

ADDRESSING THE DEVASTATING EFFECTS OF SINGLE-TASK DATA POISONING IN EXEMPLAR-FREE CONTINUAL LEARNING

Stanisław Pawlak^{1*} Bartłomiej Twardowski^{2,3} Tomasz Trzcíński^{1,2} Joost van de Weijer³

¹ Warsaw University of Technology, Poland ² IDEAS Research Institute, Poland

³ Computer Vision Center, Universitat Autònoma de Barcelona, Spain

ABSTRACT

Our research addresses the overlooked security concerns related to data poisoning in continual learning (CL). Data poisoning – the intentional manipulation of training data to affect the predictions of machine learning models – was recently shown to be a threat to CL training stability. While existing literature predominantly addresses scenario-dependent attacks, we propose to focus on a more simple and realistic single-task poison (STP) threats. In contrast to previously proposed poisoning settings, in STP adversaries lack knowledge and access to the model, as well as to both previous and future tasks. During an attack, they only have access to the current task within the data stream. Our study demonstrates that even within these stringent conditions, adversaries can compromise model performance using standard image corruptions. We show that STP attacks are able to strongly disrupt the whole continual training process: decreasing both the stability (its performance on past tasks) and plasticity (capacity to adapt to new tasks) of the algorithm. Finally, we propose a high-level defense framework for CL along with a poison task detection method based on task vectors. The code is available at: <https://github.com/stapaw/STP.git>.

1 INTRODUCTION

Continual learning (CL) aims to accumulate knowledge from a stream of data, allowing models to adapt to new information while preventing forgetting of already acquired knowledge (Delange et al. (2021); Masana et al. (2022)). This continuous process seeks to strike a balance between stability and plasticity, enabling the continual learner to maintain performance on previous data while adapting to new data. CL training is based on the implicit assumption that a stream of training data is representative of the data encountered during inference time – samples are from the same domain or share similar characteristics. Data poisoning attacks violate this assumption by manipulation of the training data aimed at degradation of model performance at inference time (Cinà et al. (2023)).

Recent works (Li & Ditzler (2022); Umer et al. (2020); Kang et al. (2023); Abbasi et al. (2024); Han et al. (2023)) demonstrate successful data poisoning attacks on CL, highlighting the importance of the security challenges present in various CL environments. Attacks against CL proposed so far differ in terms of considered threat model, learning scenarios and poisoning methods. With so many hidden assumptions, it is much harder to develop and propose adequate defenses. Moreover, most attacks use a non-restrictive threat model which is not realistic: adversaries have access to multiple tasks within the learning stream (Umer et al. (2020); Li & Ditzler (2022; 2023); Han et al. (2023)) and/or have the ability to access the trained model during training to produce adversarial samples as a poison (Li & Ditzler (2022; 2023); Han et al. (2023)).

To address these issues, we propose a new, more simple and realistic Single-Task Poisoning (STP) setup where an adversary has only access to a current CL task data, and access to the trained model is not allowed. To assess how much damage could be done when poisoning a single task, we propose two types of poisoning attacks: uniform and classifier bait. The former contaminates the task data with one of the standard image corruptions (see Figure 2). The latter corrupts only a subset of classes within a task to create data distribution differences between classes and uses corruptions as a bait for the classifier to focus on those artificially introduced spurious correlations rather than real class-discriminatory features. We want to highlight, that the goal of the STP framework is not to assess the damage done to the poisoned task performance, but rather to investigate the performance on previous and future tasks in the sequence during the CL training session. The common aim of the attacks proposed so far was to affect the stability of the model by increasing the forgetting on previous tasks. To the best of our knowledge, we are the first to investigate

* Corresponding author: stanislaw.pawlak.dokt@pw.edu.pl

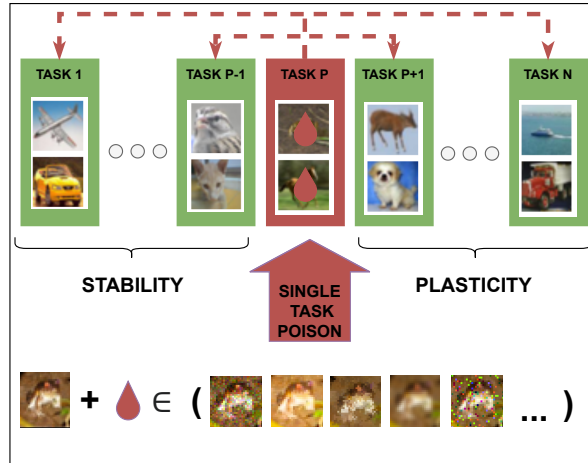


Figure 1: **The Single-Task Poison (STP) setting.** We propose a new, more realistic CL setup to investigate data poisoning attacks, where an adversary has only access to a single task (T_p) with no knowledge about other tasks in the sequence and no access to the trained model. We show that STP attacks can strongly affect the stability and the plasticity of the model.

how data poisoning attacks can affect model performance on future tasks (plasticity) and the first to show an attack that is affecting both the stability and plasticity of the model at the same time.

We summarise the contributions of this work as follows:

- We focus on class-incremental setup with exemplar-free methods and propose the Single-Task Poison – a more realistic adversarial setup for CL, where adversary’s data access is restricted only to one task in the sequence, with no knowledge about previous tasks or access to the trained model.
- We show that even under this restricted threat model, data poisoning attacks can have devastating effects. To the best of our knowledge, we are the first to investigate the consequences of poisoning on future tasks performance. We show that an attack can increase catastrophic forgetting and impede future training at the same time.
- We present a high-level defense framework for CL with a poison task detection method based on task vectors.

2 RELATED WORK

Data poisoning attacks. Data poisoning attacks are considered one of the strongest concerns for the use of offline-trained ML models in real life (Cinà et al. (2023)), and thus multiple attacks (Nguyen & Tran (2020); Yang et al. (2017); Muñoz-González et al. (2017); Shafahi et al. (2018)) and defenses (Taheri et al. (2020); Hong et al. (2020); Hayase et al. (2021); Zhu et al. (2021); Yang et al. (2022)) were proposed. The main three categories of data poisoning attacks are (Cinà et al. (2023)): indiscriminate, targeted and backdoor attacks. 1) Indiscriminate attacks involve poisoning the training data with the intention of compromising the overall performance of the model. They introduce noise or biases in a way that undermines the model’s accuracy and generalization. 2) Targeted attacks aim to manipulate the model’s behavior in particular situations or for certain inputs. Poisoned data influence the model’s predictions for certain targeted samples while maintaining normal performance for all the others. 3) Backdoor attacks are targeted attacks that poison the data with a hidden pattern (trigger) added to training samples. This trigger does not impede network performance during training, but allow adversaries to manipulate the model’s behavior during inference time. The backdoor remains hidden and passive under normal circumstances but can be activated by adding a trigger to specific samples during inference time. Using triggers adversaries can control the model’s outputs in a predefined manner (e.g. by changing the prediction of sample with backdoor to a specific label).

The issue of data poisoning in CL setups is largely unexplored, with only a few recent works discussing specific attacks (Li & Ditzler (2022); Umer et al. (2020); Kang et al. (2023); Abbasi et al. (2024)) and defenses (Umer & Polikar (2023); Wang et al. (2022)). Table 1 compares data poisoning methods in CL. Most focus on Task Incremental scenarios, where task IDs are available during inference, while our work emphasizes the understudied Class Incremental scenario. Backdoor attack methods (Umer et al. (2020); Umer & Polikar (2022)) require data access during training and inference to trigger the attack. The most common attack technique relies on accessing the trained model to create

Table 1: **Comparison of data poisoning attacks in CL.** CL scenarios: DI – Domain Incremental, TI – Task Incremental, CI – Class Incremental (van de Ven & Toliás (2019)). Data Access: assumption of the adversary’ access to a stream of data. Model Access: assumption of the adversary’ access to the continually trained model. Our approach is the first data poisoning attack considering CI with single access to the training data stream and without access to the trained model.

Work	CL data poisoning attacks			Threat Model	
	CL Scenario			Data Access	Model Access
	DI	TI	CI		
Umer et al. (2020)	✗	✓	✗	Multiple	✗
Li & Ditzler (2022)	✗	✓	✗	Multiple	✓
Li & Ditzler (2023)	✓	✓	✗	Multiple	✓
Han et al. (2023)	✗	✓	✗	Multiple	✓
Kang et al. (2023)	✗	✓	✗	Single	✗
Abbasi et al. (2024)	✗	✓	✗	Single	✓
Ours	✗	✗	✓	Single	✗

adversarial samples for previous tasks (Li & Ditzler (2022; 2023); Han et al. (2023)) or to revert the model for task information (Abbasi et al. (2024)). We share a similar goal with Kang et al. (2023) in making the threat model more realistic. However, while Kang et al. (2023) focuses on exploiting vulnerabilities in generative replay (a generative model is unable to generate a label-changing pattern that is present during first classifier training, but is not generated during replay phase, which leads to a postponed label-switching attack and increase in the forgetting), our focus differs.

Data poisoning defenses Most of the standard defenses against data poisoning (like training data sanitization and outlier detection methods) assume that only a small fraction of the training data is poisoned Cinà et al. (2023). Although this assumption is correct for a standard stationary training, we argue that in CL this assumption is not certain. We show how an adversary can even fully poison a small task in a CL sequence, without being noticed immediately (see 10). Thus, most standard defenses against data poisoning are not suited to CL training. Although data poisoning in CL have recently gained increasing attention resulting in multiple new attacks proposed, there is a significant gap in terms of discussing and proposing effective defensive measures. Defensive approaches were recently proposed only for a specific subset of poisoning attacks, namely backdoor attacks (Umer & Polikar (2023); Wang et al. (2022)). A subset of methods based on robust training, model inspection and model sanitization (Cinà et al. (2023)) could be adapted to prevent data poisoning in CL, but as of now they are still not evaluated under CL setup and thus there is no off-the-shelf defense yet.

CL Threat model. We consider the simplest and most realistic scenario of the attack: an adversary has access only to the data of one task in the sequence. To our knowledge, this is the most constrained threat model yet considered. Unlike previous models, our 1) adversary lacks access to the CL model (Li & Ditzler (2022; 2023); Abbasi et al. (2024); Han et al. (2023)), and 2) adversary has no access to or knowledge of data from previous (Li & Ditzler (2022; 2023); Han et al. (2023); Wang et al. (2022)) or 3) future tasks (Umer et al. (2020); Umer & Polikar (2022; 2023)).

3 SINGLE-TASK DATA POISONING

In this section, we present Single-Task Poisoning (STP) framework to investigate the effects of data poisoning attacks in continual learning. We start with an overview and justification of STP, followed by the description of data poisoning methods, attacks, and defenses proposed in this work.

3.1 SINGLE TASK DATA POISONING: WHAT IT IS AND WHY IT MATTERS

STP is a setup to investigate the effects of data poisoning in CL. It puts more restrictions on the adversary, who has only access to a single task data and cannot access the trained model. STP is motivated by the fact that in real-life scenarios we have access neither to the model nor to all available training data stream. In the following, we show these constraints, explain and justify the STP setting by considering a defensive party that is continually training a model on a data stream obtained from multiple sources (see Figure 1).

Firstly, STP employs a more restrictive threat model, because the potential adversary data provider has no access to the data of other providers, nor to the trained model, and he or she does not know when the data would be used for training. Secondly, STP considers poisoning high-percentages of data samples, which is not common for the data poisoning settings. In stationary training, a single data provider is not able to poison a high percentage of training dataset due to the data volume. Contrary, CL training is divided into multiple stages and poisoning heavily one small

task is possible. Finally, STP provides a simple threat model to investigate how much damage could be made by a single poisoned CL task. It is true that the data from a single provider could be used within the one task or split into several tasks which could introduce additional vulnerabilities. While many previous attacks consider multi-stage attacks, we argue to focus firstly on a simpler scenario. STP provides a framework to investigate the devastating effects of poison from a single task and is a step forward to build defensive mechanisms that could be used for more secure CL training in general.



Figure 2: **Example of image corruption: Gaussian blur with increasing severity (1-5 from left to right).**

3.2 DATA POISONING CONTINUAL LEARNING WITH STP

3.2.1 TASK POISONING WITH IMAGE CORRUPTIONS.

We use the CIFAR-C (Hendrycks & Dietterich (2019)) set of image corruptions to transform images of the poisoned task. These corruptions are well established and commonly used, with clear publicly available algorithm and code to create them. We use them to propose the simplest possible poisoning attack, without any additional knowledge, as the STP threat model is very restrictive towards the attacker and he or she cannot use other poisoning methods proposed so far (they rely on the access to the model or information about other tasks in the sequence for optimization). Additionally, proposed different severity levels were used in the experiments, to show the strength vs. detectability trade-off faced by the attacker Frederickson et al. (2018).

While CIFAR-C dataset is usually used for test-time adaptation scenarios and thus benchmark consists of corrupted test data, we use those predefined transformations to create corrupted images for training the poisoned task. Adapting these predefined transformations results in a clean-label, non-adversarially optimized data poisoning (poisoning based on a pixel distribution change after transformations like e.g. gaussian blur). For examples of corruption severities, see 2. All corruption types are shown in the Appendix 15.

3.2.2 ATTACKS.

Let T_p be the task that is poisoned during CL training with D_p being the dataset used in that task. Let D_{clean} be the original dataset for the attacked task, $D_{poisoned}$ be the poisoned copy of D_{clean} .

Transformation from D_{clean} into $D_{poisoned}$ is defined by three parameters (see Figure 3):

- pcp – poisoned class percentage, defining how many classes are affected by corruptions in $D_{poisoned}$
- pn – number of poisons (i.e. different corruptions) used to poison data in $D_{poisoned}$
- pp – the percentage of exemplars poisoned (corrupted) in each class of $D_{poisoned}$

The poisoning attack is changing the distribution of T_p data by applying one of the corruptions to pp % of images from pcp classes of D_{clean} . Let I_n be the subset of n indices randomly chosen from D_{clean} , $I_n \subseteq \{1, 2, \dots, |D_{clean}|\} = I$. Then D_p is a set containing n poisoned and $|D_{clean}| - n$ clean samples:

$$D_p = D_{poisoned}[I_n] \cup D_{clean}[I/I_n] \quad (1)$$

Naming convention for attack evaluation. We evaluate STP with four kinds of attacks:

1. BASE – the simplest attack with single corruption added to all classes in T_p ,
2. BAIT – use single corruption as a bait for classifier: we poison only half of the classes in T_p to attract the model to discriminate rather based on introduced data distribution differences than the class differences within the same data distribution,
3. MULTIBASE – we poison all classes in T_p with different corruptions to check if the attacker can obtain additional gains by increasing the number of poisons. Classes are poisoned proportionally with pn corruptions.
4. MULTIBAIT – we poison half of the classes in T_p . Classes are poisoned proportionally as in MULTIBASE attack.

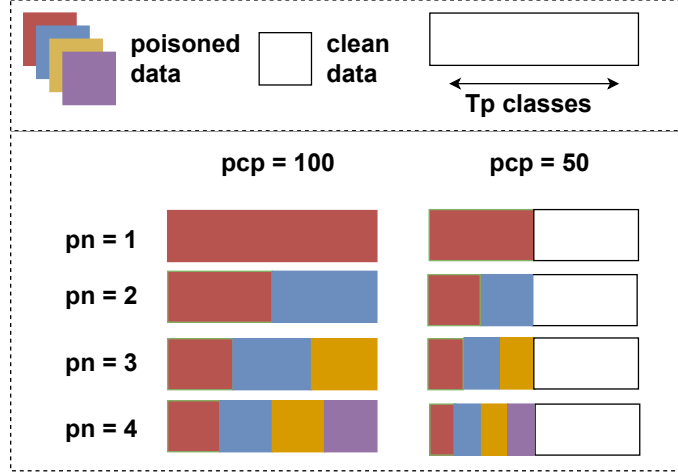


Figure 3: **Main parameters of STP attacks: pcp and pn .** pcp defines how many classes are poisoned in poisoned task T_p , pn how many poisons (different corruptions) are used.

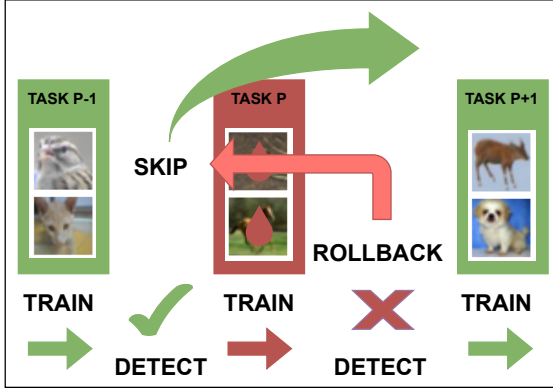


Figure 4: **The high-level defense framework.** We propose to add after each CL task a poison-detection phase to roll back changes in model weights upon detection of poison.

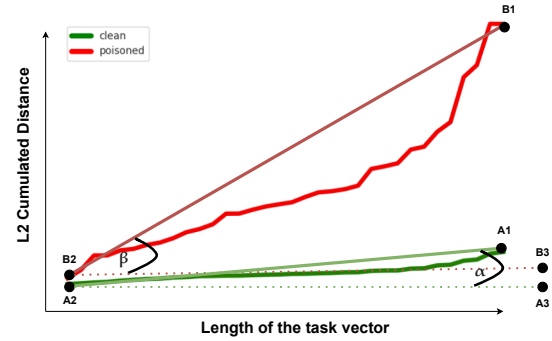


Figure 5: **Difference between clean and poisoned task vectors increase with more layers contributing to task vectors.** Consequently, we detect the poison when $\beta > \alpha$.

3.2.3 DEFENSE.

The defense framework involves three steps: checkpointing the model, poison task detection, and (optional) weight rollback (see Figure 4). First, we save the model weights before training on a new task. After training, a poison detection test is conducted. If the task is detected as poisoned, the model’s weights are rolled back to the pre-task state, the poisoned task is skipped, and training proceeds with the next task.

Detection using task vectors. The core of our framework is poisoned task detection. We propose a simple task-vector-based method that uses one clean task at the start of training to calculate a threshold. This threshold is then applied to all tasks in the stream to determine if a task is poisoned. Task vectors (Ilharco et al. (2023)) are vector representations in the parameter space of a neural network that encapsulate the changes required to adapt the model from one task to another. They are computed by taking the difference between the model parameters fine-tuned on different tasks, effectively capturing the task-specific adaptations needed. We propose to use them for the poison task detection: by comparing task vectors for the model before and after the task training. For both clean and poisoned task distance between *before* and *after* vectors is increasing when we use more network layers to create task vectors, but Figure 5 shows that for poisoned task distance increase is more significant. Thus, we define the angle α with single clean task data to denote how fast the distance increase when trained on clean data (see algorithm 1). Then, we test the β obtained for each of the following tasks. If $\beta > \alpha$ we denote the task as poisoned.

Algorithm 1: Calculation of Threshold Angle α

Input: A sequence of tasks $\{T_1, T_2, \dots, T_n\}$, continually trained model M_{θ_i} , where θ_i are the weights of the model after training on task T_i

Output: Threshold angle α

```

1 Select a task  $T_c$  from the beginning of the sequence (assumed to be clean);
2 foreach seed  $s_i \in \{s_1, s_2, \dots, s_k\}$  do
3   Train  $M_{\text{before}}$  ( $M_{\theta_{c-1}}$ ) on  $T_c$  to obtain  $M_{\text{after}}$ ;
4   Compute activations  $A'_{\text{before}}$  and  $A'_{\text{after}}$  using  $T_c$  data;
5   Aggregate activations using mean into single vectors  $A_{\text{before}}$  and  $A_{\text{after}}$ ;
6   Compute task vector:  $v = \|A_{\text{after}} - A_{\text{before}}\|_2$ ;
7   Compute cumulative sum vector:  $c = \text{cumsum}(v)$ ;
8   Let  $A_1 = c_{\text{last}}$ ,  $A_2 = c_{\text{first}}$ ;
9   Define  $A_3$  with same x-coordinate as  $A_1$  and y-coordinate as  $A_2$ ;
10  Compute angle  $\alpha'_i = \angle A_1 A_2 A_3$ ;
11 Aggregate all angles:  $\alpha = \text{Aggregate}(\{\alpha'_1, \alpha'_2, \dots, \alpha'_k\})$ ; // e.g., 90th percentile or MAX
12 return  $\alpha$ ;
```

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

CL setup. In the main experiments, we use a class-incremental setup (van de Ven & Tolias (2019)) that divides datasets into tasks with disjoint classes and lacks task identifiers during inference. We investigate data poisoning attacks on classical exemplar-free CL methods, LwF (Li & Hoiem (2017)) and EWC (Kirkpatrick et al. (2017)). Basic experiments use CIFAR10 split into five tasks of two classes each. For most experiments, we use CIFAR100, split into six tasks, with the first task having 50 classes and the others 10 classes each, to assess attacks on tasks with more than two classes.

Model and training details. We use the Resnet32 model for all experiments in the main section of the work. We use Resnet18 for additional experiments on Tiny-Imagenet (Le & Yang (2015)), added in the appendix (see A.1). We use implementations from FACIL (Masana et al. (2022)) framework for CL training. For CIFAR10 we start each task with learning rate equal to 0.05 with weight-decay 0.0001 and momentum 0.9. We train 200 epochs for each task. Additional information about learning parameters are in the appendix B.

Datasets We test data poisoning attacks on CIFAR10 and CIFAR100 datasets. Both contain 32x32 images (Alex (2009)). Additional experiments with bigger images on Tiny-Imagenet (64x64) (Le & Yang (2015)) are in the appendix A.1.

Defense evaluation. To evaluate proposed poison task detection method, we prepared a dataset of tasks created from different subsets of CIFAR100, each having ten classes. In realistic scenario clean tasks should be much more frequent than poisoned ones. Thus, we used a dataset with 92% clean and 8% poisoned tasks.

4.2 ATTACKS EVALUATION

We investigate STP attacks with these questions:

1. What is the difference between poisoning a single task in exemplar-free CL compared to poisoning single task in joint training?
2. What is the decrease in the model performance after poisoning task T_p with proposed STP attacks? How poisoning affects continual learner stability – performance on tasks before T_p – and plasticity – performance on the following tasks?
3. Is it better for the attacker to poison all data with the same corruption (i. e. using BASE, BAIT attacks), or use multiple corruptions (MULTIBASE, MULTIBAIT attacks) for different classes?
4. How can the adversary address the strength vs detectability trade-off (Frederickson et al. (2018)) of the data poisoning attack? What percentage of T_p data should be poisoned for STP attack to be effective?

Table 2: **Average accuracy when poisoning a single T_p task in CL sequence.** Grey (CLEAN) rows show results of runs without poisoning T_p task, while black (BASE, BAIT, MULTIBASE, MULTIBAIT) rows show results for different attacks. Difference between CLEAN and poisoned runs is shown in the parenthesis. As most data poisoning works in CL we report accuracies after the poisoned task (ACC AT POISONING TIME). Additionally we report accuracies at the end of the full CL sequence (FINAL ACC) to show the impact of poisoning for further training. Experimental setup: CIFAR100(50/10), CIFAR10(2/2), BASE: $pcp = 100$, $pn = 1$, $pp = 100$, BAIT: $pcp = 50$, $pn = 1$, $pp = 100$, BASE: $pcp = 100$, $pn = 5$, $pp = 100$, BASE: $pcp = 50$, $pn = 5$, $pp = 100$.

ROW	ATTACK	METHOD	DATASET	P	ACC AT POISONING TIME		FINAL ACC			
					T_p	BEFORE T_p	T_p	BEFORE T_p	AFTER T_p	TOTAL
1	CLEAN	JOINT	CIFAR10	3	95.2	95.6	93.6	93.4	92.4	93.1
2	BASE	JOINT	CIFAR10	3	0.2 (-95.0)	96.6 (1.0)	0.1 (-93.5)	94.1 (0.7)	93.1 (0.7)	74.9 (-18.2)
3	BAIT	JOINT	CIFAR10	3	0.3 (-94.9)	96.5 (0.9)	46.2 (-47.4)	93.8 (0.4)	92.7 (0.3)	83.8 (-9.3)
4	CLEAN	JOINT	CIFAR100	4	73.0	71.2	71.3	68.0	66.5	68.1
5	BASE	JOINT	CIFAR100	4	1.2 (-70.0)	72.2 (1.0)	1.0 (-70.3)	69.7 (1.7)	67.1 (0.6)	57.4 (-10.7)
6	BAIT	JOINT	CIFAR100	4	36.8 (-36.2)	71.9 (0.7)	36.0 (-35.3)	68.7 (0.7)	66.6 (0.1)	62.6 (-5.5)
7	MULTIBASE	JOINT	CIFAR100	4	4.5 (-68.5)	72.2 (1.0)	4.3 (-67.0)	69.5 (1.5)	66.7 (0.2)	57.7 (-10.4)
8	MULTIBAIT	JOINT	CIFAR100	4	40.3 (-32.7)	71.9 (0.7)	39.7 (-31.6)	69.0 (1.0)	66.4 (-0.1)	63.3 (-4.8)
9	CLEAN	LWF	CIFAR10	3	67.2	71.7	50.2	62.3	63.9	60.5
10	BASE	LWF	CIFAR10	3	56.3 (-10.9)	70.1 (-1.6)	39.9 (-10.3)	59.1 (-3.2)	65.4 (1.5)	57.8 (-2.7)
11	BAIT	LWF	CIFAR10	3	45.8 (-21.4)	48.5 (-23.2)	41.9 (-8.3)	43.8 (-18.5)	50.2 (-13.7)	46.0 (-14.5)
12	CLEAN	LWF	CIFAR100	4	78.9	50.8	61.7	30.6	69.7	48.8
13	BASE	LWF	CIFAR100	4	64.4 (-14.5)	36.3 (-14.5)	45.5 (-16.2)	17.3 (-13.3)	68.6 (-1.1)	39.1 (-9.7)
14	BAIT	LWF	CIFAR100	4	52.2 (-26.7)	41.5 (-9.3)	44.3 (-17.4)	23.5 (-7.1)	62.3 (-7.4)	39.9 (-8.9)
15	MULTIBASE	LWF	CIFAR100	4	37.1 (-41.8)	46.8 (-4.0)	28.7 (-33.0)	27.4 (-3.2)	66.1 (-3.6)	40.5 (-8.3)
16	MULTIBAIT	LWF	CIFAR100	4	54.5 (-24.4)	45.7 (-5.1)	47.3 (-14.4)	27.1 (-3.5)	65.3 (-4.4)	43.2 (-5.6)
17	CLEAN	EWC	CIFAR10	3	34.3	22.2	20.0	11.6	24.0	18.2
18	BASE	EWC	CIFAR10	3	19.7 (-14.6)	25.7 (3.5)	4.2 (-15.8)	15.0 (3.4)	29.5 (5.5)	18.6 (0.4)
19	BAIT	EWC	CIFAR10	3	40.6 (6.3)	6.3 (-15.9)	16.7 (-3.3)	3.7 (-7.9)	24.9 (0.9)	14.8 (-3.4)
20	CLEAN	EWC	CIFAR100	4	82.3	27.8	35.9	17.8	57.2	33.9
21	BASE	EWC	CIFAR100	4	48.2 (-34.1)	17.7 (-10.1)	25.6 (-10.3)	18.9 (1.1)	57.4 (0.2)	32.9 (-1.0)
22	BAIT	EWC	CIFAR100	4	51.9 (-30.4)	15.2 (-12.6)	22.4 (-13.5)	18.2 (0.4)	56.8 (-0.4)	31.8 (-2.1)
23	MULTIBASE	EWC	CIFAR100	4	37.2 (-45.1)	20.1 (-7.7)	13.1 (-22.8)	18.0 (0.2)	58.9 (1.7)	30.8 (-3.1)
24	MULTIBAIT	EWC	CIFAR100	4	56.6 (-25.7)	20.5 (-7.3)	25.1 (-10.8)	17.5 (-0.3)	56.6 (-0.6)	31.8 (-2.1)

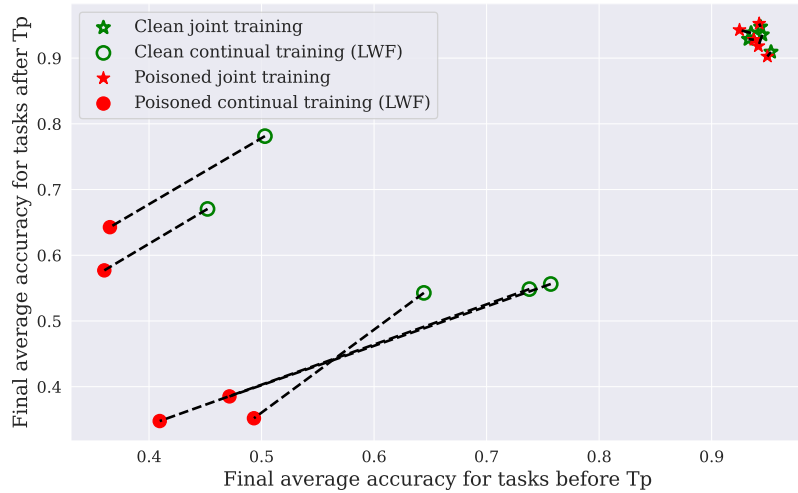


Figure 6: **Poisoning single task (T_p) training data has devastating effect for the performance in CL.** In CL single task poisoning may decrease the performance for past and future classes, while in joint training, performance on non-poisoned subset of classes is not affected. Experimental setup: CIFAR10, LWF, Resnet32, $p = 3$.

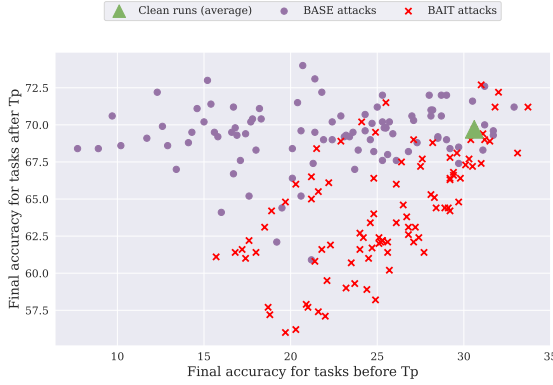


Figure 7: **Comparison of corresponding BASE and BAIT attacks impact on the performance on past and future tasks.** BAIT attacks are stronger than BASE while poisoning much less exemplars of T_p . Experimental setup: CIFAR100, LWF, Resnet32.

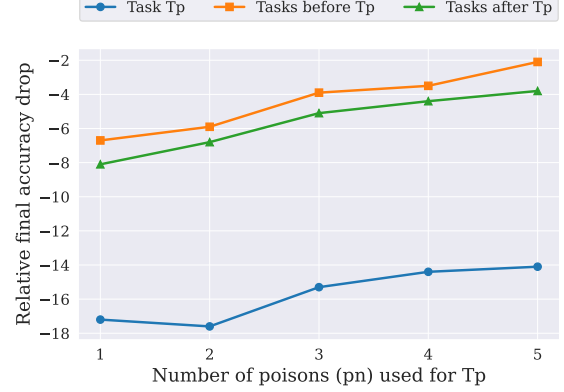


Figure 8: **Relative accuracy drop for BAIT ($pn = 1$) and MULTIBAIT ($pn > 1$) attacks using different number of poisons.** Increasing pn decrease poison effectiveness. Experimental setup: CIFAR100, LWF, Resnet32.

We answer them with the following observations:

1) Poisoning has more severe consequences in CL than in joint training. Figure 6 shows the difference between poisoning part of the data in joint training on CIFAR10 and the same part of the data treated as single task in CL sequence. The former does not affect the performance on the rest of the classes (poisoned red-filled stars cover the same area of the plot as clean green open stars). The latter affects both the performance on classes from past and future tasks (poisoned red-filled dots are in the left-down direction from clean green-open dots). The reason behind the above observation is that in JOINT case the model is trained on clean and poisoned classes at the same time, while in CL sequence we train a model on isolated poisoned data without the access to clean samples from previously seen classes.

2A) A poisoning single task in CL sequence decreases performance not only on T_p , but also on past and future tasks. As expected, the performance decrease is usually the biggest for the poisoned task (see T_p columns in Table 2). However, Table 2 presents further evidence for observation 1): there is almost no difference in accuracy for BEFORE T_p and AFTER T_p in JOINT training (for rows 2-3 and 5-8 differences between CLEAN and poisoned training in the parenthesis are close to 0), while there is a significant difference when poisoning CL sequence (for rows 10-11, 13-16 differences are negative). At the same time, we see that the relative performance drop is much bigger for tasks before T_p than after: overriding patterns to distinguish old classes without access to the data is a bigger problem than struggles in learning new patterns in tasks after T_p .

2B) BAIT attacks affect performance on tasks after T_p stronger than BASE attacks. Figure 7 expands evidence presented in Table 2 and show the effects of BASE and BAIT attacks with different corruptions on CIFAR100. Interestingly, BAIT attacks are stronger than BASE while poisoning much less exemplars of T_p . BASE attacks primarily reduce accuracy on past tasks, affecting model stability, while BAIT attacks decrease performance on both past and future tasks. We hypothesize that spurious correlations ('baits') in T_p data impair future task learning more than single corruption-based data distribution shifts in BASE attacks.

3) MULTIBASE and MULTIBAIT attacks are less effective than their counterparts with single corruption. Comparing results for multi-corruption attacks and single corruption attacks in Table 2 we see a bigger performance decrease for the latter ones: impact of one corruption is stronger than multiple, when we poison constant number of classes (with constant pcp). Results presented on Figure 8 support this observation. The bigger the number of used corruptions, the smaller the poisoning effect on both past and future tasks.

4A) Adversary may use different corruption severities to adjust the attack strength. Figure 9a) shows relative final accuracy drop for BAIT attack with image corruptions of different severity (see Figure 2). Corruptions with severity equal to 1 are negligible and can be seen as augmentation. Therefore only accuracy of T_p is slightly affected. However, increasing severity correlates with increasing performance drop on all tasks. Thus, corruption severity can be seen as a kind of poisoning 'adversarial budget' for the attacker. Smaller severity means smaller performance drop, but also smaller chance of poison detection.

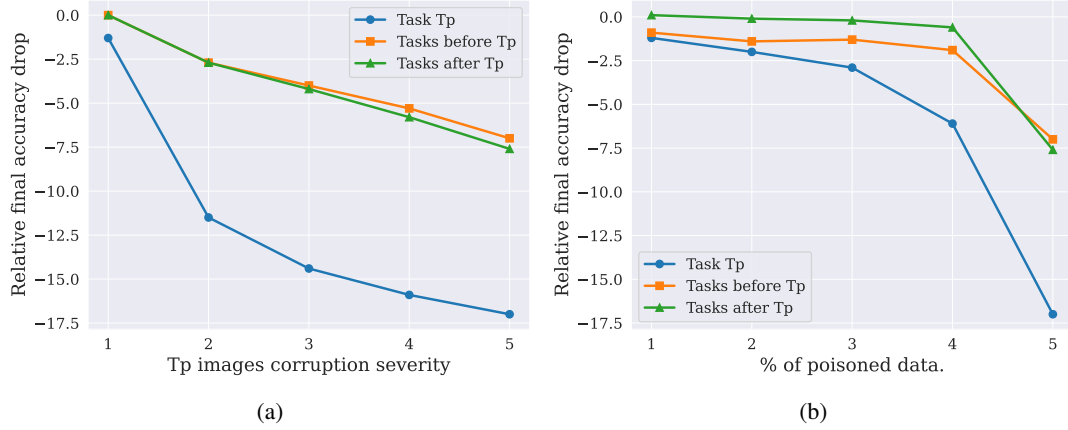


Figure 9: **BAIT attack effect for different corruption severities (a) and % of poisoned data (b).** (a) Accuracy drop for past and future tasks correlates with image corruption severity. (b) Significant part of T_p must be poisoned to create an effective attack. Note that we report % of full training data for CIFAR100, 5% corresponds to fully poisoning ($pp = 100$) half of classes in T_p ($pcp = 50$), which corresponds to strongest possible BAIT attack. Experimental setup: CIFAR100, LWF, Resnet32.

4B) Significant part of T_p must be poisoned to create an effective attack. Our additional experiments show that adversary must poison large fraction of samples in T_p (Figure 9b)) to obtain a significant effect. Poisoning 4% of data (i.e. 80% of exemplars in half of the classes in T_p) is much less effective than poisoning 5% (i.e. 100% of exemplars in half of the classes in T_p). This may be seen as an attack limitation (because the attacker cannot control the strength of the attack by poisoning small part of the T_p data).

4.3 DEFENSE EVALUATION

We explore defense against STP attacks via these questions:

5. Is the defending party able to actively detect poisoning during learning task T_p based on the validation accuracy?
6. How effective is the proposed poison task detection?

We answer them with the following observations:

5) STP attacks obtain a reasonable validation accuracy on task T_p . An important characteristic of a successful data poisoning attack is its stealthiness, that ensures the attack is not easy to spot during training on the T_p task. As the data poisoning attacks usually impede the training, the most basic defense mechanism is to monitor the accuracy of the classifier on a hold-out clean validation dataset representative to all classes. The defense against data poisoning attack is more challenging in CL. With no access to previous tasks' data the defender party can only monitor the accuracy for the current task. Moreover, in realistic CL setup, the defender party does not have access to a clean validation dataset: it is obtained by splitting the data from a new task into train and validation sets. Then, if the task is poisoned with STP it affects also the validation set and the defender party cannot easily discover the attack only by monitoring current task accuracy. Figure 10 supports these claims and show that the adversary obtains reasonable accuracy for poisoned task for all attacks. Note, that the defender party is not able to tell beforehand what should be the accuracy for an unknown task during CL training. While there is almost no difference between validation and test data in clean data case, the difference is significant for poisoned validation set and clean test showing the data distribution change caused by image corruptions. Note that clean test data is not available for the defender party during training and here used only for evaluation purposes.

6A) Poison task detection using task vectors correctly identify most of the poisoned tasks. Defender party is looking for a trade-off between the number of correctly detected attacks (True Positives) and false alarms when task is not poisoned (False Positives). This trade-off is connected with angle threshold selection (see Figure 5). The bigger the angle, the less false alarms, but also less correct attack detections. These trade-off can be presented on Precision-Recall curve (see Figure 11). While selecting the threshold and further defense evaluation is discussed in detail in the appendix A.2, we report results for the most straightforward approach taking the maximum clean angle

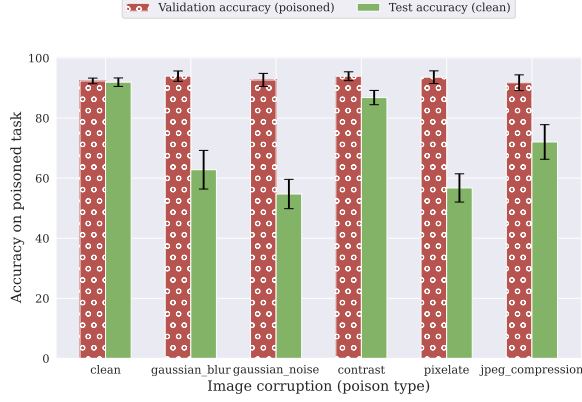


Figure 10: **STP attacks obtain a reasonable validation accuracy, while real performance on test data is much worse.** While there is no difference between evaluation when T_p is not poisoned ('clean' bars), there is significant drop for test data after attack. However, during CL training defender party has only access to (possibly) poisoned validation data, which cannot be simply used to detect the attack. Experimental setup: CIFAR100, LWF, Resnet32, $T_p = 4$.

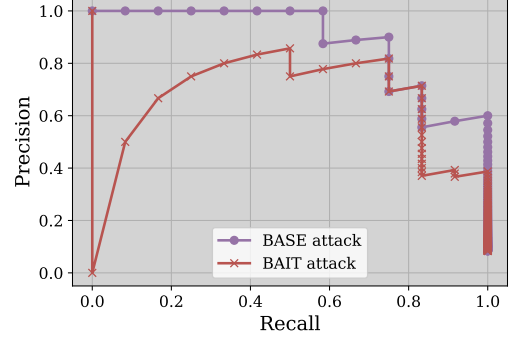


Figure 11: **Precision-Recall curve for our defense method.** BAIT attacks are harder to detect than BASE attacks, due to poisoning less exemplars. Experimental setup: CIFAR100, LWF, Resnet32, $T_c = 2$, $T_p = 4$.

calculated from one of the tasks from the beginning of CL sequence. For that threshold we obtain 92.3% for BASE and 86.2% of total accuracy for BAIT attacks (predicting correctly respectively 80% and 60% of attacks).

6B) BAIT attacks are harder to detect than BASE attacks. Results presented on Figure 11 confirm that BASE attacks are easier to detect than BAIT. The main reason is that the latter poison only half of the classes in T_p , thus change in T_p data distribution is more difficult to detect.

Limitations In this work we present the STP framework to investigate and mitigate data poisoning threats in CL. It is beneficial to start building defensive mechanisms from simplest possible setup, where adversary can poison only one task in the sequence, but it comes with a cost. STP is useful when investigating memory-free CL methods, but less with those using exemplars. In STP attack only one task can be poisoned, thus CL methods with memory have (a non-realistic) guarantee during T_p task that all stored previous samples are clean, which counter STP attack by making it more similar to the case of joint training (see results in appendix A.1).

5 CONCLUSION

In this work, we propose Single-Task Poison (STP) to investigate data poisoning attacks in Continual Learning. In contrast to previously proposed poisoning settings, in STP adversaries do not have knowledge and access to a model, previous and future tasks data, having access only to a single task in the data stream. Our study demonstrates that even within these stringent conditions adversaries can compromise model performance by poisoning data with basic image corruptions. We show that STP attacks strongly disrupt the training: decrease the performance for the past tasks and impede training on future tasks. Finally, we propose a high-level defense framework for CL and evaluate a poison task detection method based on task vectors. Our research reveals CL vulnerabilities to data poisoning attacks and supports further investigation of such threats to address security challenges present in continual learning with adequate defense methods.

Acknowledgements. This research was supported by Warsaw University of Technology (Poland) within the Excellence Initiative Research University (IDUB) programme. We additionally acknowledge projects PID2022-143257NB-I00, financed by MCIN/AEI/10.13039/501100011033 and FSE+, funding by the European Union ELLIOT project, and the Generalitat de Catalunya CERCA Program. Bartłomiej Twardowski acknowledges the grant RYC2021-032765-I and National Centre of Science (NCN, Poland) Grant No. 2023/51/D/ST6/02846. Stanisław Pawlak acknowledges National Centre of Science (NCN, Poland) Grant No. 2023/51/D/ST6/01609.

REFERENCES

- Ali Abbasi, Parsa Nooralinejad, Hamed Pirsiavash, and Soheil Kolouri. Brainwash: A poisoning attack to forget in continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 24057–24067, June 2024.
- Krizhevsky Alex. Learning multiple layers of features from tiny images. <https://www.cs.toronto.edu/kriz/learning-features-2009-TR.pdf>, 2009.
- Antonio Emanuele Cinà, Kathrin Grosse, Ambra Demontis, Sebastiano Vascon, Werner Zellinger, Bernhard A. Moser, Alina Oprea, Battista Biggio, Marcello Pelillo, and Fabio Roli. Wild patterns reloaded: A survey of machine learning security against training data poisoning. *ACM Computing Surveys*, 55(13s):1–39, July 2023. ISSN 1557-7341. doi: 10.1145/3585385. URL <http://dx.doi.org/10.1145/3585385>.
- Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2021.
- Prithviraj Dhar, Rajat Vikram Singh, Kuan-Chuan Peng, Ziyang Wu, and Rama Chellappa. Learning without memorizing, 2019.
- Christopher Frederickson, Michael Moore, Glenn Dawson, and Robi Polikar. Attack strength vs. detectability dilemma in adversarial machine learning, 2018.
- Gyojin Han, Jaehyun Choi, Hyeong Gwon Hong, and Junmo Kim. Data poisoning attack aiming the vulnerability of continual learning. In *2023 IEEE International Conference on Image Processing (ICIP)*, pp. 1905–1909. IEEE, 2023.
- Jonathan Hayase, Weihao Kong, Raghav Somani, and Sewoong Oh. Spectre: Defending against backdoor attacks using robust statistics. In *International Conference on Machine Learning*, pp. 4129–4139. PMLR, 2021.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.
- Sanghyun Hong, Varun Chandrasekaran, Yiğitcan Kaya, Tudor Dumitraş, and Nicolas Papernot. On the effectiveness of mitigating data poisoning attacks with gradient shaping. *arXiv preprint arXiv:2002.11497*, 2020.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic, 2023. URL <https://arxiv.org/abs/2212.04089>.
- Siteng Kang, Zhan Shi, and Xinhua Zhang. Poisoning generative replay in continual learning to promote forgetting. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 15769–15785. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/kang23c.html>.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- Huayu Li and Gregory Ditzler. Targeted data poisoning attacks against continual learning neural networks. In *2022 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2022. doi: 10.1109/IJCNN55064.2022.9892774.
- Huayu Li and Gregory Ditzler. Pacol: Poisoning attacks against continual learners, 2023.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- Marc Masana, Xialei Liu, Bartłomiej Twardowski, Mikel Menta, Andrew D Bagdanov, and Joost Van De Weijer. Class-incremental learning: survey and performance evaluation on image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5513–5533, 2022.

- Luis Muñoz-González, Battista Biggio, Ambra Demontis, Andrea Paudice, Vasin Wongrassamee, Emil C Lupu, and Fabio Roli. Towards poisoning of deep learning algorithms with back-gradient optimization. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pp. 27–38, 2017.
- Tuan Anh Nguyen and Anh Tran. Input-aware dynamic backdoor attack. *Advances in Neural Information Processing Systems*, 33:3454–3464, 2020.
- Grégoire Petit, Adrian Popescu, Hugo Schindler, David Picard, and Bertrand Delezoide. Fetril: Feature translation for exemplar-free class-incremental learning, 2023. URL <https://arxiv.org/abs/2211.13131>.
- Ameya Prabhu, Philip Torr, and Puneet Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *The European Conference on Computer Vision (ECCV)*, August 2020.
- Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. *Advances in neural information processing systems*, 31, 2018.
- Rahim Taheri, Reza Javidan, Mohammad Shojafar, Zahra Pooranian, Ali Miri, and Mauro Conti. On defending against label flipping attacks on malware detection systems. *Neural Computing and Applications*, 32:14781–14800, 2020.
- Muhammad Umer and Robi Polikar. False memory formation in continual learners through imperceptible backdoor trigger. *arXiv preprint arXiv:2202.04479*, 2022.
- Muhammad Umer and Robi Polikar. Adversary aware continual learning, 2023.
- Muhammad Umer, Glenn Dawson, and Robi Polikar. Targeted forgetting and false memory formation in continual learners through adversarial backdoor attacks. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8. IEEE, 2020.
- Gido M. van de Ven and Andreas S. Tolias. Three scenarios for continual learning, 2019.
- Shuaiqi Wang, Jonathan Hayase, Giulia Fanti, and Sewoong Oh. Towards a defense against backdoor attacks in continual federated learning. *arXiv preprint arXiv:2205.11736*, 2022.
- Chaofei Yang, Qing Wu, Hai Li, and Yiran Chen. Generative poisoning attack method against neural networks. *arXiv preprint arXiv:1703.01340*, 2017.
- Yu Yang, Tian Yu Liu, and Baharan Mirzasoleiman. Not all poisons are created equal: Robust training against data poisoning. In *International Conference on Machine Learning*, pp. 25154–25165. PMLR, 2022.
- Liuwan Zhu, Rui Ning, Chunsheng Xin, Chonggang Wang, and Hongyi Wu. Clear: Clean-up sample-targeted backdoor in neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16453–16462, 2021.

6 APPENDIX

A ADDITIONAL EXPERIMENTS

A.1 EXTENDED INVESTIGATION OF STP ATTACKS

We state the following questions to further investigate STP attacks:

7. How plasticity-stability regularization affect the outcome of the attack? Which task in the sequence is the easiest to poison?
8. How does connection between different severities of the corruption and different percentage of poisoned exemplars (pp) affect STP attack strength?
9. Are STP attacks successful for a stream containing more complex images than those in CIFAR10/100?
10. Are STP attacks successful when using memory-based CL methods?

We answer them with the following observations:

7A) Regularization strength matters in STP attacks. Increasing the regularization in CL improves model performance for past tasks. Thus, it should also affect the outcome of the STP attack. Figure 12 shows that increasing regularization increase also the poisoning impact on past tasks, while undermines the impact of poisoning for future tasks performance. The relative performance drop for past tasks increases from -3.1 at $\lambda = 1$ to -7.1 at $\lambda = 10$. Conversely, for future tasks, the relative performance drop decreases from -12.6 at $\lambda = 1$ to -7.4 at $\lambda = 10$. This result is consistent with the fact, that the ability to poison the model is connected with the plasticity of the model. For indiscriminate attacks, the lower the plasticity, the lower is the data poisoning risk. In an extreme case, when the plasticity is equal to zero (model is not updated, the weights are fixed), poisoning the data has no effect.

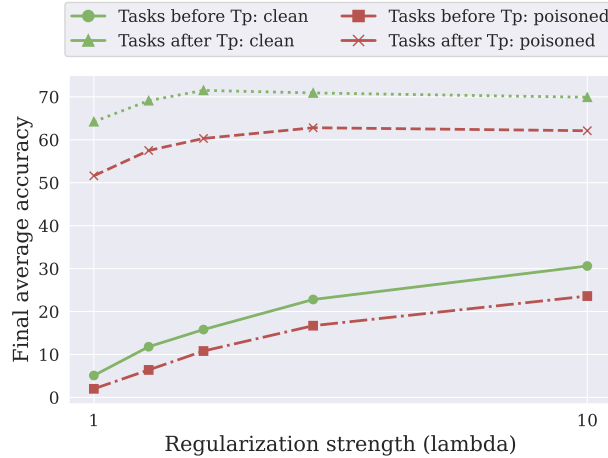


Figure 12: **Regularization strength matters in STP attack.** The stronger regularization, the smaller is the impact of poisoning for future tasks performance, while bigger for past tasks.

7B) Poisoning first task in a sequence has the most severe consequences. Figure 13 shows average accuracies for 5 tasks in CL split-CIFAR10 sequence with BAIT attacks on T_p , where $p \in \{1, 5\}$. When changing the poisoned task index p , we can observe that T_p is the most affected task, consistently with our previous claim in observation 2A. Moreover, poisoning the first task in the sequence is far the most influential for further training: the model trained on poisoned data from the start struggles to learn significant patterns later on. While interesting, the significance of this fact is not crucial: the defender party has much bigger control on the first task or model initialization than any other task in the sequence, so additional defensive mechanism can be used to eliminate this threat.

8) Poisoning a vast percentage of class exemplars is more important than corruption severity. Figure 14 is an extension of Figure 9 presented in the main body of this work. It shows (in accordance to observations stated in 4)) that two conditions are necessary for a successful STP attack. Firstly, the corruption of the image itself cannot be

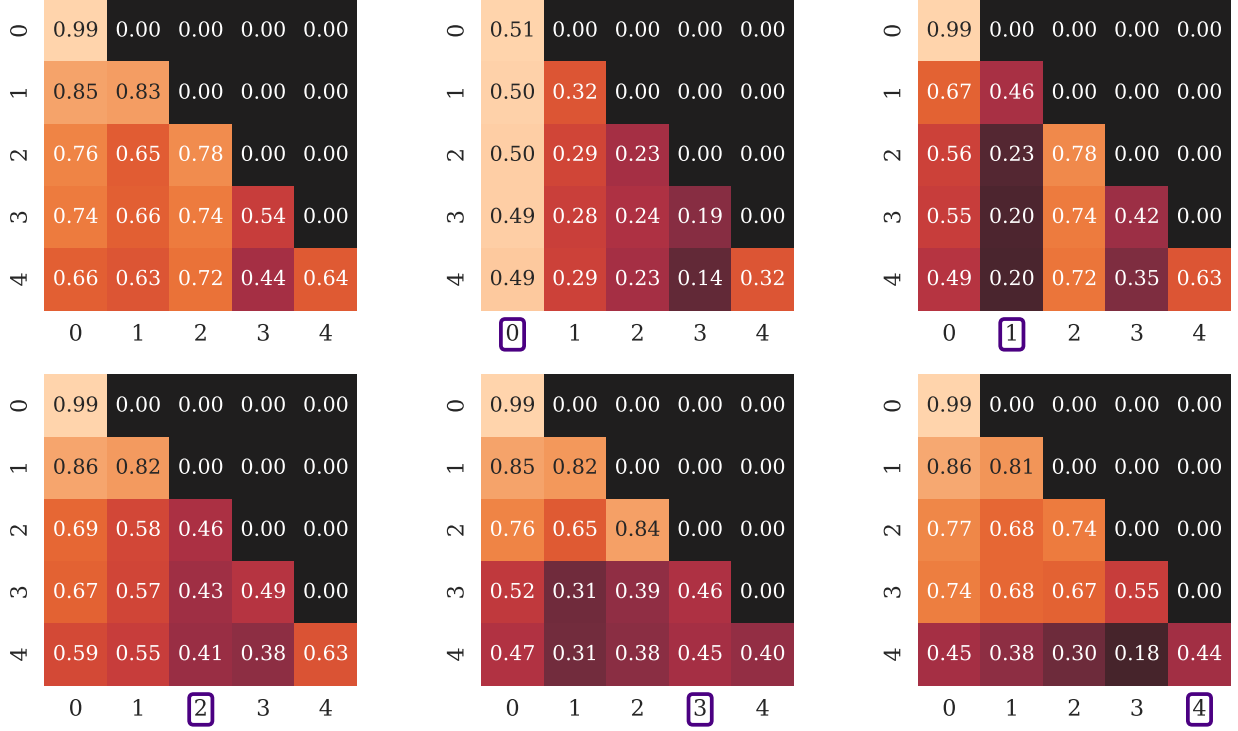


Figure 13: **Poisoning different tasks on CIFAR10: BAIT attack with gaussian blur corruption.** Top left image is LwF class incremental accuracy on clean data. Other plots have poisoned T_p number marked with purple frame.

Table 3: **Average accuracy when poisoning a single T_p task in CL sequence.** CLEAN row shows results of runs without poisoning T_p task, while black (BASE, BAIT, MULTIBASE, MULTIBAIT) rows show results for different attacks. Difference between CLEAN and poisoned runs is shown in parentheses. We report accuracies after the poisoned task (ACC AT POISONING TIME). Additionally, we report accuracies at the end of the full CL sequence (FINAL ACC) to show the impact of poisoning for further training.

ATTACK	METHOD	DATASET	P	ACC AT POISONING TIME		FINAL ACC			
				T_p	BEFORE T_p	T_p	BEFORE T_p	AFTER T_p	TOTAL
CLEAN	LWF	TINY	4	75.8	39.8	47.3	23.0	65.9	41.3
BASE	LWF	TINY	4	43.5 (-32.3)	21.2 (-18.6)	23.3 (-24.0)	12.5 (-10.5)	58.6 (-7.3)	29.7 (-11.6)
BAIT	LWF	TINY	4	45.2 (-30.6)	32.5 (-7.3)	33.9 (-13.4)	17.5 (-5.5)	62.0 (-3.9)	35.1 (-6.2)
MULTIBASE	LWF	TINY	4	64.2 (-11.6)	39.1 (-0.7)	41.3 (-6.0)	21.9 (-1.1)	64.6 (-1.3)	39.4 (-1.9)
MULTIBAIT	LWF	TINY	4	65.8 (-10.0)	39.8 (0.0)	38.4 (-8.9)	21.6 (-1.4)	59.2 (-6.7)	36.9 (-4.4)

negligible. Very low severity does not disrupt the training (4A). Secondly, the vast percentage of T_p must be poisoned (4B). From those two conditions, the latter has greater significance. When poisoning $\geq 90\%$ of class exemplars in T_p performance decrease even with smaller severities of corruption. On the other hand, when less exemplars are poisoned the accuracy drop is negligible even for the strongest severity considered.

9A) STP attacks are successful for streams with more complex images. To confirm that STP attacks pose a threat also for more complex CL benchmarks, we conduct experiments on Split Tiny Imagenet. Table 3 shows similar results to those presented in main part of this work (see Table 2): in terms of FINAL ACC, the most affected is accuracy on poisoned task (T_p), but there exists also significant drop for *BEFORE T_p* and *AFTER T_p* results. The results further reinforce the evidence supporting the effectiveness of STP attacks and validate their impact on more complex datasets beyond CIFAR-10/100.

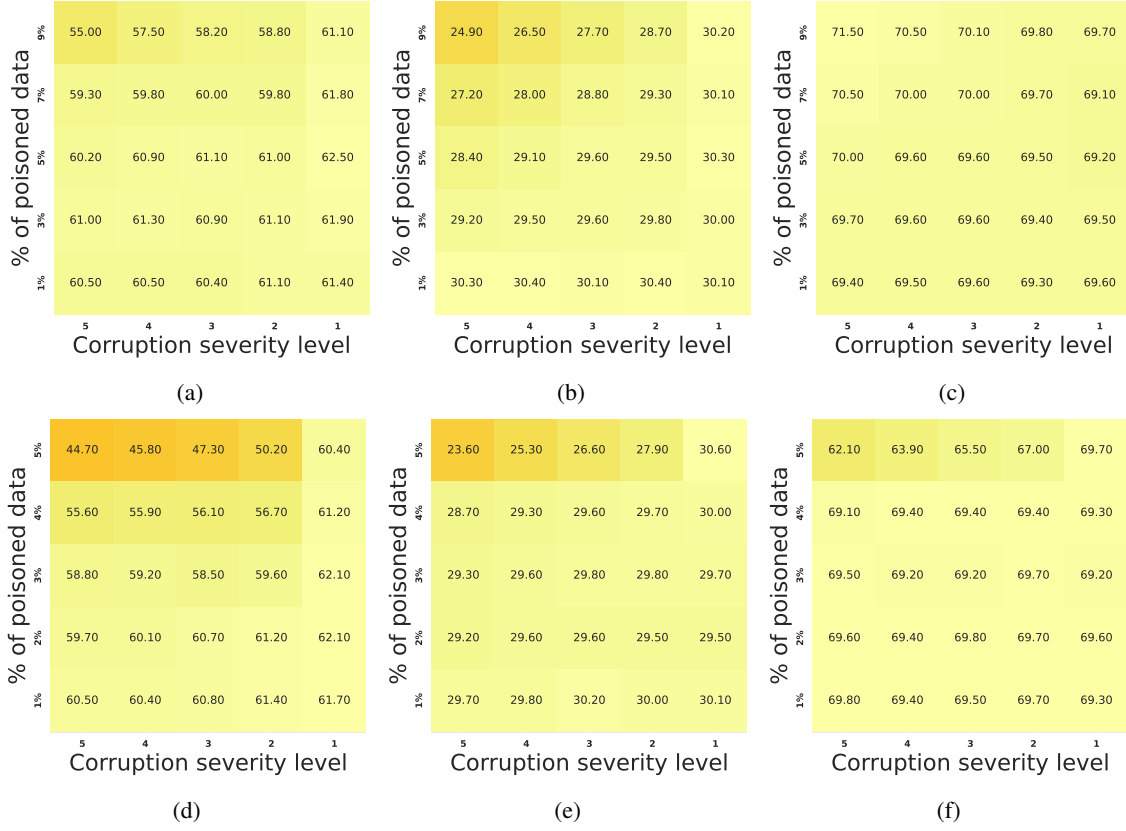


Figure 14: **The impact of corruption severity and poisoned percentage of samples on STP attack performance.** Poisoning a vast percentage of class exemplars is more important than corruption severity. We report the average accuracy for: T_p , BEFORE T_p and AFTER T_p for BASE attack (a-c) and BAIT attack (d-f). Note that the % on y axis is calculated as % of samples poisoned in the whole dataset and 10% means that all samples in T_p are poisoned.

Table 4: **Poison task detection for BASE attacks.** P90 achieves the highest overall accuracy and the top F1 score. THRESHOLD – statistic used to calculate α , ACC – total detection accuracy (on clean and poisoned tasks), CLEAN – accuracy on clean tasks, ATTACK – accuracy on poisoned tasks, F1 – F1 score.

THRESHOLD	ACC	CLEAN	ATTACK	F1
MAX+IQR	93.8	100.0	80.0	0.86
MAX	92.3	97.8	80.0	0.86
P90	93.8	97.8	85.0	0.89
MAX-IQR	92.3	91.1	95.0	0.88
P75	84.6	80.0	95.0	0.79

9B) BASE attacks are stronger with increasing number of classes in T_p and increasing classification complexity. Table 3 shows strong performance drop for BASE attack. We attribute it to two following factors: 1) the number of classes in T_p , 2) total number of classes (i.e. classification task complexity). Regarding the TOTAL FINAL ACC, the decreases become more pronounced as the number of classes and dataset complexity increase: -2.7 (with T_p having two classes for CIFAR-10), -9.7 (with T_p having 10 classes for CIFAR-100), and -11.6 (with T_p having 20 classes for Tiny ImageNet).

10) Performance on past and future tasks is not affected for STP attacks on memory-based CL methods. In the Limitations section of this work, we note that the STP framework is effective for investigating memory-free CL methods but less applicable to those utilizing exemplars. The STP framework restricts the adversary to poisoning only a single task. When only one task is poisoned, CL methods with memory maintain an unrealistic guarantee that

Table 5: **Poison task detection for BAIT attacks.** P90 achieves the highest overall accuracy and the top F1 score. THRESHOLD – statistic used to calculate α , ACC – total detection accuracy (on clean and poisoned tasks), CLEAN – accuracy on clean tasks, ATTACK – accuracy on poisoned tasks, F1 – F1 score.

THRESHOLD	ACC	CLEAN	ATTACK	F1
MAX+IQR	75.4	100.0	20.0	0.10
MAX	86.2	97.8	60.0	0.73
P90	89.2	97.8	70.0	0.80
MAX-IQR	84.6	91.1	70.0	0.74
P75	80.0	80.0	80.0	0.71

Table 6: **Average accuracy when poisoning a single T_p task in CL sequence with memory buffer of 2000 exemplars.** CLEAN rows show results of runs without poisoning T_p task, while black (BASE, BAIT) rows show results for different attacks. Difference between CLEAN and poisoned runs is shown in parentheses. We report accuracies after the poisoned task (ACC AT POISONING TIME). Additionally, we report accuracies at the end of the full CL sequence (FINAL ACC) to show the impact of poisoning for further training.

ATTACK	METHOD	DATASET	P	ACC AT POISONING TIME		FINAL ACC			
				T_p	BEFORE T_p	T_p	BEFORE T_p	AFTER T_p	TOTAL
CLEAN	FINETUNING	CIFAR10	3	97.6	84.5	74.8	73.9	86.4	79.1
BASE	FINETUNING	CIFAR10	3	0.1 (-97.5)	93 (8.5)	0.0 (-74.8)	74.7 (0.8)	86.9 (0.5)	64.6 (-14.5)
BAIT	FINETUNING	CIFAR10	3	49.0 (-48.6)	87.2 (2.7)	35.2 (-39.6)	73.8(-0.1)	86.9 (0.5)	71.3 (-7.8)
CLEAN	FINETUNING	CIFAR100	4	90.4	44.9	42.9	37.5	64.3	47.4
BASE	FINETUNING	CIFAR100	4	0.1 (-90.3)	56.9 (12.0)	0.0 (-42.9)	38.1 (0.6)	64.7 (0.4)	40.6 (-6.8)
BAIT	FINETUNING	CIFAR100	4	45.8 (-44.6)	50.2 (5.3)	20.5 (-22.4)	37.8(0.3)	64.2 (-0.1)	43.7 (-3.7)
CLEAN	GDUMB	CIFAR10	3	76.3	78.7	68.0	66.5	66.7	66.9
BASE	GDUMB	CIFAR10	3	0.7 (-75.6)	83.9 (5.3)	0.9 (-67.1)	71.4 (4.9)	69.1 (2.4)	56.4 (-6.8)
BAIT	GDUMB	CIFAR10	3	39.0 (-37.3)	80.7 (2.0)	34.4 (-33.6)	68.8(2.3)	67.1 (0.4)	61.3(-5.6)
CLEAN	GDUMB	CIFAR100	4	23.8	26.0	21.5	21.8	17.8	20.4
BASE	GDUMB	CIFAR100	4	0.2 (-23.6)	28.4 (2.4)	0.5(-21.0)	23.0 (1.2)	19.5 (1.7)	18.1 (-2.3)
BAIT	GDUMB	CIFAR100	4	13.9 (-9.9)	27.3 (1.3)	12.1 (-9.4)	22.9(1.1)	17.9 (0.1)	19.5 (-0.9)

all previous samples stored in the buffer are clean during the T_p task (since T_p is the only poisoned task). Introducing an additional buffer of clean samples undermines the concept of ‘learning the task in isolation’ and mitigates the STP attack, thereby creating a scenario more akin to joint training. Results for simple finetuning with buffer and GDUMB (Prabhu et al. (2020)) presented in Table 6 confirm that (similarly to JOINT training) attack decrease only performance on T_p , while there are no accuracy drops for *BEFORE T_p* and *AFTER T_p* results (numbers in parenthesis are positive or close to zero).

Additional results for LwM. Table 7 presents STP attacks evaluation for additional data regularization method – LwM (Dhar et al. (2019)). These results are consistent with observations described in the main part of this work.

A.2 EXTENDED INVESTIGATION OF DEFENSE AGAINST STP ATTACKS.

The threshold value that determines whether a task is poisoned plays a crucial role in our method. We establish this threshold by using a clean task at the beginning of the training sequence. In order to obtain a threshold value we follow two simple steps: 1) We calculate the angle α' (see Figure 5) on single clean task multiple times (using different random seeds for initialization), 2) We use a statistic measure to consolidate multiple possible threshold values into one α . Aggregated threshold value is then applied to all subsequent tasks in the stream to assess whether they are poisoned. The choice of an appropriate threshold involves balancing the trade-off between correctly identifying attacks

Table 7: **Average accuracy when poisoning a single T_p task in CL sequence.** CLEAN rows show results of runs without poisoning T_p task, while black (BASE, BAIT, MULTIBASE, MULTIBAIT) rows show results for different attacks. Difference between CLEAN and poisoned runs is shown in parentheses. We report accuracies after the poisoned task (ACC AT POISONING TIME). Additionally, we report accuracies at the end of the full CL sequence (FINAL ACC) to show the impact of poisoning for further training.

ATTACK	METHOD	DATASET	P	ACC AT POISONING TIME		FINAL ACC			
				T_p	BEFORE T_p	T_p	BEFORE T_p	AFTER T_p	TOTAL
CLEAN	LWM	CIFAR10	3	67.8	71.0	51.4	62.6	62.8	60.4
BASE	LWM	CIFAR10	3	55.6 (-12.2)	69.8 (-1.2)	40.0 (-11.4)	58.2 (-4.4)	64.5 (1.7)	57.1 (-3.3)
BAIT	LWM	CIFAR10	3	45.5 (-22.3)	50.5 (-20.5)	41.8 (-9.6)	45.8 (-16.8)	50.2 (-12.6)	46.8 (-13.6)
CLEAN	LWM	CIFAR100	4	79.0	51.1	62.2	31.1	70.1	49.3
BASE	LWM	CIFAR100	4	64.4 (-14.6)	37.7 (-13.4)	46.7 (-15.5)	18.1 (-13.0)	69.4 (-0.7)	40.0 (-9.3)
BAIT	LWM	CIFAR100	4	52.3 (-26.7)	40.7 (-10.4)	45.1 (-17.1)	23.6 (-7.5)	62.1 (-8.0)	40.0 (-9.3)
MULTIBASE	LWM	CIFAR100	4	36.6 (-42.4)	45.9 (-5.2)	28.6 (-33.6)	26.8 (-4.3)	65.2 (-4.9)	39.9 (-9.4)
MULTIBAIT	LWM	CIFAR100	4	54.5 (-24.5)	47.5 (-3.6)	47.2 (-15.0)	28.4 (-2.7)	65.7 (-4.4)	44.0 (-5.3)

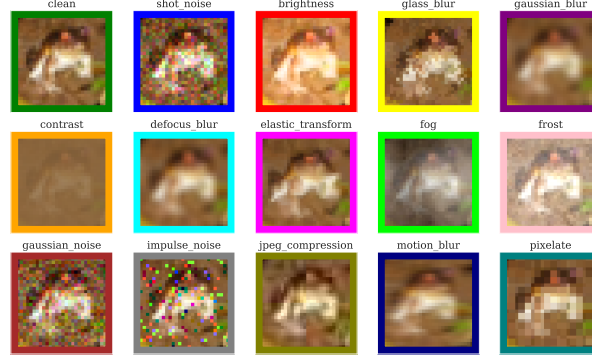


Figure 15: **Image corruptions used to create poisoned training samples.** We use CIFAR-C set of image corruptions to transform images of the poisoned task. We show clean image in the upper right corner for comparison. We present corruptions with maximum severity (equal to 5).

(True Positives) and minimizing false alarms when a task is not poisoned (False Positives). In the following section, we evaluate five different potential thresholds.

Evaluation of different threshold values We evaluate the following statistics to aggregate multiple α' values into single threshold α :

- MAX+IQR – We calculate Interquartile Range for α' values and add it to maximum α' value. This threshold is optimized for the lowest number of False Positives (false alarms).
- MAX – maximum α' value.
- P90 – We calculate 90th percentile of α' values.
- MAX-IQR – We calculate Interquartile Range for α' values and subtract it from maximum α' value.
- P75 – We calculate 75th percentile (Q3) of α' values. This threshold is optimized for the highest number of True Positives (detected attacks).

Tables 4 and 5 present the evaluation of the five proposed thresholds. The extreme options optimize performance for either clean data (MAX+IQR) or poisoned data (P75). While the choice of threshold ultimately depends on the specific needs of the defending party, metrics such as the F1 score can be used to achieve an appropriate balance between Precision and Recall. In our experiments, which considered five possible thresholds, P90 yielded the highest overall accuracy and the best F1 score, making it the most effective threshold among those evaluated.

B EXPERIMENTAL SETUP DETAILS

General comments. We use FACIL framework (Masana et al. (2022)) to conduct CL experiments for all methods. We report results averaged between five runs with different random seeds. All experiments and evaluations were performed with Nvidia RTX A5000 GPUs.

Corruptions Figure 15 shows a set of CIFAR-C based corruptions used to poison data in this work. While CIFAR-C dataset is usually used for test-time adaptation scenarios and thus benchmark consists of corrupted test data, we use those predefined transformations to create corrupted images for training the poisoned task. For all results with single poison reported in tables we use *gaussian_blur* corruption (see Figure 2). For attacks with higher number of poisons (MULTIBASE and MULTIBAIT) we use five poisons: *gaussian_blur*, *gaussian_noise*, *contrast*, *pixelate*, *jpeg_compression*. We show attacks using all kinds of corruptions (each as a single poison for BASE or BAIT attacks) on Figure 7.

CL setup. For additional experiments on Tiny Imagenet we split data into six tasks, with 100 classes in the first task, and 20 classes in each following task. Consequently, during our experiments poisoned task (T_p) consist of 2 classes for CIFAR10, ten classes for CIFAR100 and 20 classes for Tiny Imagenet.

Models. Resnet32 (used for CIFAR10/100 experiments) and Resnet18 (used for Tiny Imagenet experiments) implementations are taken from FACIL.

Additional parameters. We use $\lambda=10$ (distillation loss parameter) for both LWM and LwF methods. For additional experiments on LwM (see Table 7) we use additionally $\gamma=20$ (attention loss parameter). We use data augmentations from FeTrIL (Petit et al. (2023)).

Defense evaluation. Defense evaluation is done using CIFAR100 dataset. For threshold selection procedure we use second task (to fulfill the requirement of task in the beginning of the stream with the same number of classes as potential poisoned tasks). For poisoned tasks we use 4th task to be consistent with attack evaluation section. As clean tasks we use all 10-class tasks from CIFAR100 splits using multiple random seeds.

C BROADEN IMPACT

Our research reveals CL vulnerabilities to data poisoning attacks. To address possible negative impact and minimize the risk of misusing our work we investigate and discuss possible defensive measures. We also present a proposition for defense framework for CL with poison task detection method. Finally, our investigations using STP framework aim to facilitate closing a gap existing between the number of proposed attacks and defenses.