17.61±0.51	32.89±0.42	33.36±0.80	32.01±1.36	11.55±0.86	19.14±1.36	14.04±0.71	8.76±0.35
+TA	19.07±0.65	35.36±0.85	21.27±1.08	20.96±1.31	13.91±0.80	20.99±1.53	9.52±0.26	8.20±0.98
+WU	18.96±0.82	34.30±0.94	30.45±0.70	29.75±1.56	11.48±0.36	22.42±0.76	17.42±1.70	11.51±1.19
+TA+WU	20.52±0.85	36.79±0.70	21.88±0.81	21.84±0.33	18.15±0.70	26.10±0.75	11.64±1.04	9.17±0.20
TKD	19.39±0.41	34.58±0.34	22.06±0.46	21.13±1.17	7.88±0.08	14.64±0.33	7.96±0.47	6.02±0.54
+TA	18.30±0.50	34.62±0.92	15.22±1.25	14.72±1.28	9.05±0.64	16.66±1.66	7.17±0.53	6.88±0.36
+WU	20.77±0.59	35.72±0.09	21.14±0.88	20.93±1.62	9.55±0.53	18.01±0.59	12.63±1.04	9.90±0.99
+TA+WU	20.24±0.97	36.26±0.71	17.98±1.00	17.01±0.89	13.34±0.73	22.00±0.97	10.80±1.03	9.36±0.68
ANCL	19.21±0.75	34.32±0.41	39.97±1.81	36.74±1.38	11.04±0.40	21.84±1.33	18.31±1.64	11.79±0.51
+TA	20.63±1.04	37.52±0.78	31.07±1.00	26.77±0.49	18.96±0.91	26.47±0.41	12.69±1.49	10.33±0.57
+WU	21.53±0.38	36.06±0.73	37.34±0.53	35.34±1.52	11.81±0.24	23.81±0.88	20.05±0.95	13.18±0.34
+TA+WU	22.54±0.33	38.48±0.94	31.87±0.92	27.99±1.20	20.79±0.28	29.67±1.13	13.84±0.45	10.98±0.56

Table 8. Additional results for CIFAR100

	Basic order				Reverse order			
	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
GKD	26.71±0.53	46.74±0.94	46.57±0.89	37.47±1.07	26.61±1.07	54.88±0.62	52.71±1.01	39.34±0.41
+TA	31.16±0.75	46.38±0.74	44.09±0.68	37.84±0.98	33.61±0.40	54.97±0.35	41.89±0.64	36.63±0.58
+WU	23.96±1.53	36.37±0.95	34.78±1.80	46.72±1.22	23.60±1.33	43.02±0.56	39.25±1.36	51.69±0.76
+TA+WU	24.74±0.95	37.74±1.03	40.44±0.87	44.85±0.74	27.06±0.44	44.68±0.39	37.26±0.77	47.10±0.61

Table 11. Additional results for training on full fine-grained datasets in standard and reversed order.

	6 tasks, 20 classes each, base order				6 tasks, 20 classes each, reverse order			
	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
GKD	38.12±2.00	56.03±3.75	45.18±3.15	30.23±0.72	37.10±3.40	67.30±1.99	55.82±4.80	34.60±3.77
+TA	43.11±0.64	57.24±1.79	39.58±0.80	33.12±1.76	46.97±3.94	67.19±4.81	41.20±4.55	33.94±7.78
+WU	31.56±6.58	44.06±5.43	42.48±1.08	49.24±2.65	36.70±3.30	54.88±1.53	42.92±4.58	51.30±3.39
+TA+WU	42.85±2.34	52.96±1.20	35.97±2.75	38.77±2.45	44.05±1.94	58.92±2.47	38.27±2.75	45.16±4.57
MKD	41.32±0.81	60.12±1.59	40.30±3.98	26.42±1.60	43.96±3.89	70.65±1.67	42.83±5.11	25.52±3.43
+TA	39.93±1.95	55.77±2.23	36.22±1.07	31.43±2.28	46.94±4.42	67.15±4.48	37.01±6.67	31.07±8.38
+WU	34.65±2.72	48.01±1.97	44.86±0.61	47.07±2.40	35.49±5.51	56.83±2.79	46.14±6.67	48.15±5.95
+TA+WU	44.91±2.51	56.41±1.04	35.53±2.19	34.47±0.71	44.62±2.85	61.25±3.28	38.31±3.91	41.48±6.24
TKD	38.70±3.61	56.69±3.36	44.99±1.04	33.86±1.36	40.23±4.35	66.84±2.79	50.63±5.56	35.77±4.85
+TA	42.34±0.77	57.99±1.88	39.97±0.52	33.79±1.80	46.86±5.55	67.18±5.28	39.90±6.90	33.67±8.84
+WU	32.71±5.27	45.64±5.00	44.43±2.13	49.06±3.58	37.53±2.65	56.12±1.34	44.06±4.80	50.70±3.59
+TA+WU	43.04±1.30	53.68±1.56	37.14±1.52	39.73±2.92	44.52±2.14	59.37±2.39	38.32±2.53	45.65±4.61
	12 tasks, 10 classes each, base order				12 tasks, 10 classes each, reverse order			
	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
GKD	23.66±2.75	44.03±1.53	58.93±1.75	42.70±2.32	26.84±4.11	51.68±2.03	59.80±5.93	47.19±2.43
+TA	36.46±0.94	50.80±1.27	43.53±1.96	37.58±3.01	41.15±5.50	59.51±4.40	42.61±6.36	36.14±6.20
+WU	11.02±1.11	29.93±2.68	50.32±3.34	54.06±0.60	19.55±0.93	39.99±0.23	46.40±0.80	54.24±0.27
+TA+WU	36.00±0.89	49.46±2.09	35.49±1.90	35.81±2.02	44.79±4.28	52.88±5.06	26.81±5.14	40.33±3.03
MKD	25.35±4.36	45.19±2.91	55.04±2.71	43.81±1.56	31.81±5.29	53.63±2.31	52.38±6.99	43.86±3.05
+TA	35.10±1.92	51.16±1.84	37.96±1.84	34.89±2.60	41.95±5.36	60.11±3.87	36.16±6.75	33.18±5.64
+WU	15.61±6.67	32.93±3.65	49.31±7.17	54.64±1.52	26.19±8.31	42.21±1.31	44.49±6.28	54.20±1.67
+TA+WU	34.65±4.00	50.40±3.57	36.79±1.04	35.74±3.45	45.91±4.64	54.05±4.55	27.74±4.92	40.08±3.32
TKD	24.69±3.39	46.28±1.42	56.57±2.87	43.38±0.82	31.50±3.26	54.13±2.24	52.51±5.19	43.47±2.21
+TA	35.06±1.42	51.66±1.63	39.72±2.08	35.43±2.75	41.40±4.78	59.22±3.67	37.10±5.97	34.27±5.88
+WU	15.04±1.20	33.32±1.74	50.10±1.89	53.59±1.52	26.10±4.04	42.34±1.31	43.50±3.29	52.66±0.59
+TA+WU	35.81±1.27	50.12±2.31	34.73±0.52	34.64±1.99	43.91±5.45	52.38±5.54	26.77±4.76	40.35±3.60
	24 tasks, 5 classes each, base order				24 tasks, 5 classes each, reverse order			
	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
GKD	9.41±2.78	30.14±4.00	60.16±4.76	51.95±2.28	15.33±2.59	33.65±0.76	64.43±3.19	57.18±2.13
+TA	26.92±2.31	42.86±1.82	44.80±3.25	39.01±1.55	29.52±2.33	48.34±1.51	48.32±2.43	41.58±3.30
+WU	7.24±1.03	24.91±1.44	43.78±2.98	51.18±1.15	6.15±1.91	25.45±2.08	39.88±3.29	46.36±1.88
+TA+WU	27.30±1.97	45.28±2.25	34.75±4.58	33.67±2.48	34.14±9.24	42.65±1.21	19.92±11.23	32.38±8.24
MKD	11.57±2.70	32.61±0.55	63.70±2.58	55.91±1.77	15.53±2.96	34.85±1.40	62.28±3.64	56.57±2.11
+TA	27.88±1.58	45.49±1.82	37.65±1.16	34.48±0.89	31.82±1.51	49.47±1.16	38.89±1.18	35.98±2.07
+WU	7.52±4.24	24.62±1.56	44.85±6.75	52.59±2.59	4.57±3.36	21.36±4.30	40.82±4.75	51.73±0.78
+TA+WU	26.71±2.29	45.80±2.48	33.52±4.66	32.48±2.15	30.63±6.78	41.09±4.04	20.99±4.34	31.64±3.87
TKD	9.99±0.47	32.17±2.71	60.29±3.39	51.72±1.95	16.28±2.59	34.68±1.62	59.14±4.12	54.83±1.69
+TA	26.24±3.62	43.95±2.60	36.28±2.17	33.28±1.86	30.81±1.61	48.95±0.95	38.10±1.53	35.45±1.59
+WU	3.79±1.21	22.99±1.33	36.45±3.12	42.30±4.03	8.62±4.87	24.91±6.55	39.22±7.44	49.71±2.49
+TA+WU	20.24±4.08	37.22±6.11	23.39±2.94	24.53±3.01	25.21±3.57	47.46±5.21	28.40±4.04	26.15±2.93

Table 12. Additional results for the fine-grained classification datasets, extended by the experiments with reversed order of the datasets.

6 tasks 50 classes with $\lambda = 10$								
	Trained from scratch				Pre-trained on Imagenet			
	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
GKD	3.43±0.45	18.63±0.27	31.05±0.42	23.27±0.26	27.69±0.09	43.27±0.10	41.61±0.71	36.83±0.88
+TA	6.73±0.82	19.55±0.42	36.26±0.55	27.22±0.21	30.12±0.23	43.52±0.17	43.38±0.39	38.33±0.41
+TA+WU	5.83±0.70	18.68±0.36	33.00±0.59	31.12±0.38	28.83±0.26	39.00±0.12	42.64±0.56	49.44±0.37
ANCL	7.33±0.91	19.58±0.46	33.28±0.18	25.63±0.20	25.58±0.62	42.90±0.84	45.01±1.19	37.28±1.86
+TA	8.85±0.25	20.34±0.40	39.70±0.38	30.73±0.44	26.95±0.07	42.67±0.51	41.62±0.52	38.56±1.24
+TA+WU	9.65±0.40	20.39±0.09	37.57±0.12	32.83±0.66	29.15±0.98	42.44±0.80	34.23±0.85	37.09±1.33
12 tasks 25 classes with $\lambda = 10$								
	Trained from scratch				Pre-trained on Imagenet			
	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
GKD	2.48±0.41	14.45±0.25	29.01±0.63	29.04±0.37	17.69±0.84	35.98±0.96	50.12±1.00	43.00±1.66
+TA	4.88±0.46	16.25±0.46	36.47±1.06	33.03±0.19	28.02±0.68	38.89±0.52	40.87±0.64	41.10±0.42
+TA+WU	6.37±1.39	16.22±0.45	27.45±0.70	32.63±0.11	27.97±0.38	34.24±0.42	30.86±0.33	44.17±0.37
ANCL	3.97±0.50	14.82±0.41	32.80±0.61	33.46±0.40	14.13±1.35	33.34±0.55	55.72±0.16	49.05±0.65
+TA	6.54±0.51	17.19±0.06	39.77±0.76	37.40±0.38	23.42±1.03	35.81±0.18	40.86±0.85	43.84±0.23
+TA+WU	8.06±0.47	18.12±0.07	36.07±0.13	36.71±0.05	27.83±1.36	36.73±1.78	31.00±0.39	39.29±0.49
6 tasks 50 classes with $\lambda = 1$								
	Trained from scratch				Pre-trained on Imagenet			
	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
GKD	5.52±2.11	18.80±0.75	39.48±1.12	32.74±0.60	19.93±0.47	36.56±0.12	58.83±0.96	56.75±0.89
+TA	11.42±0.30	20.84±0.24	42.92±0.49	34.98±0.19	25.50±0.86	39.01±0.60	56.75±0.90	54.99±0.93
TKD	4.87±0.97	18.94±0.36	32.95±0.78	26.60±0.29	27.44±0.15	42.18±0.41	42.47±0.65	41.23±1.05
+TA	6.76±1.03	19.67±0.19	37.43±0.57	29.91±0.09	28.96±0.32	42.51±0.43	43.91±0.16	42.00±0.59
+TA+WU	5.88±0.82	18.65±0.27	33.50±0.96	32.33±0.49	29.55±0.30	40.62±0.07	42.08±0.55	46.86±0.39
MKD	3.24±1.10	18.74±0.52	26.83±0.83	19.10±0.40	30.41±0.47	45.70±0.30	32.96±0.29	27.48±0.60
+TA	3.99±0.38	18.04±0.16	29.16±0.29	22.70±0.25	27.47±0.22	42.91±0.04	31.05±0.31	29.43±0.04
+TA+WU	6.26±0.02	19.18±0.13	31.36±1.04	27.81±0.49	30.10±1.07	42.20±0.65	34.50±0.41	38.01±0.40
12 tasks 25 classes with $\lambda = 1$								
	Trained from scratch				Pre-trained on Imagenet			
	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$	$Acc_{Final} \uparrow$	$Acc_{Inc} \uparrow$	$Forg_{Final} \downarrow$	$Forg_{Inc} \downarrow$
GKD	2.88±0.25	14.38±0.82	38.89±3.42	39.78±2.23	9.62±0.56	27.99±0.40	64.74±0.21	59.93±1.13
+TA	6.80±0.28	17.02±0.19	47.37±0.82	43.44±0.16	19.00±0.83	33.90±0.62	58.20±0.71	54.97±0.71
TKD	3.39±0.43	12.28±0.99	41.22±0.28	33.32±0.33	22.66±1.59	37.56±1.03	41.65±1.73	41.21±1.97
+TA	5.10±0.47	16.84±0.42	34.03±0.61	33.46±0.41	29.39±0.37	38.85±0.31	34.51±0.54	39.91±0.39
+TA+WU	6.27±0.52	16.40±0.56	26.51±0.57	32.87±0.34	29.26±2.05	36.81±0.97	30.00±1.48	41.98±0.95
MKD	2.32±0.39	13.45±0.53	23.92±1.02	27.23±0.98	24.55±0.20	39.14±0.21	38.09±1.02	36.53±1.00
+TA	3.60±0.51	15.30±0.35	27.32±0.37	28.71±0.26	27.49±0.94	37.84±0.35	27.99±1.16	34.09±0.39
+TA+WU	5.46±1.20	15.89±0.96	22.23±1.98	30.62±0.48	30.34±1.32	37.79±1.11	24.81±1.05	37.00±1.84

Table 13. Additional results for DomainNet.