

ErrorEraser: Unlearning Data Bias for Improved Continual Learning

Xuemei Cao
School of Computing and Artificial
Intelligence, Southwestern University
of Finance and Economics
Chengdu, China
caoxuemei.pqz@gmail.com

Hanlin Gu*
WeBank
Shenzhen, China
allengu@webank.com

Xin Yang*
School of Computing and Artificial
Intelligence, Southwestern University
of Finance and Economics
Chengdu, China
yangxin@swufe.edu.cn

Bingjun Wei
School of Computing and Artificial
Intelligence, Southwestern University
of Finance and Economics
Chengdu, China
dutu66cai@gmail.com

Haoyang Liang
School of Computing and Artificial
Intelligence, Southwestern University
of Finance and Economics
Chengdu, China
hpsigma@gmail.com

Xiangkun Wang
School of Computing and Artificial
Intelligence, Southwestern University
of Finance and Economics
Chengdu, China
xiangkunwang18@gmail.com

Tianrui Li
School of Computing and Artificial
Intelligence, Southwest Jiaotong
University
Chengdu, China
trli@swjtu.edu.cn

ABSTRACT

Continual Learning (CL) primarily aims to retain knowledge to prevent catastrophic forgetting and transfer knowledge to facilitate learning new tasks. Unlike traditional methods, we propose a novel perspective: CL not only needs to prevent forgetting, but also requires intentional forgetting. This arises from existing CL methods ignoring biases in real-world data, leading the model to learn spurious correlations that transfer and amplify across tasks. From feature extraction and prediction results, we find that data biases simultaneously reduce CL's ability to retain and transfer knowledge. To address this, we propose ErrorEraser, a universal plugin that removes erroneous memories caused by biases in CL, enhancing performance in both new and old tasks. ErrorEraser consists of two modules: Error Identification and Error Erasure. The former learns the probability density distribution of task data in the feature space without prior knowledge, enabling accurate identification of potentially biased samples. The latter ensures only erroneous knowledge is erased by shifting the decision space of representative outlier samples. Additionally, an incremental feature distribution learning strategy is designed to reduce the resource overhead during error identification in downstream tasks. Extensive

experimental results show that ErrorEraser significantly mitigates the negative impact of data biases, achieving higher accuracy and lower forgetting rates across three types of CL methods. The code is available at <https://github.com/diada/ErrorEraser>.

CCS CONCEPTS

• **Computing methodologies** → **Machine learning**;

KEYWORDS

Continual learning; Machine unlearning; Selective Forgetting; Catastrophic Forgetting; Knowledge Transfer

ACM Reference Format:

Xuemei Cao, Hanlin Gu, Xin Yang, Bingjun Wei, Haoyang Liang, Xiangkun Wang, and Tianrui Li. 2025. ErrorEraser: Unlearning Data Bias for Improved Continual Learning. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V.2 (KDD '25)*, August 3–7, 2025, Toronto, ON, Canada. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/00000000.000000>

1 INTRODUCTION

Continual Learning (CL) or Lifelong Learning [6, 12, 15, 17] has been proposed in recent years to address the challenges posed by non-stationary and dynamically changing data for neural networks. Its primary objectives are: (1) Preventing catastrophic forgetting: This aims to retain learned knowledge, preventing the forgetting of previously learned tasks when learning new ones [43]; (2) Facilitating knowledge transfer: This seeks to leverage knowledge from previous tasks to guide new tasks, thereby improving accuracy or saving learning time [15]. As CL methods gain widespread attention, they have evolved into three scenarios based on different

*Co-Corresponding Author.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
KDD '25, August 3–7, 2025, Toronto, ON, Canada.
© 2025 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-XXXX-X/18/06
<https://doi.org/10.1145/00000000.000000>

learning settings [43]: task incremental, class incremental, and domain incremental. Additionally, current popular CL methods can be categorized into three types [25]: memory replay-based, parameter isolation-based, and regularization-based.

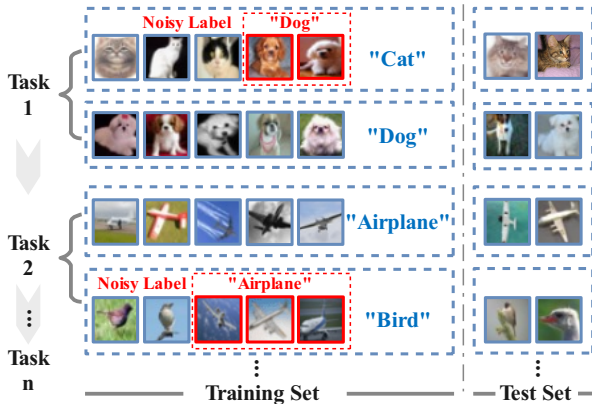


Figure 1: Diagram of the data bias.

Although current CL methods have achieved significant progress [22, 45], they heavily rely on the quality of training data [16], often overlooking more realistic settings where data biases exist in real-world applications. During data collection, due to the large scale and the high cost of expert annotation, data often contains noisy labels, as shown in Figure 1. This can cause the model to learn incorrect knowledge, making predictions unreliable [3]. This presents a greater challenge for CL, as errors can accumulate and transfer across task sequences, even amplifying [4], leading to reduced performance in both old and new tasks. This is similar to humans stubbornly retaining incorrect memories, which obstructs the acquisition of new and valuable knowledge, especially when new information conflicts with old, incorrect memories. To address this, we provide the first in-depth experimental analysis of the specific impact of noisy labels on various CL methods (see Sect. 3.1), laying the foundation for handling data biases.

Unfortunately, existing methods cannot address the dual challenges in CL mentioned above. A direct solution is to increase the training iterations for new tasks to continually mask the impact of noisy labels from old tasks. However, this approach incurs additional resource overhead and may lead to overfitting on new tasks, resulting in degraded performance on both new and old tasks (see Sect. 5.3). Another line of solutions is to handle noisy labels such as label cleaning [18, 47] and repair [20, 44]. However, these methods assume the entire dataset is available for noise purification, which is infeasible in CL where historical data cannot be revisited. Furthermore, a few replay-based CL methods were proposed to address noisy labels, which is based on contrastive self-supervised learning [16] or pre-warmed model losses [14]. Nevertheless, these approaches incur the higher training costs (see Sect. 5.3). For a detailed comparison, please refer to Sect. 2.3.

In this work, inspired by the human ability to eliminate incorrect learned knowledge through correction or forgetting [36], we propose a novel approach: achieving CL that both prevents forgetting and intentionally forgets by actively erasing already learned

errors in the model. However, combining CL with existing erasure methods presents two key challenges. **Challenge 1:** The erroneous samples in the data stream are unknown, so before applying any forgetting method, it is essential to accurately identify these biases and their corresponding memories in the model. **Challenge 2:** In CL, erasing the contribution of only a small number of noisy samples in the model can easily lead to a sharp decline in the performance of old tasks. The few existing methods focus on forgetting entire classes rather than specific samples [13], [39]. Therefore, accurately forgetting only the erroneous memories in the model, without affecting the correct knowledge, is highly challenging.

These challenges motivate us to propose the **innovative ErrorEraser plugin**, which consists of two modules: Error Identification and Error Erasure. The ErrorEraser plugin offers **four main advantages**. **First**, it effectively characterizes the training data distribution by learning the probability density in a compact feature space, without requiring preprocessing or prior knowledge. This accurately identifies potential biased samples based on distribution margins (addressing Challenge 1, see Sect. 4.1). **Second**, it efficiently forgets erroneous knowledge in the model by selectively moving the decision space of only a few representative noisy label samples. The new decision space is moved away from the decision space of normal data to avoid the loss of correct knowledge (addressing Challenge 2, see Sect. 4.2). **Third**, the incremental feature distribution learning strategy reduces resource consumption during error identification in new tasks. Additionally, selectively fine-tuning the model with representative samples is efficient and effective (see Sect. 5.3). **Fourth**, ErrorEraser is independent of specific CL learning strategies, making it compatible with various CL methods and ensuring good generalizability (see Sect. 5.2). Our main contributions are as follows:

(1) We are the first to deeply investigate the dual impact of noisy labels on CL. We analyze the intrinsic reasons for performance degradation caused by noisy labels from the perspectives of model feature extraction and prediction outcomes across three different types of CL methods.

(2) We are the first to enhance CL methods using intentional forgetting and propose a universal plugin, ErrorEraser. It effectively identifies and accurately erases potential errors already learned by the model, without the need to access historical data.

(3) We applied the ErrorEraser plugin to three types of CL methods (replay-based, regularization-based, and optimization-based). Extensive experiments on benchmark datasets validated the generality and effectiveness of ErrorEraser in forgetting erroneous knowledge caused by noisy labels.

2 RELATED WORK

2.1 Continual Learning

Continual learning (CL) encompasses task incremental learning [2, 39], class incremental learning [11, 26], and domain incremental learning [32, 43]. Existing mainstream CL methods are categorized into three types: regularization-based [1, 19], replay-based [40, 46], and model optimization-based methods [27, 35]. However, these methods assume that all learned knowledge in the model is beneficial, neglecting the errors introduced by data biases in real-world scenarios. These errors may also be amplified in downstream tasks

of a task sequence [4], further degrading the performance of CL. Although recent research has begun to address this issue by introducing filtering techniques to identify clean samples within data biases [14, 16], allowing the model to train on correct data as much as possible, fully eliminating all problematic samples remains challenging. As a result, the model continues to be affected by data biases. Furthermore, these methods are primarily applicable to replay-based CL approaches, lacking generalizability.

To address this, this paper focuses on the correction of learned erroneous knowledge in three types of CL methods under the task-incremental learning setting (where the IDs of historical tasks are known). We propose a universal plugin method that, for the first time, tackles this issue from two perspectives: learning the correct knowledge under real-time constraints and actively erasing the erroneous knowledge that has already been learned.

2.2 Machine Unlearning

Machine Unlearning (MU) is an emerging data forgetting method developed in response to growing attention on model privacy protection and data privacy regulations, such as the European Union’s General Data Protection Regulation (GDPR) [13]. The primary goal of MU is to adjust a trained model to forget information learned from specified subsets of data [9]. Recent research has largely targeted privacy preservation in deep neural networks, categorized into class-based [8, 39] and instance-based data forgetting [7, 10]. Class-based forgetting methods, more prevalent, aim to forget all instances of a specific class while retaining the performance on other classes. Instance-based forgetting seeks to forget specific instances without affecting the performance on the remaining data, presenting a greater challenge with fewer proposed methods.

Currently, a limited body of work has integrated MU into CL to address privacy protection issues [38, 39]. However, these methods assume that the data to be forgotten is known in advance, making them unsuitable for tackling the unknown data biases we aim to address, presenting an additional challenge.

Table 1: Summary of Existing Methods

Methods	CL	MU	Assumption	Where are data bias
CNLL [14]	✓	✗	-	New task
SPR [16]	✓	✗	-	New task
SCL-SSR [30]	✓	✗	-	New task
LSF [38]	✓	✓	Known data bias	Old task
LSFM [39]	✓	✓	Known data bias	Old task
ErrorEraser	✓	✓	No	Old task

2.3 Oversights of Existing Methods

Table 1 presents existing methods relevant to our work, with extensive experimental comparisons provided in Sect. 5.3. While LSF [38] and LSFM [39] integrate continual learning (CL) and machine unlearning (MU), they assume that forgetting data are known. This assumption conflicts with the privacy constraints in CL, where the forgetting data in the old task is inaccessible due to privacy

concerns. Our method, however, *is capable of identifying such biases without requiring access to the original data from the old task*. Furthermore, CNLL [14], SCL-SSR [30], and SPR [16] address issues related to data biases in the new task. In this paper, *we focus on the data biases present in the old task*, such as those arising from large-scale data and the high cost of expert annotation. Often, data contains noisy labels that lead to the model learning these biases.

3 PROBLEM STATEMENT

3.1 Problem Phenomenon: Bias-Induced CL Degradation

Data is the cornerstone of deep learning. When data contains noisy label biases, the model may learn incorrect knowledge, affecting its performance. This issue is particularly pronounced in CL. Specifically, we have identified **two key impacts** of noisy labels in data on CL: **(i) Knowledge Retention**: Noisy labels diminish the model’s ability to retain knowledge, exacerbating catastrophic forgetting (CF) of previous tasks. **(ii) Knowledge Transfer**: Noisy labels cause the model to transfer erroneous knowledge, leading to degraded learning performance on new tasks.

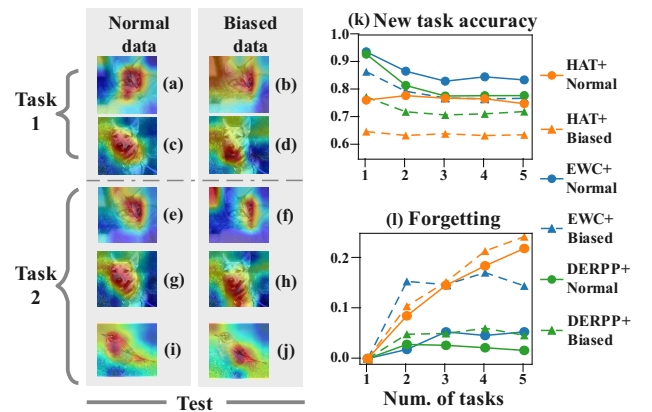


Figure 2: The impact of noisy label biases on feature extraction (left) and prediction results (right).

We conducted impact analysis experiments on CIFAR-10 with label bias and clean data, using three popular CL methods: EWC [19], DERPP [5], and HAT [37]. Specifically, we analyzed the deep reasons behind the impact of noise on CL by visualizing model feature extraction attention and using line plots of prediction results. To facilitate better analysis, we set Task 1 to include noisy bias, while keeping the data in downstream tasks clean. The attention maps highlight the areas of focus during image prediction after the model is trained, providing an intuitive way to demonstrate whether the feature extraction was correct. Figure 2 presents the experimental results.

Impact (i): We analyze the extent of CF in the model by comparing Task 1 test results during each task for models trained on normal data versus noisy biased data. The feature extraction attention maps show that, after learning Task 2, noisy bias causes a significant shift in the model’s attention to key features of Task 1, narrowing the focus (compare (f)(h) with (b)(d)). This indicates

significant changes in the parameter space of the feature extraction layer. In contrast, the model trained on normal data remains stable. Coupled with the forgetting curve (Figure 2(l)), the model trained with noisy data has higher and steeper forgetting rates. Therefore, we conclude that data biases exacerbate CF and weaken CL’s ability to retain knowledge in the parameter space. Works [28] and [48] support the idea that changes in the feature space contribute to CF.

Impact (ii): We analyze the impact on knowledge transfer by comparing the performance of models trained on two types of data on the new task. The attention maps show that the model trained on normal data primarily focuses on the animal’s facial area (Figures (a)(c)(i)), while the model trained on noisy biased data focuses more on other body parts or the background, with more dispersed attention (Figures (b)(d)(j)). This suggests that noisy bias interferes with the model’s ability to extract key features. Accuracy analysis (Figure 2(k)) further confirms this, as the model trained on noisy biased data performs significantly worse on the new task compared to the model trained on normal data. This indicates that erroneous knowledge introduced by noisy biases continuously transfers during CL, weakening the transfer of correct knowledge.

Therefore, erasing the erroneous knowledge learned by the model to address the dual negative impact of noisy biased data on CL is a challenging and urgent problem that needs to be solved.

3.2 Problem Formulation

Continual Learning: Let $\{D_1, D_2, \dots, D_T\}$ denote a sequence of datasets corresponding to T sequential tasks, where each task dataset is defined as $D_t = \{(x_t^i, y_t^i)\}_{i=1}^{n_t}$, with $x_t^i \in \mathcal{X}$ representing the input and $y_t^i \in \mathcal{Y}$ its associated class label. The goal of Continual Learning (CL) is to incrementally train a model $\mathcal{F}_{\theta_t} : \mathcal{X}_t \rightarrow \mathcal{Y}_t$ on tasks 1 through t , and then update it to $\mathcal{F}_{\theta_{t+1}}$ for task $t+1$ by leveraging knowledge from \mathcal{F}_{θ_t} , such that it can correctly predict labels for any input x encountered in tasks 1 through $t+1$.

Data Biases and Erroneous Knowledge: Data biases may randomly emerge in tasks 1 through t . Specifically, for task t , the dataset D_t contains k samples with noisy labels, denoted as $D'_t = \{x_t^i \mid 1 \leq i \leq k\}$. The true labels $\{y_t^i \mid 1 \leq i \leq k\} \in \mathcal{Y}_t$ are incorrectly annotated as $\{\hat{y}_t^i \mid 1 \leq i \leq k\}$, where $y_t^i \neq \hat{y}_t^i$. When trained on such noisy data, the model \mathcal{F}_{θ_t} may capture incorrect or conflicting representations in the feature space, which we define as *erroneous knowledge*, denoted by $\bar{\mathbb{K}}_t$.

Objective: The goal is to selectively erase the erroneous knowledge $\bar{\mathbb{K}}_t$ learned by the model \mathcal{F}_{θ_t} during tasks 1 through t , while learning task $t+1$. Simultaneously, the model should retain and transfer correct knowledge to task $t+1$, thereby enhancing the training effectiveness of $\mathcal{F}_{\theta_{t+1}}$.

Constraints: (1) The noisy labeled samples are randomly distributed and not known in advance. (2) Historical data from previous tasks is inaccessible. (3) The solution must be computationally efficient with low resource overhead.

4 PROPOSED METHOD

In this section, we introduce ErrorEraser, a universal plugin designed to forget erroneous knowledge in CL. As shown in Figure 3, its core components are Error Recognition and Error Erasure. Section 4.1 describes the process of identifying noisy labels based

on the learned data distribution. Section 4.2 introduces how to forget erroneous knowledge in the model. The implementation of ErrorEraser is in Algorithm 1.

Algorithm 1: ErrorEraser plugin.

Input : Task set $\{D_{train}^t, D_{valid}^t, D_{test}^t\}_{t=1}^T$, threshold δ ;
Output : Trained parameters θ and prediction model \mathcal{F}_{θ_t} .

```

1 Initialize  $\theta$ ;
2 for  $task = 1$  to  $T$  do
3   // Error Recognition
4   for each training epoch do
5     Calculate the feature space  $h_b(D_{train}^t)$  of  $\mathcal{F}_{\theta_t}$ ;
6     Construct the distribution learning model  $\mathcal{M}$  on
        $h_b(D_{train}^t)$  by Eq. (6);
7     Train  $\mathcal{F}_{\theta_t}$  on  $\{D_{train}^t, D_{valid}^t\}$  by loss  $\mathcal{L}$ ;
8   end
9   Calculate the probability density  $p$  value of  $D_{train}^t$  on  $\mathcal{M}$ 
       by Eq. (2);
10  // Error Erasure
11  Select representative sample  $D_s$  according to  $p$  by Eq. (11);
12  for each forgetting epoch do
13    Fine-tune the model  $\mathcal{F}_{\theta_t}$  on  $D_s$  by Eq. (12);
14  end
15 end
```

4.1 Error Recognition

This subsection focuses on identifying potential noisy labels in the data stream to forget erroneous knowledge in the model. It begins by constructing a continual distribution learning model using an incremental feature probability density strategy to conserve training resources. A new loss function is then used to collaboratively train the backbone classification model and the distribution learning model, enabling accurate data distribution learning and reducing the impact of noisy labels in the current task. Finally, potential noisy label samples are identified based on the learned distribution.

4.1.1 Incremental Feature Probability Density Distribution Learning.

We begin by proposing an incremental feature probability density learning strategy. This strategy transforms complex data distributions in a compact feature space into a specified multivariate normal distribution in the latent space through a series of function transformations. Assume the structure of the classification network model is $h = \{h_1, \dots, h_n\}$, where $h_a = \{h_1, \dots, h_l\}$ represents continuous feature extraction layers such as convolutional and pooling layers. $h_b = \{h_{l+1}, \dots, h_{n-1}\}$ represents intermediate continuous feature extraction layers, such as linear layers. h_c denotes the shared feature layer connected to the classification head, indicating the extraction of core features.

We first introduce the Normalizing Flows (NF) model [34], which can learn the probability density of data through distribution transformations. However, directly learning the distribution on raw data poses two challenges: first, the high-dimensional nature of the data space makes the learning process time-consuming and

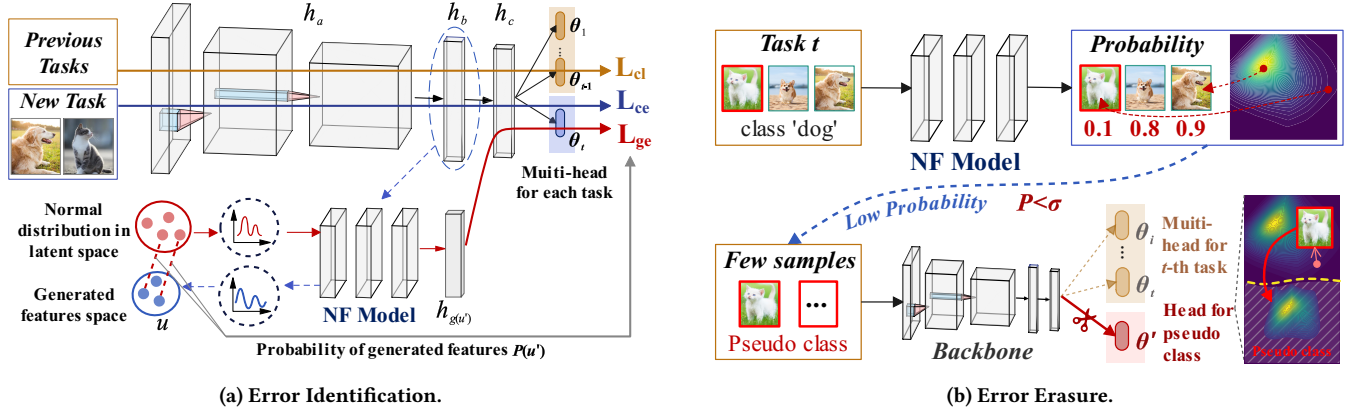


Figure 3: Illustrations of our ErrorEraser plugin for classifier training in CL.

challenging; second, the sparsity of raw data makes it difficult to capture its distribution complexity, leading to overfitting on limited training samples and poor generalization to the overall data distribution. Therefore, we use the compact, low-dimensional feature space within the classification network model structure as input to construct the NF model. This approach not only accelerates the NF modeling process but also enables the model to better capture the distribution characteristics of the data.

Here, we use the feature space extracted by h_b as the learning object. We first transform the real feature space in h_b into a simpler prior data distribution, such as a multivariate normal distribution or a standard Gaussian distribution, through a reversible smooth mapping $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$, with the inverse mapping $f^{-1} = g$. Assume the prior distribution is $u \sim p_u(u)$, and the data for the t -th task is $D_t = \{(x_t^i, y_t^i)_{i=1}^{n_t}\}$, with the corresponding $h_b(x) \sim p_{h_b(x)}(h_b(x))$. The distribution of the random variable $h_b(x)$ is calculated as follows:

$$p_u(u) = p_{h_b(x)}(h_b(x)) \left| \det J_f(g(h_b(x))) \right|, \quad (1)$$

where J_f is the Jacobian matrix of the transformation f .

Then, a complex distribution can be constructed through a series of reversible mappings $u = f_k \circ \dots \circ f_2 \circ f_1(h_b(x))$. At this point, using this series of reversible transformations, the probability density of $h_b(x)$ in the prior distribution $p_u(u)$ can be calculated as follows:

$$\log p(x) = \log p_u(u) + \log \left| \det \frac{\partial u}{\partial x} \right| \quad (2)$$

$$= \log p_u(f(x)) + \sum_{l=1}^{k-1} \log \left| \det \frac{\partial f^{l+1}}{\partial f^l} \right|, \quad (3)$$

where $u = f(h_b(x))$ is the mapping learned through the parameters θ that transforms the distribution $p_{h_b(x)}$. f^l denotes the input of the l -th transformation in the series of mappings.

Additionally, we use the label y as conditional information to calculate the likelihood estimate $p_{h_b(x)}(h_b(x)|y)$. The objective of our modeling series mapping function is to maximize the likelihood of samples from the data D_t . The loss function of the current task

distribution learning is as follows:

$$\mathcal{L}_p(f; h_b(x)) = -\frac{1}{n} \sum_{i=1}^{n_t} \log p_{h_b(x_i)}(h_b(x_i)|y_t^i) \quad (4)$$

$$= -\frac{1}{n} \sum_{i=1}^n \left(\log p_u(f(h_b(x_i))|y_i) + \sum_{l=1}^{k-1} \log \left| \det \frac{\partial f_i^{l+1}}{\partial f_i^l} \right| \right). \quad (5)$$

To simplify the distribution learning process in CL and avoid learning an independent distribution mapping for each task, we designed a strategy for incrementally learning the feature space distribution. This strategy uses the feature space distributions from previous tasks as constraints to guide the current task's feature space distribution learning. Thus, in CL, the improved distribution learning modeling loss based on Eq. (4) is as follows:

$$\mathcal{L}_p(f; D_t) = -\frac{1}{|D_t|} \sum_{x_t^i, y_t^i \sim D_t} \log p_x(h_b(x_t^i)|y_t^i) \quad (6)$$

$$- \frac{1}{|G_{h_b'}|} \sum_{h_b', y_j' \sim G_{h_b'}} \log p_{h_b'}(h_b'|y_j'), \quad (7)$$

where y_j represents the class labels from tasks 1 to $t-1$. G_{h_b} represents the feature set of previous tasks sampled according to the condition y_j in the distribution learning model, which has the same shape as the h_b layer. This method allows for continuous learning of the feature space distributions for multiple tasks and prevents forgetting the feature distributions learned in old tasks.

4.1.2 Distribution Learning Optimization. Here, the focus is on obtaining high-quality feature representations from data containing noisy labels to input into the distribution learning model for a more accurate description of data distribution. We achieve this by designing a new loss function \mathcal{L} , which collaboratively trains the backbone classification model and the distribution learning model, allowing for real-time adjustment of input features. Specifically, the distribution learning model constrains the classifier's training based on the learned probability density. The updated features from the classifier are then used as input for a new round of distribution learning.

This loss function \mathcal{L} comprises three components: classification loss \mathcal{L}_{ce} , continual learning loss \mathcal{L}_{cl} , and distribution learning loss \mathcal{L}_{ge} , as shown in Figure 3. The classification loss \mathcal{L}_{ce} focuses on learning new task knowledge, the continual learning loss \mathcal{L}_{cl} aims to prevent catastrophic forgetting of previously acquired knowledge, and the distribution learning loss \mathcal{L}_{ge} seeks to mitigate the impact of noisy labels. The combined loss function \mathcal{L} is defined as follows:

$$\mathcal{L} = \underbrace{\mathcal{L}_{ce} + \mathcal{L}_{ge}}_{\text{new task}} + \underbrace{\mathcal{L}_{cl}}_{\text{previous tasks}}. \quad (8)$$

For the distribution learning loss \mathcal{L}_{ge} (\mathcal{L}_{ce} and \mathcal{L}_{cl} , see Sect. 4.2.3), the key lies in assessing the accuracy of the feature space within the classification model. Considering that the outlying features relative to the current task might be unreliable, we quantify the correlation between the features u' generated by the NF in the latent space and the prior distribution features u . For ease of computation, the prior distribution in the NF adopts a multivariate normal distribution. Finally, we measure the correlation by calculating probabilities and design the distribution learning loss \mathcal{L}_{ge} accordingly. We map the data x to an NF latent space, resulting in $U_t = \{u_i = f(h_b(x_i))\}_{i=1}^{n_t}$, and sample a subset of generated features $U' = \{u'_i\}_{i=1}^{n'}$, then compute the correlation using the following formula:

$$q_{D_t}(u'_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_t|^{1/2}} \exp\left(-\frac{1}{2}(u'_i - \mu_t)^T \Sigma_t^{-1} (u'_i - \mu_t)\right), \quad (9)$$

where $\mu_t = \frac{1}{|D_t|} \sum_{i=1}^{|D_t|} u_i$ represents the mean vector, and $\Sigma_t = \frac{1}{|D_t|} \sum_{i=1}^{|D_t|} \text{diag}(u_i - \mu_t) \cdot \text{diag}(u_i - \mu_t)$ denotes the covariance matrix. The $\text{diag}(\cdot)$ turns the vector into a diagonal matrix. Finally, the distribution learning loss \mathcal{L}_{ge} is defined as follows:

$$\mathcal{L}_{ge}(h; U') = \frac{1}{n'} \sum_{u' \in U'} q_{D_t}(u') \mathcal{L}_{ce}(h_b(f^{-1}(u')), y'_i). \quad (10)$$

4.1.3 Distribution Analysis. Considering that abnormal features relative to the current task may be caused by noisy label samples, we evaluate the probability density value of each sample's corresponding features by Eq. (2) after completing the distribution learning model. The likelihood of a sample being a noisy label is measured by the probability density value. Finally, samples with low probability density values at the distribution margins are detected, enabling the identification of potential noisy labels.

4.2 Error Erasure

The error erasure of erroneous knowledge is illustrated in Figure 3(b). The core of this process involves altering the decision space of representative noisy label samples.

4.2.1 Representative Sample Selection. After completing the t -th task, any data point $x \in D_t$ is mapped to the NF latent space distribution through the feature extraction of network layers $h_{a,b}$, and the corresponding probability $p(x)$ is obtained using Eq. (2) and (3). A high probability indicates that the data is more likely to be a correct sample, while a low probability suggests the sample is less reliable and most likely a noisy labeled sample. As shown in Figure 3(b), we visualize the distribution of the original data in the NF latent

space. The central area highlighted in yellow represents the high-probability density region, consisting predominantly of correct samples. Conversely, the edges of the distribution correspond to the low-probability regions, where the sample feature space is far from the normal area, and can be considered as noisy labeled and outlier samples.

We consider the subset of samples with the lowest probability density, D_s , as the representative samples for erroneous knowledge forgetting. It is defined as follows:

$$D_s = \{(x_i, y_i) \in D_t \mid p(x_i) < \delta\}, \quad (11)$$

where δ is a predefined threshold that can be customized based on the data volume and other factors. As δ increases, more samples are included in the forgetting process. It is important to note that the selected subset of samples should be much smaller than the current task data (i.e., $D_s \ll D_t$) to avoid excessive overhead. In our implementation, δ is set based on a density $p(x)$ percentile.

4.2.2 Decision Boundary Shift. We disperse the activation of D_s in the model through boundary expansion, remapping, and model pruning to forget erroneous knowledge while retaining correct knowledge. We intentionally assign these low-probability data D_s to an additional pseudo-class in the classification model, moving its decision space to a new area distinct from the model's original decision boundaries, as illustrated in Figure 3(b). As analyzed earlier, since low-probability samples are at the margins of data distribution, altering their decision space does not affect the classification of regular samples.

Specifically, we assign a pseudo-label \hat{y} to the selected D_s and form new training pairs (x_i, \hat{y}_i) . First, we expand the model's decision space by adding a new neuron to the original classification model's output layer. We then fine-tune the model with these new data pairs, mapping the decision space of these likely noisy label samples to the pseudo-class. The loss of model fine-tuning is:

$$\theta'_t = \arg \min_{\theta} \sum_{(x_i, \hat{y}_i) \in D_s} \mathcal{L}_{ce}(x_i, \hat{y}_i, \theta_t), \quad (12)$$

where θ'_t represents the fine-tuned model parameters.

During fine-tuning, we use the new output neuron to collect activations of noisy samples, recording and analyzing their behavior in the new neuron. Following this, we perform pruning operations to remove the new neuron and decision space associated with the noisy samples, actively forgetting the erroneous knowledge in the model. After pruning, the classification model's size remains consistent with the original, maintaining the same efficiency. The constraint of the our entire forgetting process is:

$$\mathcal{F}_{\theta'_t}(x_i) \approx 0. \quad (13)$$

Notably, the forgetting process is applied only to D_s and does not allocate any samples from existing classes, so the activations of the remaining data $D_t - D_s$ do not undergo significant changes. Consequently, the model can actively forget incorrect knowledge without affecting the correct classification performance of the remaining data, thereby enhancing overall model performance.

4.2.3 Continual Learning for New Tasks. After the forgetting process is completed, the learning of a new task begins. The learning loss remains as \mathcal{L} in Eq. (8). The \mathcal{L}_{ce} and \mathcal{L}_{cl} are not limited to

Table 2: Results on CIFAR-10, CIFAR-100, and MNIST. Label noise is 50%. Bold indicates that the results of applying the ErrorEraser plugin outperform the original CL method on noisy labels.

	Methods	CIFAR10				CIFAR100				MNIST			
		#Taks:5, #Class:10				#Taks:5, #Class:100				#Taks:5, #Class:10			
		S ↑	A ₁ ↑	A ₂ ↑	F ↓	S ↑	A ₁ ↑	A ₂ ↑	F ↓	S ↑	A ₁ ↑	A ₂ ↑	F ↓
Normal data	EWC	57.90	71.15	74.88	14.51	26.93	32.50	32.97	1.56	81.72	98.25	98.63	1.46
	MAS	58.47	71.56	74.74	12.02	19.85	24.33	25.13	4.02	40.05	49.42	53.70	13.37
	LWF	61.12	74.29	77.40	9.06	22.43	27.32	28.09	3.32	78.96	95.14	98.29	5.56
	DERPP	63.65	77.63	80.17	10.37	16.76	21.66	28.07	19.36	81.34	97.83	98.78	2.27
	HAT	56.77	69.62	74.78	14.96	20.81	25.40	26.88	4.18	80.72	97.09	98.77	3.05
	TAT	51.89	64.26	70.75	20.01	13.56	18.55	29.48	35.68	71.22	86.26	88.19	6.82
Noisy labeled	EWC	53.28	65.12	69.49	12.07	26.89	32.48	32.78	1.56	71.84	87.87	90.99	13.65
	MAS	54.55	66.28	69.65	8.59	19.48	23.84	24.83	3.96	69.57	84.33	87.95	8.97
	LWF	51.68	63.16	68.61	13.00	20.20	24.77	25.87	4.45	72.68	88.58	90.11	9.94
	DERPP	47.12	58.79	68.39	25.9	15.82	20.71	26.13	19.72	71.97	87.64	90.56	11.00
	HAT	45.87	55.42	54.55	1.34	12.40	19.36	20.75	38.90	56.37	67.77	68.29	1.29
	TAT	50.72	61.98	64.98	10.13	12.97	17.64	28.93	34.45	70.85	85.72	86.44	4.99
Noisy labeled	EWC+ErrorEraser	64.78	78.30	78.31	3.42	39.30	48.37	51.38	10.86	74.26	90.26	93.87	10.84
	MAS+ErrorEraser	61.27	75.15	76.42	11.40	28.03	34.72	37.53	10.06	72.26	87.13	87.59	3.00
	LWF+ErrorEraser	58.71	71.72	72.67	8.83	24.19	29.65	32.56	7.05	80.52	96.93	97.12	1.98
	DERPP+ErrorEraser	56.62	69.66	68.57	9.26	16.58	21.07	27.72	16.58	73.97	88.91	88.29	0.29
	HAT+ErrorEraser	50.16	63.35	73.44	33.23	15.86	21.50	33.04	37.28	59.35	70.58	82.41	7.34
	TAT+ErrorEraser	54.52	65.86	66.43	3.16	13.82	18.81	29.01	33.58	72.61	87.35	88.14	2.10

specific loss functions but depend on the CL method employed. Specifically, common \mathcal{L}_{ce} in CL methods include cross-entropy loss and mean squared error loss, among others. The \mathcal{L}_{cl} can correspond to regularization loss, replay loss, or parameter isolation loss in CL methods. Our focus is on designing \mathcal{L}_{ge} to be compatible with any \mathcal{L}_{ce} and \mathcal{L}_{cl} , constraining the model to learn the correctly labeled samples and their probability density of the current task data as accurately as possible. \mathcal{L}_{ce} and \mathcal{L}_{cl} are as follows:

$$\mathcal{L}_{ce}(x; D_t) = \frac{1}{|D_t|} \sum_{i=1}^{|D_t|} \mathcal{L}_{ce}(x_i, y_i), \quad (14)$$

$$\mathcal{L}_{cl}(x; D_t) = \frac{1}{|D_t|} \sum_{i=1}^{|D_t|} \sum_{k=1}^{t-1} \Omega_k \mathcal{L}_k(x_i, y_i), \quad (15)$$

where Ω_k is a weight factor that determines the importance of retaining the previous knowledge of the k -th task, and \mathcal{L}_k represents the loss associated with the k -th task. This loss can be based on various methods, such as the regularization loss in EWC: $\mathcal{L}_k = \frac{\lambda}{2} \sum_j (\theta_j - \theta_j^k)^2$, where λ is a hyperparameter controlling the strength of the regularization, θ_j are the current parameters, and θ_j^k are the parameters learned from the k -th task.

5 EXPERIMENTS

5.1 Setting

Datasets and Implementation Details: We conducted experiments on two types of noisy label data: synthetic (MNIST [23], CIFAR-10 [21], CIFAR-100 [21]) and real (WebVision [41]). For the synthetic setup, following [39], we partitioned dataset classes into tasks to simulate CL with 5 and 10 tasks, applying asymmetric noise labeling to specific tasks [24]. For WebVision with a 20% noise rate, we divided it into 5 tasks. A Convolutional Neural Network (CNN)

[23] was used as the classification model for all tasks, with the final layer modified to a multi-head architecture, as shown in Figure 3.

Baselines: We selected three representative CL methods to validate the effectiveness of our plugin: regularization-based methods (EWC [19] and MAS [1]), replay-based methods (LWF [27] and DERPP [5]), and optimization-based methods (TAT [29] and HAT [37]). Additionally, we compare our method with state-of-the-art data forgetting approaches LSFM [39] and ADV+IMP [7], as well as noisy label CL methods CNLL [14] and SPR [16].

Evaluation Metric: In the ErrorEraser setting, our goal is to forget erroneous knowledge to enhance CL performance. Traditional metrics do not adequately evaluate performance in this context, as they overlook the combined retention and transfer capabilities of knowledge. Following [39], we introduce a new metric, S , which is the harmonic mean of three standard CL metrics: average accuracy on new tasks (A_1), average historical accuracy (A_2), and average forgetting rate (F_k). This metric provides a comprehensive measure of model effectiveness in a CL framework:

$$S = \frac{A_1 * A_2 * n}{A_1 + (A_2 * n) + F_k}, \quad (16)$$

where n is the total number of tasks, $A_1 = \frac{1}{n} \sum_{k=1}^n a_{k,k}$, $A_2 = \frac{1}{n} \sum_{k=1}^n a_{k,n}$, $F_k = \frac{1}{n} \sum_{k=1}^n f_k^n$, with $f_k^n = \frac{1}{(n-k)} \sum_{i=k}^n (a_{k,i} - a_{k,n})$, and $a_{k,n}$ denoting the accuracy of task k after learning task n . The values of A_1 , A_2 , and F_k all range between $[0,1]$. We randomly shuffled the task order and conducted experiments five times, reporting the overall average for each metric.

5.2 Effectiveness Guarantee

Overall performance: Table 2 and Table 3 present performance comparisons on synthetic and real data, respectively. Table 1 shows the results for all methods on both clean and noisy label data. It

is clear that, regardless of the CL method used, including the ErrorEraser plugin significantly improves the model’s performance in terms of accuracy and forgetting rate. In some cases, the results even surpass those obtained with clean data. Notably, (1) some CL methods exhibit lower forgetting on noisy data compared to clean data. This is because, as defined by the forgetting metric F , these methods produce similar accuracy at learning time and at final evaluation, resulting in small gaps and thus artificially low forgetting scores. (2) On certain datasets, even with our plugin applied, the HAT-based method shows slightly higher forgetting. This is because, although our plugin improves accuracy, HAT’s rigid task-specific masking limits backward transfer, which can amplify forgetting. Furthermore, the results in Table 3 further demonstrate that the ErrorEraser plugin effectively addresses the erroneous knowledge introduced by real data. This improvement is mainly due to the two modules in our plugin, which are closely aligned with the two goals of CL. The error identification module helps discard conflicting or irrelevant knowledge during the knowledge retention process, while the error erasure module further eliminates potential internal biases in the model. The appendix provides detailed experimental results for different noise label rates, task numbers, and the number of categories within tasks.

Table 3: Results on the real-world noisy dataset WebVision

Methods	$S \uparrow$	$A_1 \uparrow$	$A_2 \uparrow$	$F \downarrow$
EWC	42.69	52.08	50.70	3.70
LWF	41.05	50.21	48.36	3.69
HAT	41.86	50.77	50.93	3.36
EWC+ErrorEraser	43.72	53.44	51.25	3.54
LWF+ErrorEraser	42.02	51.38	48.70	2.89
HAT+ErrorEraser	42.75	51.92	51.21	3.01

Visualization of Decision Space: To better understand how decision boundaries change with noisy label data, we visualized the decision boundaries for each task and historical task 1 after each task in Figure 4. With the application of the ErrorEraser plugin, we observed that the boundaries between the two classes of Task 1 data remained clear after each task, illustrating that the Error Identification module effectively prevents catastrophic forgetting. Additionally, the decision boundaries for new tasks remained compact, indicating that the Error Erasure module further eliminates erroneous knowledge, facilitating the transfer of correct knowledge to guide the learning of new tasks.

Attention Map: To better illustrate the impact of noisy labels, we visualized attention maps of CL models trained on the CIFAR-10 dataset with noisy labels, both with and without the ErrorEraser plugin. The highlighted regions indicate key areas for concept prediction. Each column in Figure 5 represents different CL methods, and each row shows tests after completing various tasks. We observed that after completing Task 1, methods like EWC and DERPP focused heavily on background information, with attention still misaligned from facial features. After Task 2, attention on Task 1 test data (cat) sharply decreased, focusing only on minimal facial regions. For Task 2 test data (bird), the models failed to capture comprehensive facial information. In contrast, with ErrorEraser,

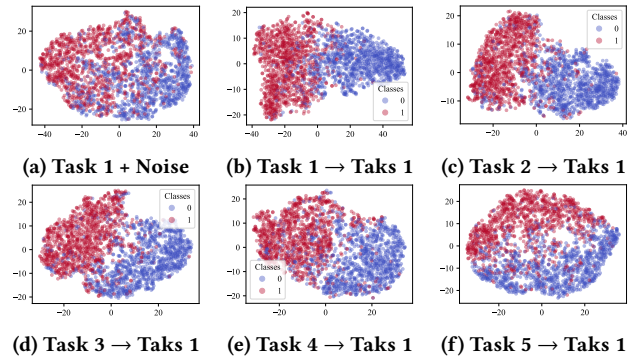


Figure 4: t-SNE visualization of feature extraction for Task 1 (training data contains 50% noisy labels) in the CIFAR-10 test dataset. (a) shows feature extraction for Task 1 using the original EWC model. (b)–(f) show feature extraction for Task 1 after completing Tasks 1 to 5 using ErrorEraser.

the models consistently focused on the core facial regions after completing both tasks. This further validates the effectiveness of the ErrorEraser plugin in erasing erroneous knowledge.

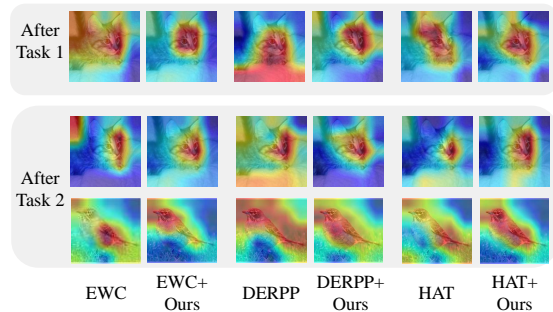


Figure 5: The attention map of three CL methods before and after applying ErrorEraser on noisy labels.

5.3 Comparison with other methods

To evaluate the effectiveness and efficiency of our method, we also compared it with two traditional noise-handling CL methods, CNLL and SPR, which focus on constraining the model’s learning of errors, while we emphasize correcting errors after the model has learned them. In particular, we counted the running time, memory overhead, and model size of each method for training a single task. Furthermore, we also compared other forgetting strategies.

Effectiveness: As shown in Table 4, regarding data forgetting methods, LSFM targets entire classes rather than specific samples, making it less effective at selectively forgetting erroneous knowledge. The ADV+IMP method uses adversarial perturbations for sample-level forgetting, but this negatively impacts the performance of normal samples.

Compared to CNLL and SPR, our method outperforms in preventing catastrophic forgetting and enhancing new task learning. CNLL and SPR mainly focus on mitigating catastrophic forgetting in old

Table 4: Results of other data forgetting methods and noisy label CL methods. Label noise is 50%.

Type	Methods	$S \uparrow$	$A_1 \uparrow$	$A_2 \uparrow$	$F \downarrow$
Data forgetting	LSFM	28.21	36.54	44.62	29.30
	ADV+IMP	39.42	48.95	59.72	23.21
Noisy label	CNLL	57.66	71.44	71.95	14.47
	SPR	58.41	71.49	73.86	11.20
Ours	ErrorEraser	64.78	78.31	78.30	3.42

tasks but overlook the impact of noisy labels on new tasks. Specifically, CNLL’s approach of differentiating noisy labels based on loss from a warm-up model is insufficient. Additionally, CNLL and SPR, being replay-based methods, lack the generalizability of our approach. This further underscores the innovation of our method in correcting the erroneous knowledge already learned by the model.

Efficiency: The resource overhead of these methods is shown in Table 5. “Original” represents the baseline CL method without additional measures. “Cover” refers to the method of repeatedly training new task data (three times) to mitigate the impact of noisy labels. The results show that our method is closest to “Original” in terms of time, memory, and model size. Specifically, “Cover” consumes three times the resources of “Original,” proportional to the number of additional training iterations. CNLL and SPR show significant overhead in time and memory, with SPR being particularly resource-intensive due to its extensive data augmentation for contrastive self-supervised learning. Thus, our plugin is more efficient.

Table 5: Resource cost comparison

Methods	Time (s)	Memory (MB)	Model Size (MB)
Original	$\sim 4 \times 10^1$	$\sim 3 \times 10^2$	41.7
Cover	$\sim 1.2 \times 10^2$	$\sim 9 \times 10^2$	125.10
CNLL	$\sim 7 \times 10^2$	$\sim 2 \times 10^3$	42.9
SPR	$\sim 2 \times 10^4$	$\sim 5 \times 10^3$	92
ErrorEraser	$\sim 7 \times 10^1$	$\sim 4 \times 10^2$	51.7

5.4 Ablation Study

We conducted an ablation experiment using EWC as the representative method under a 30% noisy label setting, with results shown in Table 6. The findings indicate that the complete ErrorEraser plugin performs best, highlighting the contribution of each component to the method’s overall effectiveness. First, the \mathcal{L}_{ge} loss in Module 1 helps CL efficiently learn the data distribution during each task and restricts training to high-probability density samples, reducing the impact of noisy labels. Module 2’s active forgetting strategy further improves performance by guiding new task learning and preventing catastrophic forgetting of old tasks. This validates the effectiveness of our method in forgetting erroneous knowledge.

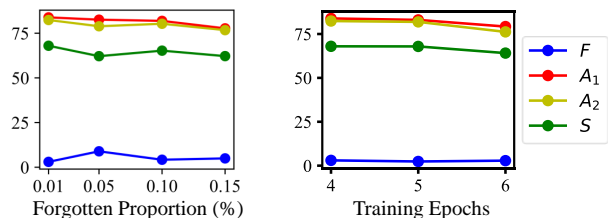
5.5 Sensitivity Analysis

We analyzed the sensitivity of hyperparameters, as shown in Figure 6, including the amount of data used for forgetting (i.e., the number of samples selected by Eq.(11), which correspond to the lower

Table 6: Performance comparison with different component combinations. Module 1 is error identification, and Module 2 is error erasure.

Components				CIFAR10			
Module 1			Module 2	$S \uparrow$	$A_1 \uparrow$	$A_2 \uparrow$	$F \downarrow$
\mathcal{L}_{ce}	\mathcal{L}_{ge}	\mathcal{L}_{cl}					
✓				65.04	86.17	62.85	15.91
✓	✓			66.31	86.94	64.73	13.72
✓	✓		✓	65.56	85.43	65.29	13.48
✓		✓		66.14	81.69	80.00	12.34
✓	✓	✓		68.17	82.67	81.64	4.09
✓	✓	✓	✓	69.25	83.86	82.39	3.01

percentile of the probability density distribution) and the number of fine-tuning epochs. As the proportion of data used for forgetting increases, overall performance slightly decreases, as the likelihood of including normal samples increases, which may result in the erasure of correct knowledge. Similarly, increasing the number of fine-tuning epochs causes a small performance drop, as excessive epochs lead to overfitting on noisy labels. Overall, the impact of these parameters is minimal, demonstrating that the ErrorEraser plugin is robust to parameter sensitivity.

**Figure 6: Performance on different parameters.**

6 CONCLUSION

This paper explores the dual impact of noisy labels on CL and introduces ErrorEraser, a universal plugin designed to address this issue. ErrorEraser learns the probability density distribution of data to effectively identify potential noisy labels in the task stream. By selectively adjusting the decision space of representative noisy samples, it accurately erases erroneous knowledge in the model, mitigating the impact of noisy labels on knowledge retention and transfer in CL. Additionally, we propose an incremental data distribution learning strategy to continuously capture each task’s data distribution, optimizing resource usage. We apply ErrorEraser to three popular CL methods and conduct extensive experiments on benchmark datasets, demonstrating its effectiveness, efficiency, and generality.

REFERENCES

- [1] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. 2018. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European conference on computer vision*. 139–154.
- [2] Rahaf Aljundi, Klaas Kelchtermans, and Tinne Tuytelaars. 2019. Task-free continual learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 11254–11263.

- [3] Nourhan Bayasi, Jamil Fayyad, Alceu Bissoto, Ghassan Hamarneh, and Rafeef Garbi. 2024. Biaspruner: Debaised continual learning for medical image classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 90–101.
- [4] Florian Peter Busch, Roshni Kamath, Rupert Mitchell, Wolfgang Stammer, Kristian Kersting, and Martin Mundt. 2024. Where is the Truth? The Risk of Getting Confounded in a Continual World. *arXiv preprint arXiv:2402.06434* (2024).
- [5] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. 2020. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems* 33 (2020), 15920–15930.
- [6] Xuemei Cao, Xin Yang, Shuyin Xia, Guoyin Wang, and Tianrui Li. 2024. Open Continual Feature Selection via Granular-Ball Knowledge Transfer. *IEEE Transactions on Knowledge and Data Engineering* 36, 12 (2024), 8967–8980.
- [7] Sungmin Cha, Sungjun Cho, Dasol Hwang, Honglak Lee, Taesup Moon, and Moontae Lee. 2024. Learning to unlearn: Instance-wise unlearning for pre-trained classifiers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 11186–11194.
- [8] Min Chen, Weizhuo Gao, Gaoyang Liu, Kai Peng, and Chen Wang. 2023. Boundary unlearning: Rapid forgetting of deep networks via shifting the decision boundary. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7766–7775.
- [9] Tao Fan, Hanlin Gu, Xuemei Cao, et al. 2025. Ten Challenging Problems in Federated Foundation Models. *IEEE Transactions on Knowledge and Data Engineering* (2025), 1–20.
- [10] Hanlin Gu, Gongxi Zhu, Jie Zhang, Xinyuan Zhao, Yuxing Han, Lixin Fan, and Qiang Yang. 2024. Unlearning during Learning: An Efficient Federated Machine Unlearning Method. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*. 4035–4043.
- [11] Yiduo Guo, Bing Liu, and Dongyan Zhao. 2023. Dealing with cross-task class discrimination in online continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11878–11887.
- [12] Jindong Han, Hao Liu, Shui Liu, Xi Chen, Naiqiang Tan, Hua Chai, and Hui Xiong. 2023. iETA: A Robust and Scalable Incremental Learning Framework for Time-of-Arrival Estimation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4100–4111.
- [13] Alvin Heng and Harold Soh. 2023. Selective amnesia: A continual learning approach to forgetting in deep generative models. *Advances in Neural Information Processing Systems* 36 (2023), 17170–17194.
- [14] Nazmul Karim, Umar Khalid, Ashkan Esmaeili, and Nazanin Rahnavard. 2022. Cnll: A semi-supervised approach for continual noisy label learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3878–3888.
- [15] Zixuan Ke, Bing Liu, Nianzu Ma, Hu Xu, and Lei Shu. 2021. Achieving forgetting prevention and knowledge transfer in continual learning. *Advances in Neural Information Processing Systems* 34 (2021), 22443–22456.
- [16] Chris Dongjoo Kim, Jinseo Jeong, Sangwoo Moon, and Gunhee Kim. 2021. Continual learning on noisy data streams via self-purified replay. In *Proceedings of the IEEE/CVF international conference on computer vision*. 537–547.
- [17] Sein Kim, Namkyeong Lee, Donghyun Kim, Minchul Yang, and Chanyoung Park. 2023. Task Relation-aware Continual User Representation Learning. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 1107–1119.
- [18] Youngdong Kim, Junho Yim, Juseung Yun, and Junmo Kim. 2019. Nlnl: Negative learning for noisy labels. In *Proceedings of the IEEE/CVF international conference on computer vision*. 101–110.
- [19] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences* 114, 13 (2017), 3521–3526.
- [20] Jan Kremer, Fei Sha, and Christian Igel. 2018. Robust active label correction. In *International conference on artificial intelligence and statistics*. 308–316.
- [21] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. *Technical report* (2009).
- [22] Guannan Lai, Yujie Li, Xiangkun Wang, Junbo Zhang, Tianrui Li, and Xin Yang. 2025. Order-robust class incremental learning: Graph-driven dynamic similarity grouping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [23] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [24] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. 2019. Learning to learn from noisy labeled data. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5051–5059.
- [25] Miaomiao Li, Jiaqi Zhu, Xin Yang, Yi Yang, Qiang Gao, and Hongan Wang. 2023. Cl-wstc: Continual learning for weakly supervised text classification on the internet. In *Proceedings of the ACM Web Conference 2023*. 1489–1499.
- [26] Yujie Li, Xin Yang, Hao Wang, Xiangkun Wang, and Tianrui Li. 2024. Learning to Prompt Knowledge Transfer for Open-World Continual Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 13700–13708.
- [27] Zhizhong Li and Derek Hoiem. 2018. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence* 40, 12 (2018), 2935–2947.
- [28] Sen Lin, Peizhong Ju, Yingbin Liang, and Ness Shroff. 2023. Theory on forgetting and generalization of continual learning. In *International Conference on Machine Learning*. 21078–21100.
- [29] Sihao Lin, Hongwei Xie, Bing Wang, Kaicheng Yu, Xiaojun Chang, Xiaodan Liang, and Gang Wang. 2022. Knowledge distillation via the target-aware transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10915–10924.
- [30] Yiwei Luo and Min Jiang. 2023. Solving Continual Learning with Noisy Labels by Sample Selection and Replay. In *2023 International Joint Conference on Neural Networks*. 1–8.
- [31] Imad Eddine Marouf, Subhankar Roy, Enzo Tartaglione, and Stéphane Lathuilière. 2024. Weighted ensemble models are strong continual learners. In *European Conference on Computer Vision*. 306–324.
- [32] M Jehanzeb Mirza, Marc Masana, Horst Possegger, and Horst Bischof. 2022. An efficient domain-incremental learning approach to drive in all weather conditions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 3001–3011.
- [33] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. 2017. Making deep neural networks robust to label noise: A loss correction approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1944–1952.
- [34] Danilo Rezende and Shakir Mohamed. 2015. Variational inference with normalizing flows. In *International conference on machine learning*. 1530–1538.
- [35] Tim GJ Rudner, Freddie Bickford Smith, Qixuan Feng, Yee Whye Teh, and Yarín Gal. 2022. Continual learning via sequential function-space variational inference. In *International Conference on Machine Learning*. 18871–18887.
- [36] Tomás J Ryan and Paul W Frankland. 2022. Forgetting as a form of adaptive engram cell plasticity. *Nature Reviews Neuroscience* 23, 3 (2022), 173–186.
- [37] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. 2018. Overcoming catastrophic forgetting with hard attention to the task. In *International conference on machine learning*. 4548–4557.
- [38] Lianlei Shan, Wenzhang Zhou, Wei Li, and Xingyu Ding. 2024. Lifelong Learning and Selective Forgetting via Contrastive Strategy. *arXiv preprint arXiv:2405.18663* (2024).
- [39] Takashi Shibata, Go Irie, Daiki Ikami, and Yu Mitsuzumi. 2021. Learning with Selective Forgetting. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence*. 989–996.
- [40] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. 2017. Continual learning with deep generative replay. *Advances in neural information processing systems* 30 (2017).
- [41] Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. 2022. Learning from noisy labels with deep neural networks: A survey. *IEEE transactions on neural networks and learning systems* 34, 11 (2022), 8135–8153.
- [42] Zhicheng Sun, Yadong Mu, and Gang Hua. 2023. Regularizing second-order influences for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 20166–20175.
- [43] Gido M Van de Ven, Tinne Tuytelaars, and Andreas S Tolias. 2022. Three types of incremental learning. *Nature Machine Intelligence* 4, 12 (2022), 1185–1197.
- [44] Andreas Veit, Neil Alldrin, Gal Chechik, Ivan Krasin, Abhinav Gupta, and Serge Belongie. 2017. Learning from noisy large-scale datasets with minimal supervision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 839–847.
- [45] Xin Yang, Hao Yu, Xin Gao, Hao Wang, Junbo Zhang, and Tianrui Li. 2024. Federated continual learning via knowledge fusion: A survey. *IEEE Transactions on Knowledge and Data Engineering* 36, 8 (2024), 3832–3850.
- [46] Jaehong Yoon, Divyam Madaan, Eunho Yang, and Sung Ju Hwang. 2021. Online coreset selection for rehearsal-based continual learning. *arXiv preprint arXiv:2106.01085* (2021).
- [47] Qing Yu and Kiyoharu Aizawa. 2020. Unknown class label cleaning for learning with open-set noisy labels. In *2020 IEEE International Conference on Image Processing*. IEEE, 1731–1735.
- [48] Zhen Zhao, Zhizhong Zhang, Xin Tan, Jun Liu, Yanyun Qu, Yuan Xie, and Lizhuang Ma. 2023. Rethinking gradient projection continual learning: Stability/plasticity feature space decoupling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3718–3727.

7 EXPERIMENTS

7.1 Setting

MNIST has 10 classes with 60,000 training images and 10,000 test images. CIFAR-10 contains 50,000 training images and 10,000 test images overall. CIFAR-100 comprises 100 classes of 50,000 images for training and 10,000 for testing.

Noisy label setting: On one hand, we introduced artificially mis-labeled data in existing datasets. This refers to the asymmetric noisy label settings from [16, 33], mapping similar classes. Specifically, for CIFAR-10, the mappings are as follows: TRUCK \rightarrow AUTOMOBILE, BIRD \rightarrow AIRPLANE, DEER \rightarrow HORSE, and CAT \rightarrow DOG. For MNIST, the mappings are: 2 \rightarrow 7, 3 \rightarrow 8, 5 \rightarrow 6, 7 \rightarrow 1. Following the mapping relationships for each dataset, we randomly selected a certain proportion of samples from the classes on the left side of the mappings to be relabeled as the corresponding classes on the right side. For CIFAR-10, this involved changing the labels of 50%, 30%, and 10% of the samples from TRUCK, BIRD, DEER, and CAT to AUTOMOBILE, AIRPLANE, HORSE, and DOG, respectively. The same approach was applied to the MNIST and CIFAR-100 datasets. Compared to previous work, our noise rate settings are more challenging, as evidenced by the comparative experiments in Section 5.3 of the main text. On the other hand, we also considered the real-world noisy label dataset webvision, which has an estimated noise rate of 20% [41]. We used the top 14 largest classes by data size, resulting in a total of 47,784 images. We curated seven tasks with randomly paired classes.

Implementation Details: We train the network for 100 epochs for each task. For each dataset, both new and old tasks are configured with a mini-batch size of 512, applied to both the training and testing sets. We utilize the Adam optimizer with a weight decay of $5e-4$ and a learning rate of 0.0001. The model parameters are initialized using the default random initialization.

Table 7: Performance Comparison on CIFAR-10 and MNIST

Methods	CIFAR-10				MNIST			
	S	A_1	A_2	F	S	A_1	A_2	F
CoFiMA [31]	57.56	70.34	72.87	10.55	72.21	88.25	90.43	12.11
CoFiMA + ErrorEraser	62.64	75.61	76.95	4.02	74.94	91.04	93.29	9.15
ICL [42]	54.78	67.12	69.34	10.97	70.26	85.33	87.65	8.61
ICL+ ErrorEraser	56.95	69.41	71.25	8.51	71.89	86.96	87.14	4.32

7.2 Comparison with other SOTA CL methods.

The experimental results in Table 7 demonstrate the effectiveness of ErrorEraser in improving continual learning (CL) performance across CIFAR-10 and MNIST. By integrating ErrorEraser with CoFiMA [31] and ICL [42]), we observe consistent improvements in overall accuracy (S) and task-wise accuracy (A_1, A_2), indicating enhanced feature retention and better adaptation. Notably, ErrorEraser significantly reduces forgetting rates (F), with CoFiMA+ErrorEraser lowering F from 10.55 to 4.02 on CIFAR-10 and from 12.11 to 9.15 on MNIST, showcasing its ability to mitigate the impact of erroneous knowledge. Similar improvements are seen with ICL, where F drops from 10.97 to 8.51 and 8.61 to 4.32, respectively. These results suggest that ErrorEraser effectively refines feature representations, stabilizes model updates, and mitigates

the negative effects of noise, making it a robust and generalizable plugin for improving CL performance.

It is important to note that our goal is to propose a general plugin, ErrorEraser, suitable for various CL methods, that can accurately identify and forget erroneous knowledge learned by the model, much like humans do, rather than designing a state-of-the-art (SOTA) CL method.

7.3 Different Noise Label Rates

Table 8 presents the results on CIFAR-10 under different noise label rates, specifically set at 10% and 30%. The results indicate that as the noise label rate increases, the performance of standard CL methods deteriorates significantly. In most cases, the average accuracy on new tasks (A_1) and historical tasks (A_2) shows a marked decline. Interestingly, the overall forgetting (F) appears to decrease. However, this is not due to an improvement in the model but rather because the accuracy on each task worsens, thereby reducing the computed catastrophic forgetting. Therefore, considering the comprehensive metric (S), performance degradation is positively correlated with the noise rate.

Table 8: Results on CIFAR-10 with varying noise label rates.

	Methods	Noise rate = 10%				Noise rate = 30%			
		#Taks:5, #Class:10				#Taks:5, #Class:10			
		S	A_1	A_2	F	S	A_1	A_2	F
Normal data	EWC	57.90	71.15	74.88	14.51	57.90	71.15	74.88	14.51
	LWF	61.12	74.29	77.40	9.06	61.12	74.29	77.40	9.06
	HAT	56.77	69.62	74.78	14.96	56.77	69.62	74.78	14.96
	MAS	58.47	71.56	74.74	12.02	58.47	71.56	74.74	12.02
	TAT	51.89	64.26	70.75	20.01	51.89	64.26	70.75	20.01
	DERPP	63.65	77.63	80.17	10.37	63.65	77.63	80.17	10.37
Noisy labeled	EWC	55.82	68.84	73.46	17.35	54.35	66.97	71.65	16.23
	LWF	56.89	69.91	74.10	14.27	52.24	63.84	68.82	12.50
	HAT	42.85	51.52	51.03	0.03	44.14	53.33	51.28	0.01
	MAS	56.18	69.08	72.76	14.40	53.56	65.22	69.30	10.25
	TAT	50.42	62.42	68.95	19.63	51.37	63.67	70.42	20.68
	DERPP	63.68	77.29	79.76	10.21	58.08	70.76	73.82	9.80
Noisy labeled	EWC+ErrorEraser	68.33	82.36	83.53	3.39	66.72	80.65	81.511	4.41
	LWF+ErrorEraser	58.32	71.66	76.07	15.35	57.14	70.21	73.64	14.01
	HAT+ErrorEraser	55.79	69.92	75.25	25.41	54.86	68.84	75.34	27.15
	MAS+ErrorEraser	55.83	69.27	75.22	21.25	55.91	68.77	74.15	16.51
	TAT+ErrorEraser	53.01	64.36	65.58	5.80	56.50	68.07	69.12	2.66
	DERPP+ErrorEraser	61.81	74.80	75.86	4.88	57.00	68.36	67.69	1.06

With the addition of our ErrorEraser plugin, the impact of noisy labels is mitigated across most methods, enhancing overall performance. Particularly with methods like EWC, the performance with the ErrorEraser plugin exceeds the performance of EWC on clean data. This further demonstrates the superiority of our method in forgetting erroneous or irrelevant knowledge, enabling the model to learn the correct core knowledge.

7.4 Different Number of Tasks

Table 9 presents the results on CIFAR-100 with different numbers of tasks. We set the number of tasks to 2, 5, and 10, respectively. When the number of tasks is 2, each task contains 50 classes. The results indicate that as the number of tasks decreases and the number of classes per task increases, the performance of CL methods on normal data significantly diminishes. This is because when a task contains too many classes, such as 50, the classification model struggles to accurately classify each class. Even under these conditions, noisy labels further degrade the performance of CL methods.

Table 9: Results on CIFAR-100 for varying number of tasks.

	Methods	#Tasks:2, #Class:100				#Tasks:10, #Class:100			
		S	A ₁	A ₂	F	S	A ₁	A ₂	F
Normal data	EWC	14.65	22.0	21.93	0	37.74	41.87	42.48	4.52
	LWF	12.32	18.87	19.275	1.61	30.80	34.33	36.16	7.03
	HAT	12.21	18.55	18.78	0.91	31.60	35.10	36.73	5.51
	MAS	11.35	17.29	17.55	1.03	28.60	31.90	33.84	7.06
	TAT	16.10	28.04	33.42	21.52	13.52	16.49	28.62	46.17
	DERPP	18.21	22.03	22.28	1.01	32.14	36.16	41.52	15.77
Noisy labeled	EWC	15.29	22.95	22.96	0.06	37.63	41.73	42.12	4.19
	LWF	11.63	17.84	18.24	1.62	31.16	34.64	36.17	5.69
	HAT	12.53	16.3	17.1	9.36	9.56	10.53	10.56	0.15
	MAS	11.84	16.87	17.02	0.59	30.41	33.85	35.15	5.82
	TAT	15.38	27.04	32.66	22.44	12.25	15.02	28.43	49.01
	DERPP	14.02	22.17	22.37	3.79	30.99	35.01	39.63	16.42
Noisy labeled	EWC+ErrorEraser	21.87	33.6	34.42	3.27	49.43	54.60	54.83	2.67
	LWF+ErrorEraser	12.06	21.02	20.11	8.84	42.65	47.50	49.69	8.97
	HAT+ErrorEraser	11.59	21.62	28.57	27.8	18.43	24.04	25.98	15.47
	MAS+ErrorEraser	12.23	20.27	22.66	9.55	36.48	40.89	45.32	13.82
	TAT+ErrorEraser	18.04	29.06	31.36	9.20	12.23	16.07	29.44	30.21
	DERPP+ErrorEraser	14.20	23.10	22.14	4.65	28.24	34.07	36.80	3.90

However, with the application of our ErrorEraser plugin, the performance of most CL methods improves across all task settings. This further validates the effectiveness and generality of our method.

7.5 Per-Task Performance

To illustrate the performance changes for each task, we visualize the new task accuracy and old task forgetting rates on each task after noise label introduction in Figure 7. Here, we specifically set the second task (Task ID=1) to contain 50% noise labels to observe the results for each task.

We found that with the application of ErrorEraser, the accuracy of the second task significantly improved after its completion, and the performance of each subsequent new task also showed noticeable improvement. Regarding forgetting, the introduction of noise labels caused the EWC method to have a forgetting rate of 10.80% for the first task after the second task was completed. By applying ErrorEraser, the forgetting rate for the first task was reduced to 1.45%. This demonstrates that our method effectively forgets erroneous knowledge introduced by noisy labels, enhancing new task accuracy while preventing catastrophic forgetting of old tasks.

7.6 Visualization of Decision Space

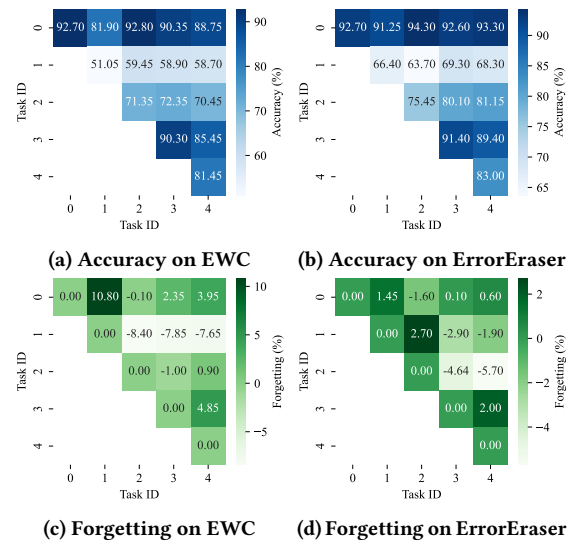
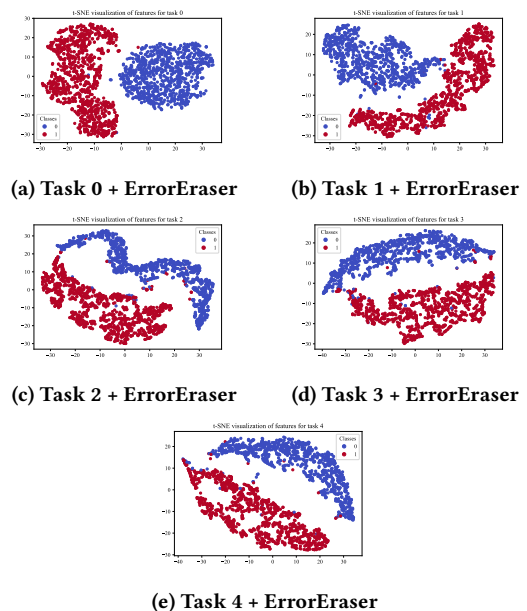
In Figure 8, we present the t-SNE visualization of the feature extraction for the 1st task on the MNIST dataset after applying ErrorEraser across five tasks. The results indicate that after each task completion, the decision boundaries between classes in the task are clear. This further verifies that ErrorEraser effectively forgets erroneous knowledge to ensure performance on new tasks.

8 DISCUSSION

The CL scenario studied in this paper is based on task incremental learning, where the task ID is known. Each task involves learning from image data, without involving multi-task learning across different data types. The classes within each task are independent and non-overlapping. Additionally, this study focuses solely on the data bias introduced by noisy labels.

In the future, we plan to explore the impact of noisy labels on CL across all three scenarios, aiming to develop more comprehensive

and general solutions. We will also investigate the challenge of overlapping classes between tasks, where the impact of noisy labels may be more severe and complex. Finally, we intend to address more realistic data bias situations, such as entire classes being incorrectly labeled, noisy data (e.g., blurred or damaged images), to develop robust CL methods. We aim to continuously contribute to the continual learning community.

**Figure 7: Accuracy and forgetting before and after ErrorEraser.****Figure 8: t-SNE visualization of the extracted feature outputs for the each task (1-th training data contains 50% noisy labels) in the MNIST test dataset.**