# T-CIL: Temperature Scaling using Adversarial Perturbation for Calibration in Class-Incremental Learning

## Supplementary Material

## 1. Detailed Algorithms

We provide full algorithms of Algorithm 2 (`MagSearch`) for magnitude search, and Algorithm 3 (`TempOpt`) for temperature optimization.

---

**Algorithm 2:** The magnitude search algorithm for perturbation (`MagSearch`).

---

**Input:** Model parameters $\theta = \{w, v\}$, exemplar set from new task $\mathcal{M}_{t,new}$, target temperature $T_{target}$, set of feature means $\mu$, tolerance $\delta$

1   $\epsilon_{low} \leftarrow 0.0, \epsilon_{high} \leftarrow 1.0$
2   **while** $\epsilon_{high} - \epsilon_{low} > \delta$ **do**
3      $\epsilon = \frac{\epsilon_{low} + \epsilon_{high}}{2.0}$
4      $\mathcal{M}_{t,new}^{\epsilon} = \{ \}$
5      **for** $(\mathbf{x}_e, y_e) \in \mathcal{M}_{t,new}$ **do**
6         $y_e' = \underset{c:\mu_c \in \mu, c \neq y_e}{\arg\max} \|\phi_v(\mathbf{x}_e) - \mu_c\|$
7         $\mathbf{x}_e^{adv} = \mathbf{x}_e - \epsilon \, \mathrm{sign}(\nabla_{\mathbf{x}_e} \mathcal{L}_{CE}(\mathbf{x}_e, y_e'; \theta))$
8         $\mathcal{M}_{t,new}^{\epsilon} \leftarrow \mathcal{M}_{t,new}^{\epsilon} \cup \{(\mathbf{x}_e^{adv}, y_e)\}$
9      $T = \texttt{TempOpt}(\mathcal{M}_{t,new}^{\epsilon}, \theta)$
10      **if** $T < T_{target}$ **then**
11         $\epsilon_{low} \leftarrow \epsilon$
12      **else**
13         $\epsilon_{high} \leftarrow \epsilon$

**Output:** $\frac{\epsilon_{low} + \epsilon_{high}}{2.0}$

---

**Algorithm 3:** The temperature optimization algorithm (`TempOpt`).

---

**Input:** Dataset $\mathcal{D}$, model parameters $\theta$

1   $T \leftarrow 1$ ;         // Initialize $T$
2   **while** *not converge* **do**
3      **for** $(\mathbf{x}, y) \in \mathcal{D}$ **do**
4         Update $T$ to minimize $\mathcal{L}_{CE}(\mathbf{x}, y; \theta, T)$

**Output:** $T$

---

## 2. Computational Complexity Analysis

In this section, we analyze the computational complexity of T-CIL. T-CIL consists of four main components: temperature optimization, feature means calculation, perturbation magnitude search, and memory update. The complexity of temperature optimization is $O(M)$, where $M$ represents the memory size, since we optimize the temperature on the set of perturbed exemplars for memory. Similarly, calculating feature means requires $O(M)$ operations. The adversarial search process takes $O(M)$ time (with a constant factor $k$ for perturbation iterations), and applying perturbations also has a complexity of $O(M)$. Therefore, the overall computational complexity of T-CIL for a single incremental task is also $O(M)$.

As conventional class-incremental learning approaches operate with $O(T(N_{\text{new}} + M))$ complexity, T-CIL maintains a lighter $O(M)$ complexity, where $T$ is the number of tasks, $N_{\text{new}}$ is the number of new task data points, and $M$ is the memory size. Consequently, integrating T-CIL with existing class-incremental frameworks preserves the asymptotic computational efficiency, as $M$ remains constant and substantially less than $N_{\text{new}}$. With a small overhead, T-CIL is an efficient approach that can be practically combined with existing class-incremental learning methods.

## 3. Experiments

### 3.1. Inapplicability of PerturbTS on the CIFAR-10

The reason PerturbTS [37] is not applicable on the CIFAR-10 in a class-incremental learning setup is that, with only two new classes per task, the model quickly fits to the data, making it impossible to achieve the designated accuracy reduction through perturbation. This overfitting prevents the perturbation magnitude optimization from converging.

### 3.2. Detailed Experimental Settings

To obtain the best calibration performance of T-CIL with the minimal impact on accuracy, we use a new-task validation set whose size varies depending on the class-incremental learning technique and dataset used. The new-task validation set sizes are listed in Table 4. The effect of the new-task validation set size will be explained later.

Class-incremental learning techniques require specific training parameters. For both EEIL [3] and WA [45], we set the knowledge distillation temperature to 2. EEIL and DER [39] incorporate a balanced fine-tuning phase. For this phase, we train the model for 30 epochs with 10 tasks and 100 epochs with 20 tasks.

After model training, we store a subset of new-task data used for training to the memory by uniformly sampling examples from each class. As the memory size is fixed, we remove some existing exemplars to accommodate the new-task data.

Table 4. The size of the new-task validation set for each class-incremental learning method and dataset used in the main experiments.

| Method | CIFAR-10 | CIFAR-100 | Tiny-ImageNet |
|--------|----------|-----------|---------------|
| ER | 500 | 500 | 100 |
| EEIL | 300 | 300 | 200 |
| WA | 100 | 500 | 100 |
| DER | 500 | 1000 | 100 |

We make a fair comparison between post-hoc calibration methods including T-CIL versus vanilla class-incremental learning techniques in terms of training data. In particular, whenever we take a (minimal) validation set from the training data, we only train models on the remaining training data. In comparison, the vanilla techniques always train on the full training set.

### 3.3. Additional Experiments

**Full experimental results**  We evaluate T-CIL against five calibration baselines (Cal method) in combination with four class-incremental learning techniques (CIL method) across three datasets. Table 5 presents a comprehensive comparison of all possible combinations between post-hoc calibration methods and class-incremental learning techniques.

Overall, T-CIL outperforms the five calibration baselines when integrated with four existing class-incremental learning techniques across three datasets. Notably, T-CIL consistently shows low calibration errors compared to all the baselines. While PerturbTS achieves the best calibration performances when combined with EEIL and WA on the Tiny-ImageNet, its effectiveness is inconsistent. In addition, PerturbTS is not applicable to the CIFAR-10 dataset and exhibits unusually high calibration errors on the CIFAR-100 dataset. In contrast, T-CIL demonstrates robust and superior performances across all experimental settings, consistently achieving lower calibration errors regardless of the underlying class-incremental learning technique. This comprehensive evaluation validates that T-CIL is a more reliable and versatile approach for addressing calibration challenges in class-incremental learning scenarios.

**Expansion of incremental tasks**  We evaluate the performance when expanding incremental tasks from 10 to 20 tasks, with results presented in Table 6. For all class-incremental learning techniques, we use 100 samples from each new task as a validation set on both CIFAR-100 and Tiny-ImageNet.

The results show that T-CIL outperforms most vanilla class-incremental learning techniques, with WA being the only exception. As explained in Section 6.2, WA scales the output logits corresponding to the new task only after train-
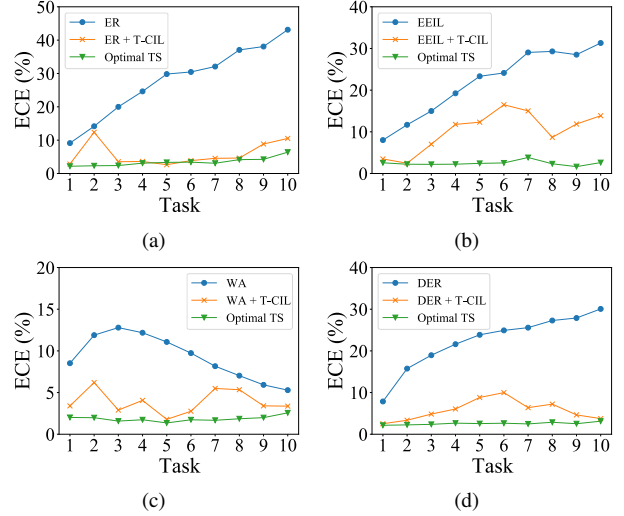


Figure 6. ECE (%) comparison after training each task on the CIFAR-100 among existing class-incremental learning techniques, their combinations with T-CIL, and the optimal TS. The existing techniques include: (a) ER, (b) EEIL, (c) WA, and (d) DER.

ing. This scaling of specific logits may not align with the insight behind T-CIL's perturbation direction policy. Nevertheless, T-CIL significantly improves the calibration performance of poorly calibrated class-incremental learning techniques.

**Varying size of memory**  We analyze how memory size affects the ECE of T-CIL compared to the vanilla method without calibration on the CIFAR-100 and Tiny-ImageNet, as shown in Figure 7. Using ER as our base class-incremental learning technique, we experiment with memory sizes of 500, 1,000, 1,500, and 2,000 samples for CIFAR-100, and 1,000, 2,000, 3,000, and 4,000 samples for Tiny-ImageNet. We use a new-task validation set sized of 100.

Our results demonstrate that T-CIL consistently achieves significantly lower ECE than the vanilla method across all memory sizes. Although smaller memory sizes lead to decreased model accuracy and consequently worse calibration performance, T-CIL with just 500 memory samples still outperforms the vanilla method using 2,000 samples in terms of calibration quality.

**Varying size of new-task validation set**  We vary the size of the new-task validation set and evaluate ECE and accuracy on CIFAR-100, using ER as a base class-incremental learning technique. We present the results in Figure 8 and Table 7.

Figure 8 demonstrates that T-CIL is effective even with a small-sized validation set. As the validation set size increases, ECE decreases. However, increasing validation set size leads to accuracy drop due to the smaller training set

size. These trends in ECE and accuracy indicate that T-CIL only requires small new-task validation set for calibrating the model effectively while minimizing the impact on accuracy.

**Additional ECE progressions**   We present the progression of ECE across tasks on CIFAR-100 and Tiny-ImageNet of vanilla, T-CIL, and Optimal TS when combined with four class-incremental learning techniques in Figure 6 and Figure 9. These figures show the progression of ECE throughout tasks, where we compare T-CIL against an ideal scenario where we run *TS* on the test set of both old and new tasks and thus obtain the best achievable calibration performance (called "Optimal TS"). When combined with various class-incremental learning techniques, T-CIL consistently demonstrates strong calibration performance across all tasks, with calibration errors approaching Optimal TS's ideal performance.

Table 5. Performance comparison between T-CIL and five baselines when integrated with four class-incremental learning techniques on three datasets.

| CIL Method | Cal Method | CIFAR-10 | | | CIFAR-100 | | | Tiny-ImageNet | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Acc. ($\uparrow$) | ECE ($\downarrow$) | AECE ($\downarrow$) | Acc. ($\uparrow$) | ECE ($\downarrow$) | AECE ($\downarrow$) | Acc. ($\uparrow$) | ECE ($\downarrow$) | AECE ($\downarrow$) |
| ER [6] | Vanilla | $65.61_{\pm 0.49}$ | $28.16_{\pm 0.33}$ | $28.12_{\pm 0.32}$ | $56.51_{\pm 0.37}$ | $27.66_{\pm 0.29}$ | $27.64_{\pm 0.28}$ | $31.95_{\pm 0.54}$ | $32.56_{\pm 0.36}$ | $32.55_{\pm 0.35}$ |
| | TS [13] | $65.86_{\pm 0.09}$ | $23.97_{\pm 0.82}$ | $23.92_{\pm 0.81}$ | $56.25_{\pm 0.62}$ | $16.28_{\pm 0.32}$ | $16.22_{\pm 0.37}$ | $31.48_{\pm 0.39}$ | $19.44_{\pm 0.75}$ | $19.44_{\pm 0.75}$ |
| | ETS [44] | $65.86_{\pm 0.09}$ | $22.46_{\pm 1.20}$ | $22.39_{\pm 1.21}$ | $56.25_{\pm 0.62}$ | $16.83_{\pm 0.47}$ | $16.80_{\pm 0.49}$ | $31.48_{\pm 0.39}$ | $19.81_{\pm 0.65}$ | $19.83_{\pm 0.65}$ |
| | IRM [44] | $65.86_{\pm 0.09}$ | $22.09_{\pm 1.73}$ | $21.83_{\pm 1.73}$ | $56.25_{\pm 0.62}$ | $17.50_{\pm 0.64}$ | $17.39_{\pm 0.61}$ | $31.48_{\pm 0.39}$ | $19.77_{\pm 0.72}$ | $19.73_{\pm 0.78}$ |
| | PerturbTS [37] | n/a | n/a | n/a | $56.25_{\pm 0.62}$ | $16.49_{\pm 2.19}$ | $16.49_{\pm 2.18}$ | $31.48_{\pm 0.39}$ | $10.60_{\pm 0.60}$ | $10.58_{\pm 0.61}$ |
| | **T-CIL** | $65.86_{\pm 0.09}$ | $\mathbf{17.70}_{\pm \mathbf{2.60}}$ | $\mathbf{17.64}_{\pm \mathbf{2.59}}$ | $56.25_{\pm 0.62}$ | $\mathbf{5.74}_{\pm \mathbf{0.53}}$ | $\mathbf{5.75}_{\pm \mathbf{0.50}}$ | $31.48_{\pm 0.39}$ | $\mathbf{8.12}_{\pm \mathbf{0.38}}$ | $\mathbf{8.12}_{\pm \mathbf{0.41}}$ |
| EEIL [3] | Vanilla | $77.67_{\pm 0.74}$ | $15.48_{\pm 0.75}$ | $15.45_{\pm 0.74}$ | $60.61_{\pm 0.33}$ | $21.96_{\pm 0.29}$ | $21.94_{\pm 0.28}$ | $37.44_{\pm 0.85}$ | $29.69_{\pm 0.36}$ | $29.68_{\pm 0.36}$ |
| | TS | $76.91_{\pm 1.27}$ | $\mathbf{10.20}_{\pm \mathbf{0.94}}$ | $\mathbf{10.16}_{\pm \mathbf{0.93}}$ | $60.74_{\pm 0.37}$ | $13.25_{\pm 0.60}$ | $13.14_{\pm 0.55}$ | $37.16_{\pm 0.91}$ | $13.16_{\pm 0.76}$ | $13.15_{\pm 0.73}$ |
| | ETS | $76.91_{\pm 1.27}$ | $10.31_{\pm 0.71}$ | $10.29_{\pm 0.72}$ | $60.74_{\pm 0.37}$ | $13.89_{\pm 0.47}$ | $13.83_{\pm 0.42}$ | $37.16_{\pm 0.91}$ | $13.52_{\pm 0.73}$ | $13.51_{\pm 0.71}$ |
| | IRM | $76.91_{\pm 1.27}$ | $10.51_{\pm 0.59}$ | $10.32_{\pm 0.53}$ | $60.74_{\pm 0.37}$ | $15.20_{\pm 0.48}$ | $15.09_{\pm 0.51}$ | $37.16_{\pm 0.91}$ | $14.66_{\pm 0.76}$ | $14.69_{\pm 0.76}$ |
| | PerturbTS | n/a | n/a | n/a | $60.74_{\pm 0.37}$ | $40.34_{\pm 3.84}$ | $40.34_{\pm 3.84}$ | $37.16_{\pm 0.91}$ | $\mathbf{8.81}_{\pm \mathbf{0.59}}$ | $\mathbf{8.81}_{\pm \mathbf{0.61}}$ |
| | **T-CIL** | $76.91_{\pm 1.27}$ | $10.49_{\pm 2.34}$ | $10.46_{\pm 2.32}$ | $60.74_{\pm 0.37}$ | $\mathbf{10.30}_{\pm \mathbf{1.10}}$ | $\mathbf{10.22}_{\pm \mathbf{1.06}}$ | $37.16_{\pm 0.91}$ | $15.58_{\pm 0.76}$ | $15.56_{\pm 0.76}$ |
| WA [45] | Vanilla | $73.06_{\pm 0.57}$ | $19.10_{\pm 0.50}$ | $19.07_{\pm 0.51}$ | $64.34_{\pm 0.40}$ | $8.89_{\pm 0.64}$ | $8.86_{\pm 0.63}$ | $39.66_{\pm 0.88}$ | $10.97_{\pm 0.41}$ | $10.96_{\pm 0.43}$ |
| | TS | $72.75_{\pm 0.47}$ | $18.22_{\pm 0.69}$ | $18.19_{\pm 0.69}$ | $64.02_{\pm 0.06}$ | $5.93_{\pm 0.24}$ | $5.88_{\pm 0.28}$ | $38.59_{\pm 0.44}$ | $13.24_{\pm 1.08}$ | $13.28_{\pm 1.09}$ |
| | ETS | $72.75_{\pm 0.47}$ | $17.96_{\pm 0.67}$ | $17.92_{\pm 0.69}$ | $64.02_{\pm 0.06}$ | $5.77_{\pm 0.39}$ | $5.75_{\pm 0.42}$ | $38.59_{\pm 0.44}$ | $13.01_{\pm 1.12}$ | $13.06_{\pm 1.12}$ |
| | IRM | $72.75_{\pm 0.47}$ | $18.03_{\pm 0.94}$ | $17.58_{\pm 0.97}$ | $64.02_{\pm 0.06}$ | $6.61_{\pm 0.47}$ | $6.51_{\pm 0.46}$ | $38.59_{\pm 0.44}$ | $10.88_{\pm 0.46}$ | $10.87_{\pm 0.48}$ |
| | PerturbTS | n/a | n/a | n/a | $64.02_{\pm 0.06}$ | $46.55_{\pm 2.20}$ | $46.54_{\pm 2.20}$ | $38.59_{\pm 0.44}$ | $\mathbf{8.63}_{\pm \mathbf{1.07}}$ | $\mathbf{8.61}_{\pm \mathbf{1.10}}$ |
| | **T-CIL** | $72.75_{\pm 0.47}$ | $\mathbf{15.61}_{\pm \mathbf{0.23}}$ | $\mathbf{15.58}_{\pm \mathbf{0.22}}$ | $64.02_{\pm 0.06}$ | $\mathbf{3.87}_{\pm \mathbf{0.52}}$ | $\mathbf{3.84}_{\pm \mathbf{0.55}}$ | $38.59_{\pm 0.44}$ | $11.43_{\pm 1.00}$ | $11.46_{\pm 1.04}$ |
| DER [39] | Vanilla | $74.53_{\pm 0.48}$ | $21.81_{\pm 0.46}$ | $21.78_{\pm 0.47}$ | $69.98_{\pm 0.69}$ | $22.38_{\pm 0.37}$ | $22.35_{\pm 0.35}$ | $46.62_{\pm 2.84}$ | $39.00_{\pm 1.72}$ | $38.99_{\pm 1.72}$ |
| | TS | $74.93_{\pm 0.35}$ | $17.27_{\pm 0.37}$ | $17.25_{\pm 0.37}$ | $69.98_{\pm 0.58}$ | $6.16_{\pm 0.25}$ | $6.04_{\pm 0.26}$ | $47.79_{\pm 0.47}$ | $11.29_{\pm 0.66}$ | $11.26_{\pm 0.66}$ |
| | ETS | $74.93_{\pm 0.35}$ | $16.82_{\pm 0.28}$ | $16.79_{\pm 0.29}$ | $69.98_{\pm 0.58}$ | $6.12_{\pm 0.36}$ | $6.03_{\pm 0.37}$ | $47.79_{\pm 0.47}$ | $7.83_{\pm 0.63}$ | $7.86_{\pm 0.58}$ |
| | IRM | $74.93_{\pm 0.35}$ | $17.04_{\pm 0.45}$ | $16.80_{\pm 0.48}$ | $69.98_{\pm 0.58}$ | $7.88_{\pm 0.33}$ | $8.03_{\pm 0.40}$ | $47.79_{\pm 0.47}$ | $10.47_{\pm 0.61}$ | $11.01_{\pm 0.61}$ |
| | PerturbTS | n/a | n/a | n/a | $69.98_{\pm 0.58}$ | $52.21_{\pm 5.97}$ | $52.21_{\pm 5.97}$ | $47.79_{\pm 0.47}$ | $15.56_{\pm 1.18}$ | $15.55_{\pm 1.18}$ |
| | **T-CIL** | $74.93_{\pm 0.35}$ | $\mathbf{12.70}_{\pm \mathbf{1.35}}$ | $\mathbf{12.68}_{\pm \mathbf{1.34}}$ | $69.98_{\pm 0.58}$ | $\mathbf{4.37}_{\pm \mathbf{0.59}}$ | $\mathbf{4.34}_{\pm \mathbf{0.60}}$ | $47.79_{\pm 0.47}$ | $\mathbf{6.91}_{\pm \mathbf{0.86}}$ | $\mathbf{6.90}_{\pm \mathbf{0.84}}$ |

Table 6. T-CIL performance combined with four existing class-incremental learning techniques on two datasets, each containing 20 incremental tasks.

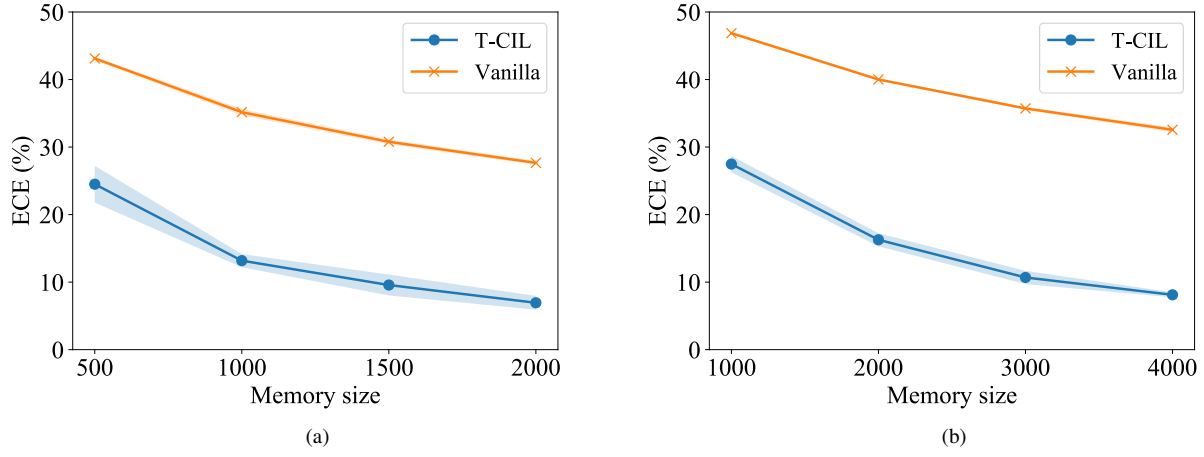| Method | CIFAR-100 | | | Tiny-ImageNet | | |
|---|---|---|---|---|---|---|
| | Acc. ($\uparrow$) | ECE ($\downarrow$) | AECE ($\downarrow$) | Acc. ($\uparrow$) | ECE ($\downarrow$) | AECE ($\downarrow$) |
| ER | $54.81_{\pm 0.70}$ | $28.95_{\pm 0.84}$ | $28.93_{\pm 0.84}$ | $30.43_{\pm 0.51}$ | $34.90_{\pm 0.16}$ | $34.88_{\pm 0.16}$ |
| ER+T-CIL | $54.48_{\pm 1.95}$ | $\mathbf{7.08}_{\pm \mathbf{0.49}}$ | $\mathbf{7.07}_{\pm \mathbf{0.46}}$ | $30.62_{\pm 0.29}$ | $\mathbf{10.72}_{\pm \mathbf{0.88}}$ | $\mathbf{10.79}_{\pm \mathbf{0.86}}$ |
| EEIL | $55.36_{\pm 1.19}$ | $25.17_{\pm 0.84}$ | $25.15_{\pm 0.84}$ | $33.54_{\pm 0.60}$ | $28.89_{\pm 0.27}$ | $28.89_{\pm 0.27}$ |
| EEIL+T-CIL | $55.86_{\pm 0.93}$ | $\mathbf{14.38}_{\pm \mathbf{1.34}}$ | $\mathbf{14.35}_{\pm \mathbf{1.35}}$ | $34.22_{\pm 0.18}$ | $\mathbf{24.45}_{\pm \mathbf{0.67}}$ | $\mathbf{24.46}_{\pm \mathbf{0.67}}$ |
| WA | $59.30_{\pm 0.64}$ | $\mathbf{7.39}_{\pm \mathbf{0.35}}$ | $\mathbf{7.37}_{\pm \mathbf{0.35}}$ | $35.99_{\pm 0.97}$ | $\mathbf{10.06}_{\pm \mathbf{0.95}}$ | $\mathbf{10.06}_{\pm \mathbf{0.99}}$ |
| WA+T-CIL | $58.89_{\pm 0.67}$ | $7.45_{\pm 0.72}$ | $7.47_{\pm 0.72}$ | $36.80_{\pm 0.42}$ | $19.56_{\pm 0.47}$ | $19.57_{\pm 0.48}$ |
| DER | $68.90_{\pm 0.31}$ | $24.05_{\pm 0.45}$ | $24.02_{\pm 0.45}$ | $48.63_{\pm 0.75}$ | $42.44_{\pm 0.56}$ | $42.43_{\pm 0.56}$ |
| DER+T-CIL | $68.84_{\pm 0.61}$ | $\mathbf{5.78}_{\pm \mathbf{0.92}}$ | $\mathbf{5.75}_{\pm \mathbf{0.91}}$ | $48.33_{\pm 1.31}$ | $\mathbf{5.90}_{\pm \mathbf{0.97}}$ | $\mathbf{5.91}_{\pm \mathbf{0.96}}$ |

(a)

(b)

Figure 7. T-CIL performance when varying the memory size on (a) CIFAR-100 and (b) Tiny-ImageNet.
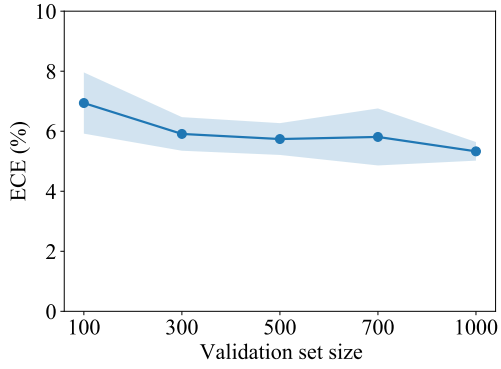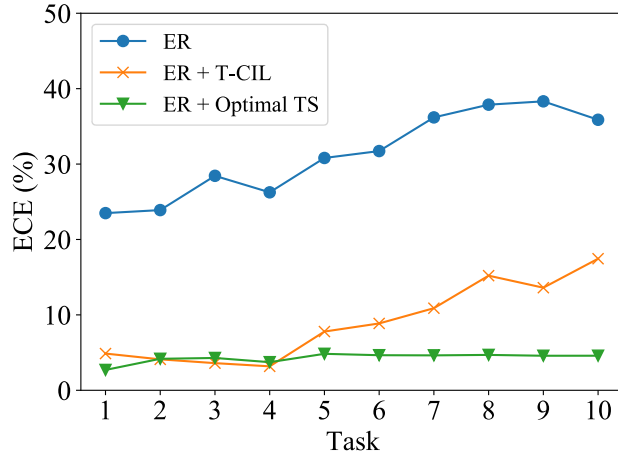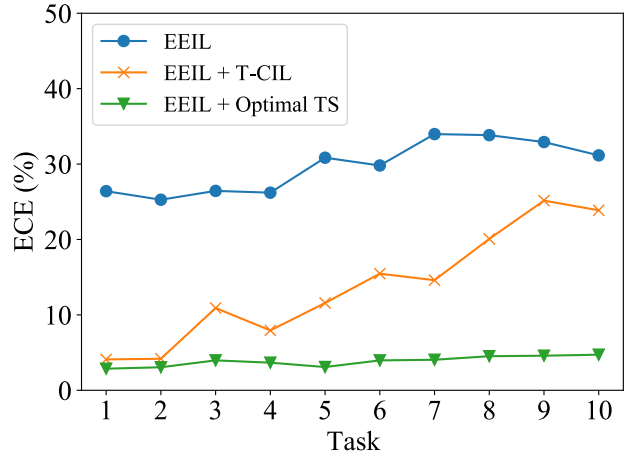


Figure 8. T-CIL performance when varying the size of the new-task validation set on the CIFAR-100.

Table 7. T-CIL performance when varying the size of the new-task validation set on the CIFAR-100.
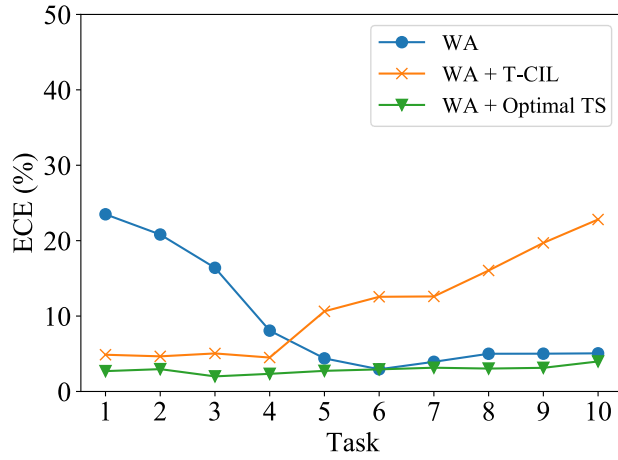
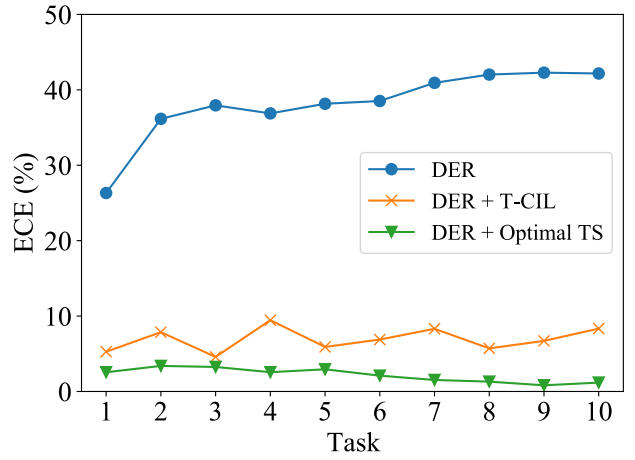| Val size | 100 | 300 | 500 | 700 | 1000 |
|---|---|---|---|---|---|
| ECE (%) | $6.94_{\pm 1.02}$ | $5.91_{\pm 0.56}$ | $5.74_{\pm 0.53}$ | $5.81_{\pm 0.95}$ | $5.33_{\pm 0.31}$ |
| Acc (%) | $56.97_{\pm 0.65}$ | $56.95_{\pm 0.34}$ | $56.25_{\pm 0.62}$ | $56.55_{\pm 0.80}$ | $56.08_{\pm 0.57}$ |

Figure 9. ECE progression after training each task on the Tiny-ImageNet among existing class-incremental learning techniques, their combinations with T-CIL, and the optimal TS. The existing techniques include: (a) ER, (b) EEIL, (c) WA, and (d) DER.