

Teacher-Forced Selective Self-Distillation for Uncurated Replay Data

Anonymous ACL submission

Abstract

Continual fine-tuning involves incrementally training a language model to acquire knowledge of new tasks. This learning paradigm introduces the challenge of catastrophic forgetting, where models tend to forget previously learned tasks as they adapt to new ones. Several techniques have been proposed to address this issue, including regularization, parameter-isolation, and replay-based approaches. Among these, replay-based methods have gained wider adoption due to their less invasive nature and ease of integration into existing continual learning pipelines. However, in real-world settings, curating ideal replay samples is a practical challenge for replay-based methods. This leads to the use of noisy replay data (either synthetic generated from exemplars or provided by task owner), which results in suboptimal task performance. To address this crucial real-world challenge, we introduce Teacher-Forced Selective Self-Distillation (TF-SSD) a novel method that employs self-distillation of the labels from the task stage model and refine the less effective samples using mixture of teachers framework. Our experiments involving challenging 16 task continual learning setting demonstrate that TF-SSD outperforms best-performing baseline by ~ 2.7 points in task performance and ~ 2.8 points in mitigating catastrophic forgetting across 2 model families: Llama2 7B and Granite3.3 2B. We plan to open-source the code of TF-SSD.

1 Introduction

Large Language Models (LLMs) have demonstrated promising performance (Dubey et al., 2024; Brown et al., 2020) across a wide range of Natural Language Processing (NLP) tasks. In practice, LLMs must be continually updated to remain effective on newly emerging tasks. To facilitate this, Continual Learning (CL) (Chen and Liu, 2018) is employed—a paradigm in which models are incrementally trained on incoming data, enabling them

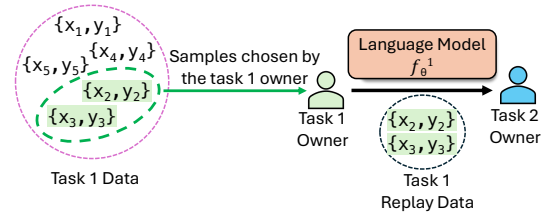


Figure 1: Demonstration of the problem setting.

to remain relevant and effective with time. Continual Fine-Tuning (CFT), a variant of CL, involves incrementally training models on new fine-tuning task data. However, this paradigm introduces a major challenge known as Catastrophic Forgetting (CF) (McCloskey and Cohen, 1989), where performance on previously learned tasks deteriorates as new tasks are acquired.

Numerous techniques have been proposed to mitigate CF, including regularization-based approaches (Aljundi et al., 2018; Kirkpatrick et al., 2017; Li and Hoiem, 2018), architecture-level methods involving task-specific parameter isolation (Wang et al., 2023b; Ke et al., 2021; Satapara and Srijith, 2024; Wang et al., 2023a; Razdaibiedina et al., 2023; Wang et al., 2024a), and data replay strategies that combine previous and current task data during training (Chaudhry et al., 2019; Rebuffi et al., 2017; Buzzega et al., 2020; He et al., 2024; M’hamdi and May, 2024; Lopez-Paz and Ranzato, 2017; Huang et al., 2024). Among these, replay-based methods have gained broader adoption due to their ease of integration into existing CL pipelines and their less invasive nature compared to architectural or regularization-based methods.

While replay-based approaches are less invasive, they typically require storing a subset of past task data. In real-world settings, often practiced and well acknowledged in industry (Google, 2022; NVIDIA, 2020), future task owners do not have control over the replay data stored by past task owners, we term such replay data as *uncurated*. As shown in Figure 1, the task 1 owner just exposes

a small part of the data (while keeping most of it private) to help task 2 owner to preserve task 1 performance. It is then upto the task 2 owner to effectively curate the available replay data to retain task 1 performance. A notable example is the Nvidia Nemotron study (NVIDIA, 2020) which released a small dataset (Nemotron-Pretraining-Dataset-sample) with representative samples ranging from 1k to 15k across different domains. The same study also has released task-centric datasets (Nemotron-Pretraining-SFT-v1) with limited samples for various tasks.

In this paper, we introduce TF-SSD (see Figure 2), a novel replay-based CFT method designed to address this challenge. TF-SSD aims to curate the available replay data. Our method distills labels for each task from its corresponding stage model within the CFT pipeline. These samples are subsequently refined: noisy samples selectively undergo teacher-forcing, where teacher (T) tokens are derived from the responses of an ensemble of open-source LLMs. The number of teacher tokens used is dynamically chosen based on the degree of defectiveness in the sample.

Our main contributions are:

1. We propose a novel replay-based method for continual learning that is less invasive than existing approaches, effectively addressing a real-world setting with uncurated replay data. Key features involve: (1) Label Distillation: distillation of labels from respective task stage model. (2) Refinement: teacher-forced selective self-distillation based on the degree of defectiveness in the sample with teacher tokens from ensemble of teachers.
2. We evaluate our method on the SuperNI multi-task dataset (Wang et al., 2022), comprising 48 subtasks across 16 tasks encompassing code, language and mathematics, using two model families: Granite3.3 (Mishra et al., 2024) and Llama2 (Touvron et al., 2023).
3. We conduct extensive analysis and ablation studies to validate the effectiveness of TF-SSD.

2 Related Work

Consistent with recent CL studies (Wang et al., 2024a; Chen and Zeng, 2025), we classify CL strategies into three distinct categories.

Notation	Description
n	Denotes number of stages in CFT pipeline. Since, each stage is associated with a task, n denotes number of tasks as well
$\{1, \dots, n\}$	Series of tasks in CFT pipeline
f_{θ}^0	Base model with parameters θ
f_{θ}^i	Model finetuned at stage i
d^i	Training data for task i
r^i	Replay data for task i from task owner i
r_p^i	Replay data from TF-SSD for task i
$r_{\theta(i)}^i$	Replay data with labels distilled from model f_{θ}^i
v	Scorer module
s	A sample in the dataset
x	Input text in the sample s
y	Label in the sample s
y_l	Label with l number of tokens
m	Number of teachers
$y^1 \dots y^m$	Intermediate labels generated by m teachers
y^{MoT}	Label generated by Mixture of Teachers precisely coming from the Aggregator
y_k^{MoT}	First k tokens of y^{MoT}
y^{TF}	Label generated by teacher-forcing the stage model
t_{θ}^i	i^{th} teacher model

Table 1: Various notations with descriptions used in TF-SSD.

Regularization. These methods (Aljundi et al., 2018; Kirkpatrick et al., 2017; Li and Hoiem, 2018) aim to mitigate forgetting of previous tasks by adding regularization terms to the loss function. However, the inclusion of multiple regularization terms may degrade model performance (Parisi et al., 2019). To address this, some works (Mok et al., 2023) combine regularization with other approaches.

Architecture. Some approaches (Wang et al., 2023b) isolate task-specific parameters via parameter isolation, while others adopt parameter-efficient fine-tuning techniques by introducing task-specific modules (Ke et al., 2021; Satapara and Srijith, 2024; Wang et al., 2023a), learning soft prompts (Razdaibiedina et al., 2023), or applying prefix tuning (Wang et al., 2024a) for newly arriving tasks.

Replay. Early methods (Chaudhry et al., 2019; Rebuffi et al., 2017) tackle forgetting by replaying a portion of past data. Later work (Rolnick et al., 2019; Buzzega et al., 2020), improved this by generating replay data with distilled labels from previous models. Further, Huang et al. (2024) synthesizing samples using the base model and refining them with stage model. He et al. (2024) introduced an architecture-level approach that distills attention from selected attention heads instead of labels. Additionally, some studies (M’hamdi and May, 2024; Lopez-Paz and Ranzato, 2017) propose techniques

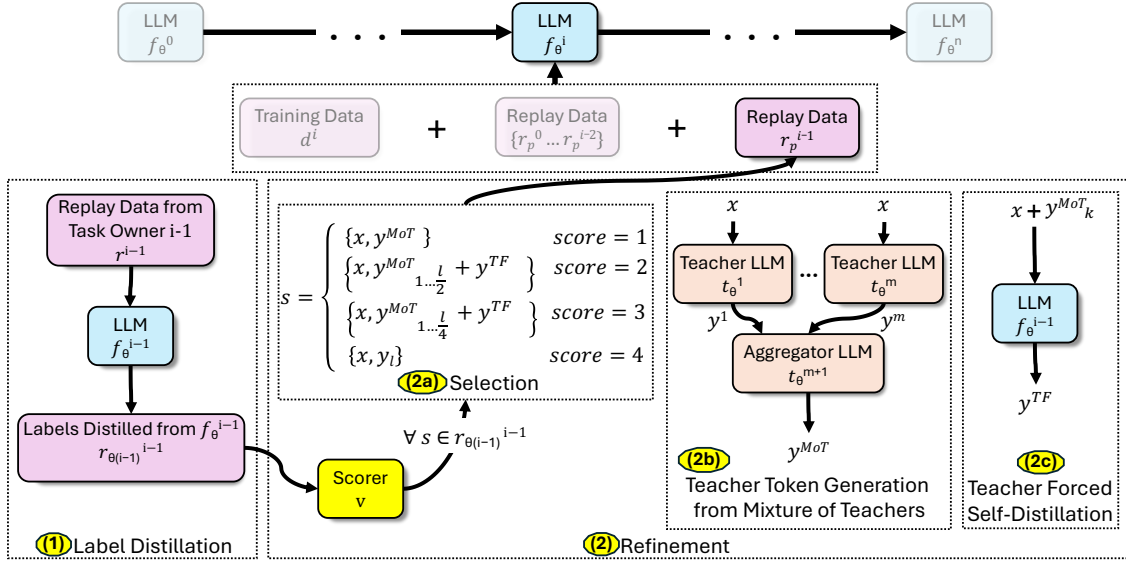


Figure 2: Overview of our method TF-SSD. Some concrete examples are provided in Tables 13 and 14 in Appendix.

for selectively sampling candidates to populate a fixed-size replay data.

Replay techniques are generally less invasive than architecture-based methods. TF-SSD falls under replay-based category and we attempt to study the crucial challenge of minimizing CF in a setting where replay data from the task owner is uncurated (see Figure 1) and yet essential for maintaining task performance. We follow (Scialom et al., 2022) setting, where CL is addressed from a broader view, enabling the model to learn diverse set of new tasks. We frame any NLP task (classification, summarization etc.) with input and output in natural language (NL). Thus, not limited to a narrow domain or a specific task CL (Biesialska et al., 2020). Since our approach belongs to the replay methods, we include prior work - original data as replay (Chaudhry et al., 2019), data distilled from a model as replay (Rolnick et al., 2019; Buzzega et al., 2020) and self-synthesis and refinement work (Huang et al., 2024). We exclude replay-based methods that require access to the model architecture (He et al., 2024) or rely on curating (M’hamdi and May, 2024; Lopez-Paz and Ranzato, 2017) the replay data from training data.

3 Method

3.1 Notations

Notations are provided in Table 1 for clarity.

3.2 Problem Formulation.

CFT aims to incrementally train the model f_θ^i over a sequence of fine-tuning tasks $\{1, \dots, n\}$, while

mitigating CF and maximizing performance across all tasks in a task properties-agnostic setting. We use the widely adopted decoder class of language models with a causal language modeling loss defined as:

$$\mathcal{L}(x, y; \theta) = - \sum_{t=1}^T \log P_\theta(y_t | x, y_{<t}) \quad (1)$$

The objective is to minimize the negative log-likelihood of each label token y_t , conditioned on the input x and the previous label tokens $y_{<t}$. Where T is total number of tokens in label y .

3.3 Overview.

TF-SSD is a replay-based method that aims to generate a refined version (r_p^i) of the uncurated replay data (r^i) provided by the task owner, aiming to mitigate CF and improving overall performance across tasks. This is achieved through two key submethods: (1) Label Distillation ((1) in Figure 2), and (2) Refinement ((2) in Figure 2). The Refinement submethod further incorporates sample scoring and score-based selection ((2a) in Figure 2), teacher forced self-distillation ((2c) in Figure 2) using tokens sampled from a mixture of teacher (MoT) models ((2b) in Figure 2).

3.4 Label Distillation

The training data for fine-tuning the model f_θ^i is composed of current task data d^i , replay data from past stages $\{r_p^0, \dots, r_p^{i-2}\}$ and the immediate previous stage r_p^{i-1} . Replay data $\{r_p^0, \dots, r_p^{i-2}\}$ is

reused from past curation while r_p^{i-1} is prepared via TF-SSD for the current stage i .

This submethod is shown as (1) in Figure 2. Task owner $i-1$ provides a portion of the uncuration task data for replay. The goal of this submethod is to generate labels y for each input x using the $i-1$ th stage model f_θ^{i-1} , thereby producing the distilled replay data $r_{\theta(i-1)}^{i-1}$, defined as:

$$y = f_\theta^{i-1}(x) \quad (2)$$

3.5 Refinement

Scoring. Each sample from $r_{\theta(i-1)}^{i-1}$ is passed through the Scorer v (see Section 4.3), which assigns a quality score in the range of 1-4. The score reflects the severity of defects in the sample’s label y , with higher scores indicate fewer defects. The score guides the degree of teacher-forcing required to improve defective samples. We adopt 1-4 scale since it covers required levels of defectiveness buckets where samples can belong to. The Scorer follows a LLM-as-a-Scorer paradigm and details on prompt template are in Appendix A.1.

Selection. Based on the score assigned by v , the sample’s label is proportionally curated ((2a) in Figure 2). If the score is 1, the full y^{MoT} is used; if 2, first half ($l/2$) of the tokens; if 3, first quarter ($l/4$); and if 4 the label is unaltered. In the cases of score 2 and 3, remaining label tokens are from y^{TF} as described later in this section.

Given the discrete nature of the score, we arrive at these rigid proportions. This is to achieve a balanced curation of label tokens from both the MoT and the stage model f_θ^{i-1} , where the proportion of MoT tokens increases with the degree of f_θ^{i-1} deficiency, as indicated by the score. A score of 2 or 3 suggests that the label generated by the stage model has partial merit but includes minor errors; thus, MoT serves as a corrective signal through prompt-level teacher-forcing. A score of 1 indicates label from f_θ^{i-1} is flawed, so y^{MoT} is fully retained. Conversely, a score of 4 implies that the stage model’s label is optimal, and is retained.

Teacher Tokens Generation through MoT. As illustrated in (2b) of Figure 2, each input x is passed through a set of Teacher models (denoted as T) to generate corresponding outputs y , following Equation (2). These outputs are then aggregated by a separate Aggregator model (denoted as A) to produce the final label y^{MoT} for each x . Details on

prompt templates used for teacher and aggregator are provided in Appendix A.1.

The paradigm of ensembling multiple open-source LLMs has been explored in latest prior work (Wang et al., 2025; Tian et al., 2025) for applications such as generating improved responses and teacher-forced chain-of-thought reasoning for smaller models. In our setup, we use MoT to selectively teacher-force label tokens by leveraging MoT tokens in the prompt (explained in the later part of the section), based on the degree of defectiveness in each sample.

Teacher-Forced Self-Distillation. This submethod is demonstrated by (2c) in Figure 2. Teacher-forcing is applied only to samples with scores 2 and 3. Each input x is augmented with k teacher tokens from y^{MoT} . Number of tokens k is explained in Selection submethod in Section 3.5. The concatenated input $x + y_k^{\text{MoT}}$ is then passed to the previous stage model f_θ^{i-1} to self-distill the remaining label tokens y^{TF} .

Through these submethods, the curated replay data r_p^{i-1} is constructed and used to train the current stage model f_θ^i .

4 Experimental Setup

4.1 Baselines

We evaluate our method against the following challenging baselines:

Base Model. The base LLM is evaluated without any additional fine-tuning. It is indicated by the corresponding model name in the results.

No Replay (NoRep). The base LLM is sequentially fine-tuned across all tasks in the CFT pipeline, without the use of any replay data.

Self-Synthesized Rehearsal (SSR) (Huang et al., 2024) Performs rehearsal from base model using replay data as few-shots with refinement using stage model.

Replay (Rep) (Chaudhry et al., 2019). The base LLM is sequentially fine-tuned with replay data consisting of all previous task replays having labels from ground-truth.

Distilled Replay from Base Model (DR_{base}) (Bhushan et al., 2025) . The base LLM is sequentially fine-tuned with replay data with labels distilled from the base model.

Task	Given a sentence, generate a new sentence by performing small changes on the sentence. Here, make sure that the changes are semantically related and syntactically similar to the input. And the generated sentence should have high commonsense plausibility, that is to have reasonable probability of it being true.		
Input	You would write a story because you have a spontaneous idea .		
Ground Truth	You would write a story because you have a detailed idea.		
	Response	ROUGE-L	LaJ
DR_{task}	You would write a story because you have a spontaneous idea.	0.91	0
TF-SSD (2 T + 1 A)	You might decide to write a story because you have a compelling idea.	0.75	2

Table 2: An example from test set illustrating the limitation of the ROUGE-L score. The top baseline DR_{task} simply repeats the input text without following the task instructions, whereas TF-SSD (2 T + 1 A) adheres to the task but is penalized by ROUGE-L due to its low overlap with the ground truth. LaJ accurately fills this gap in evaluation.

Distilled Replay from Task Stage Model (DR_{task}) (Rolnick et al., 2019; Buzzega et al., 2020). The base LLM is sequentially fine-tuned using replay data with labels distilled from the corresponding task stage model.

Multi-Task Learning (MTL). Multi-task learning involves training the model on all the tasks at once and CFT is not involved. Many works (Huang et al., 2024; He et al., 2024) consider MTL as a target topmost performance to achieve for CFT.

4.2 Ablations

We conduct the following ablation studies to isolate the contributions of MoT and teacher-forced selective self-distillation:

$DR_{task} + MoT (2 T + 1 A)^1$. The same setup as DR_{task} , with an additional refinement submethod. Each replay sample is passed through the Scorer scoring samples either 1 (accept) or 0 (reject). If rejected, the label is fully replaced with the MoT’s output; if accepted, the original label is retained.

TF-SSD (1 T)². As described in Section 3, the MoT component uses a single teacher for token generation. In other words, its an extension to $DR_{task} + MoT (2 T + 1 A)$ baseline with 1-4 scale Scorer and teacher-forced selective self-distillation of labels, instead of 0-1 scoring.

TF-SSD (2 T + 1 A). This setting includes two teachers and one aggregator model. TF-SSD and TF-SSD (2 T + 1 A) are used interchangeably.

4.3 Models

Exact sources to the models used in this paper are provided in Appendix A.3.

Models in CFT Pipeline. We choose the widely adopted Llama 2 7B (Touvron et al., 2023) and

¹Two teacher (T) models and one aggregator (A) model.

²Only one teacher (T) model and no aggregator model.

Granite 3.3 2B (Mishra et al., 2024) models to demonstrate our method on 2 diverse families and sizes. We choose the specific model versions to comply with their respective terms.

Teacher Models. We choose the models Mixtral 8x22B Instruct (MixIns) (Jiang et al., 2024) and WizardLM 2 8x22B(Wiz) (Xu et al., 2024) as teachers. Since teachers are employed in a form of synthetic data generation (SDG), we make sure to choose this specific set since we comply with their respective SDG terms. Details on prompt template used for teacher is provided in Appendix A.1.

Aggregator Model. Similarly, we use WizardLM 2 8x22B as the aggregator model. Aggregator prompt template can be found at Appendix A.1.

Scorer. We employ Prometheus 8X7B (Kim et al., 2024) as the Scorer module. Scorer prompt template can be found at Appendix A.1.

LLM-as-a-Judge. Prometheus 8X7B is used as LLM-as-a-Judge (LaJ). This model demonstrates state-of-the-art evaluation capabilities, enabling a more rigorous assessment of our method. Using the same model as both Scorer and Judge does not introduce bias, as the Scorer does not contribute data that directly influences model training, thereby leaving no scope for data leakage. To further validate this, we use DeepSeek-V3.2 (685B) as an alternative judge (Table 10). Prompt template and scoring details available in Appendix A.1.

4.4 Dataset

We use the widely adopted (Huang et al., 2024; Wang et al., 2024b) and challenging SuperNI dataset (Wang et al., 2022), which comprises 16 tasks with 3 subtasks per category. We are testing our method in 2 random orders. First order: Classification → Text Quality Evaluation → Code → Detection → Sentiment Analysis → Comprehension → Closed QA → Extraction → Dialogue → Pro-

Task	Craft one correct answer to the question given in input. To make it more interesting, try to use non-stereotypical language if possible. Make sure your correct answer is reasonably long, consistent with the context, and requires common sense (instead of explicit extraction from the context.) In your answer, use as few words as possible from the given context. Use a response that is uncommon/non-stereotypical, so that it is less predictable. To be less repetitive, please vary your language for each question.	
Context	A friend of this blog , and someone some of you know (virtually) died recently. I only received word recently and found it hard to believe . I was going to post this several days ago , but there was a contretemps on the blog and I did n't want to get that confused with this . Now that such issues are largely over , I feel able to post .	
Question	Why is it so hard to believe that someone died ?	
Ground Truth	It was sudden.	
	Response	LaJ
MTL	Because they were young and healthy.	1
TF-SSD (2 T + 1 A)	It's hard to accept that someone you knew, even if it was only virtually, is no longer alive.	2

Table 3: An example from test set is provided showcasing TF-SSD(2 T + 1 A) outperforming MTL though it is considered as a reference to topmost performance.

gram Execution → Rewriting → Open QA → Misc.
→ Generation → Summarization → Mathematics.
Second order details are provided in Appendix A.4.

4.5 Evaluation Metrics

We adopt the widely used (Huang et al., 2024) evaluation metrics for CFT. We use ROUGE-L (Lin, 2004) and LaJ score (Zheng et al., 2023) to assess the model’s task performance. Let $a_j^{(i)}$ denote the ROUGE-L score of stage i model evaluated on task j , similarly $s_j^{(i)}$ denotes LaJ score. We report scores in the following derived evaluation metrics to assess overall task performance and catastrophic forgetting in CFT. Higher values denote better performance for the all evaluation metrics.

Average ROUGE-L (AR): This metric captures the average performance of the model across all n tasks at the end of the training sequence in CFT, defined as:

$$\text{AR} = \frac{1}{n} \sum_{i=1}^n a_i^{(n)} \quad (3)$$

Average LaJ (ALaJ): Similar to AR, ALaJ is defined as:

$$\text{ALaJ} = \frac{1}{n} \sum_{i=1}^n s_i^{(n)} \quad (4)$$

Backward Transfer using ROUGE-L (BWT_r):

This metric measures how learning new tasks affects the performance on previously learned tasks. Average over all the stages except the last one where at each stage the difference between the final stage model $a_i^{(n)}$ and associated stage model $a_i^{(i)}$ task i performance is computed. A negative BWT_r value indicates CF, while positive value indicates emergent/enhanced performance. Such enhanced performance can be attributed to the replay data

used for latest tasks aiding older tasks. BWT_r is defined as:

$$BWT_r = \frac{1}{n-1} \sum_{i=1}^{n-1} (a_i^{(n)} - a_i^{(i)}) \quad (5)$$

Backward Transfer using LaJ (BWT_{LaJ}):

$$BWT_{LaJ} = \frac{1}{n-1} \sum_{i=1}^{n-1} (s_i^{(n)} - s_i^{(i)}) \quad (6)$$

Limitation of ROUGE-L. Since the tasks in our setup are highly diverse and open-ended, traditional metrics such as ROUGE-L do not fully capture (see Table 2) the quality and correctness of the generated outputs. Moreover, we observe that some samples contain defective ground truths, which can weaken evaluation reliability. Additionally, the replay data labels obtained from TF-SSD via teacher-forced selective self-distillation often include paraphrased variants. This can lead the model to generate paraphrased yet valid responses at evaluation time, which can be penalized by ROUGE-L. To address these limitations, we use LaJ to provide a more holistic and robust evaluation. Our findings on ROUGE-L also align with the recent study (Janiak et al., 2025), more details in Appendix A.5.

4.6 Training and Inference Details

Details on training and inference stack, hardware, hyperparameters, and throughput comparison with baselines are in Appendix A.2.

5 Results and Discussion

5.1 Task Performance

Performance via AR. We identify a shortcoming of the ROUGE-L score (see Section 4.5), where

	Llama 2 7B				Granite 3.3 2B			
	AR	ALaJ	BWT _r	BWT _{LaJ}	AR	ALaJ	BWT _r	BWT _{LaJ}
Baselines								
Base Model	1.65	23.69	–	–	45.11	63.67	–	–
NoRep	21.57	61.58	-55.78	-20.88	32.33	66.64	-45.59	-15.85
SSR	66.30	77.23	-7.82	-4.04	67.44	77.43	-8.13	-4.41
Rep	69.09	78.40	-4.85	-2.51	69.89	78.29	-5.38	-3.34
DR _{base}	69.87	78.44	-4.23	-2.44	70.95	78.70	-4.55	-3.11
DR _{task}	70.04	78.69	-3.70	-2.13	71.01	79.08	-4.32	-2.73
Ablations								
DR _{task} + MoT (2 T + 1 A)	66.45	81.27	-7.72	0.13	67.66	80.94	-7.93	-0.70
TF-SSD (1 T - MixIns)	67.54	81.06	-6.58	-0.11	67.74	81.59	-7.77	0.04
TF-SSD (1 T - Wiz)	67.16	81.46	-6.90	0.36	67.33	81.58	-8.23	0.00
TF-SSD (2 T + 1 A)	66.97 (-3.07)	81.51 (+2.82)	-7.16 (-3.46)	0.47 (+2.6)	66.65 (-4.36)	81.74 (+2.66)	-8.86 (-4.54)	0.29 (+3.02)
MTL	73.65	81.16	–	–	74.45	81.20	–	–

Table 4: Evaluation scores of various baselines, ablations and MTL are shown. **First-**, **second-**, and **third-** best performing methods are highlighted over all baselines and ablations. MTL is omitted from ranking since it is considered as an ideal target performance reference for a CFT model to achieve. Performance difference of TF-SSD (2T + 1A) with top method (DR_{task}) in the baselines is shown in parentheses.

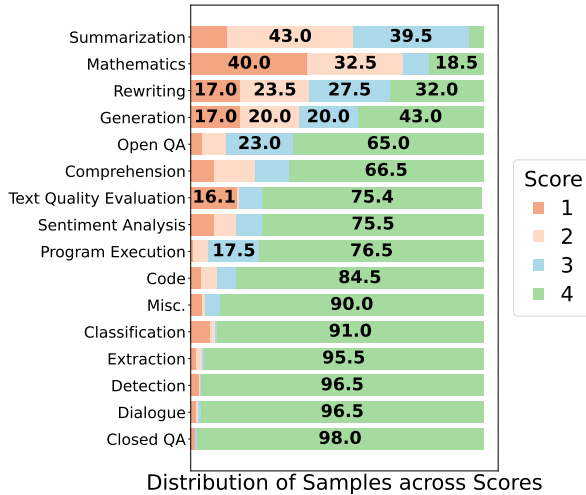


Figure 3: Distribution of the samples across Scorer’s scores for 16 tasks for Llama 2 7B. Percentage of samples across scores are annotated and insignificant percentages (below 15%) are not shown.

correct answers are penalized (see Table 2), primarily due to low overlap with the ground truth though semantically correct. Since TF-SSD involves generating a portion of the tokens from a set of open models, it is more susceptible to such penalties. This explains the observed decline of ~ 3.1 and ~ 4.4 points in AR (Table 4) across models compared to the top-performing baseline (DR_{task}). To address this, we adopt ALaJ, which evaluates responses more holistically by considering both the ground truth and the nature of the task. The higher MTL scores validate the fairness of LaJ.

Performance via ALaJ. TF-SSD consistently outperforms all baselines (see Table 4), showing an improvement of ~ 2.7 points in ALaJ over the best-

performing baseline (DR_{task}) across both model families. Moreover, TF-SSD meets the ideal target performance of MTL and slightly exceeds it by ~ 0.45 points on average across models. Results are consistent with both the judges (Tables 10, 4).

Performance over MTL. The reason behind the higher performance of TF-SSD over MTL (Table 4) is attributed to the refinement submethod which involves teacher forced self-distillation using tokens coming from an ensemble of open models outside the distribution of the replay dataset. As shown in the Table 3 where the response of the TF-SSD is well-formed and accurate while MTL provides an incorrect response not covered in the context.

5.2 Catastrophic Forgetting.

BWT_r and BWT_{LaJ} are not applicable to the MTL and Base Model baselines since these models do not undergo CFT.

CF via BWT_r. BWT_r is a derived metric of ROUGE-L, the same limitations apply attributing to an average drop of ~ 4 points (see Table 4) of TF-SSD compared to DR_{task} across models.

CF via BWT_{LaJ}. BWT_{LaJ} derived from the ALaJ score, addresses this gap. TF-SSD consistently outperforms all baselines (Table 4), with an average improvement of ~ 2.81 points over DR_{task}.

Positive BWT_{LaJ} score for TF-SSD. The last-stage model outperforms the task-stage model, suggesting it has not only retained prior knowledge but also gained enhanced capabilities. We attribute this behavior to the refinement submethod, which provides an improved replay set that was not available in training dataset of the task-stage model.

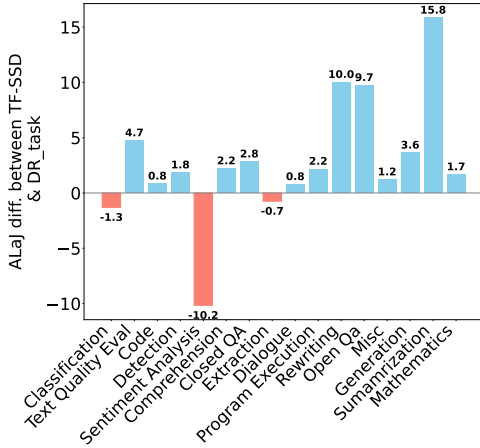


Figure 4: Task-wise difference in ALaJ performance between TF-SSD and DR_{task} using Llama2 7B

5.3 Ablation Study

Effect of self-distillation. Baselines DR_{base} and DR_{task} perform self-distillation by omitting ground-truth labels, resulting in average ALaJ improvements of ~ 0.23 points for DR_{base} and ~ 0.54 points for DR_{task} over the Rep baseline, which uses ground-truth labels. Distilling from the respective task-stage model leads to a 2X improvement compared to distillation from the base model. Our method, TF-SSD, which uses teacher-forced selective self-distillation, achieves an improvement of ~ 3.28 points over the Rep baseline that is roughly a 6X gain over the improvement from DR_{task} . A similar trend is observed for BWT_{LaJ} as well.

Effect of the Scorer. Overall, $\text{DR}_{\text{task}} + \text{MoT}$ (2 T + 1 A) and TF-SSD demonstrate the effectiveness of the Scorer. For $\text{DR}_{\text{task}} + \text{MoT}$ (2 T + 1 A), the Scorer’s classification of samples into accept (use as-is) or reject (replace labels with MoT response) improves ALaJ by ~ 2.5 points and ~ 2.2 points in BWT_{LaJ} over DR_{task} . Similarly, improvements from TF-SSD over DR_{task} are detailed in Sections 5.1 and 5.2.

We further conduct task-level analysis to understand the Scorer’s impact. While overall performance improves, DR_{task} performs better than TF-SSD on 3 tasks—Classification, Sentiment Analysis, and Extraction (see Figure 4). The drop in performance is proportional to the number of samples with scores below 4 (see Figure 3) specifically for these 3 tasks. This suggests that while the Scorer benefits most tasks, its score distribution may negatively impact a few. Our study treats tasks as black boxes and we plan to study task-specific optimizations in the future.

Effect of teacher forcing of the labels. The performance improvement of TF-SSD over $\text{DR}_{\text{task}} + \text{MoT}$ (2 T + 1 A) highlights the positive impact of selective teacher-forcing of labels. TF-SSD outperforms $\text{DR}_{\text{task}} + \text{MoT}$ (2 T + 1 A) by ~ 0.52 points in ALaJ and ~ 0.67 points in BWT_{LaJ} .

Effect of teachers in MoT. Ablations TF-SSD (1T - MixIns), TF-SSD (1T - Wiz) and TF-SSD (2T + 1A) demonstrate the effectiveness of the MoT paradigm. Individual teacher ablations fail to outperform 2T + 1A configuration (see Table 4) which is consistent across the model families. 2T + 1A configuration shows an average improvement of ~ 0.25 points in ALaJ and ~ 0.31 points in BWT_{LaJ} over single teacher configurations. Further, we investigate (see Table 9) diversity of models in MoT and findings are consistent.

LaJ We use two judges: Prometheus 8x7B and DeepSeek-V3.2 (685B) (Table 10). TF-SSD consistently outperforms the top baseline DR_{task} , indicating that our gains are agnostic of judge model and free from same scorer bias.

5.4 Task Order, Diversity, and Replay Data Size Robustness

Effect of task order and diversity. TF-SSD is robust to task order and diversity in CFT consistently outperforming (Tables 8 and 12) DR_{task} for 2 distinct random task orders and a different set of 10 tasks (Appendices A.8 and A.7).

Effect of replay data size. TF-SSD is stable (Table 11) to different replay data sizes and consistently outperforms DR_{task} (Appendix A.9).

6 Conclusion

We present TF-SSD, a novel method designed to mitigate catastrophic forgetting in continual learning, improve overall task performance, and address a critical real-world scenario of task ownership and noisy replay. TF-SSD combines (1) Label Distillation: where labels are distilled from the respective task-stage model, and (2) Refinement: applies teacher-forced selective self-distillation to improve label quality, to curate noisy replay in to a more effective replay. TF-SSD outperforms 7 challenging baselines and multi-task learning in task performance and mitigating catastrophic forgetting. We hope our work would motivate further research catering to various real-world constraints in continual learning.

569 Limitations

570 Though TF-SSD outperforms other baselines, we
571 identify the following limitations that we plan to
572 study in the future. First, we restrict the MoT com-
573 ponent in our method to open-source instruction-
574 tuned models. We acknowledge that leveraging
575 stronger proprietary models could potentially yield
576 even greater performance improvements than those
577 reported. Due to restrictive synthetic data genera-
578 tion usage terms of some model families both in
579 open and closed source, we confine to a specific set.
580 However, we perform extensive ablations showcas-
581 ing the contributions of the number and diversity of
582 the teachers in MoT. Given the observations from
583 the ablations, choosing better teachers would fur-
584 ther improve the performance of the TF-SSD in
585 proprietary or licensed settings. Second, scorer
586 plays a central role in determining the extent of
587 teacher forcing from the MoT but from our results,
588 some tasks are negatively affected. Thus, using
589 a stronger proprietary or licensed models would
590 potentially overcome the negative effect and may
591 benefit all tasks. Third, our approach is intention-
592 ally designed to be task properties agnostic. How-
593 ever, we plan to study task-specific optimizations
594 in future for additional improvements. Tokens from
595 MoT are chosen from left to right in our approach,
596 which may not be ideal in some cases where im-
597 portant tokens are scattered at different positions,
598 indentifying such tokens and incorporating them
599 in our method comes with the challenge of identi-
600 fying a right way to infix them in teacher forced
601 self-distillation submethod of TF-SSD where mod-
602 els used are auto-regressive generating tokens from
603 left to right. Finally, amount of tokens is inversely
604 proportional to the Scorer’s score which is discrete
605 in nature keeping the proportions rigid, while these
606 proportions may be made more continuous and flexi-
607 ble which we plan to study in future. Irrespectively,
608 across diverse set of tasks in a continual learning
609 setting, we see TF-SSD providing improvements
610 over the baselines and ablations validating value
611 addition of our method.

612 Ethics Statement

613 We have taken deliberate steps to filter unsafe or
614 harmful data while selecting tasks for training data
615 from SuperNI dataset. Due to scale and diversity
616 of the resources, the data may still contain sensitive
617 informations.

References

- 618
619 Rahaf Aljundi, Francesca Babiloni, Mohamed Elho-
620 seiny, Marcus Rohrbach, and Tinne Tuytelaars. 2018.
621 *Memory aware synapses: Learning what (not) to*
622 *forget*. In *Computer Vision – ECCV 2018: 15th Eu-*
623 *ropean Conference, Munich, Germany, September*
624 *8–14, 2018, Proceedings, Part III*, page 144–161,
625 Berlin, Heidelberg. Springer-Verlag.
- 626
627 Kushagra Bhushan, Yatin Nandwani, Dinesh Khandel-
628 wal, Sonam Gupta, Gaurav Pandey, Dinesh Raghu,
629 and Sachindra Joshi. 2025. *Systematic knowledge*
630 *injection into large language models via diverse aug-*
631 *mentation for domain-specific rag*.
- 632
633 Magdalena Biesialska, Katarzyna Biesialska, and
634 Marta R. Costa-jussà. 2020. *Continual lifelong learn-*
635 *ing in natural language processing: A survey*. In
636 *Proceedings of the 28th International Conference*
637 *on Computational Linguistics*, pages 6523–6541,
638 Barcelona, Spain (Online). International Committee
639 on Computational Linguistics.
- 640
641 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
642 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
643 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
644 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
645 Gretchen Krueger, Tom Henighan, Rewon Child,
646 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens
647 Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-
648 teusz Litwin, Scott Gray, Benjamin Chess, Jack
649 Clark, Christopher Berner, Sam McCandlish, Alec
650 Radford, Ilya Sutskever, and Dario Amodei. 2020.
651 *Language models are few-shot learners*. In *Ad-*
652 *vances in Neural Information Processing Systems*,
653 volume 33, pages 1877–1901. Curran Associates,
654 Inc.
- 655
656 Pietro Buzzega, Matteo Boschini, Angelo Porrello, Da-
657 vide Abati, and Simone Calderara. 2020. Dark expe-
658 rience for general continual learning: a strong, simple
659 baseline. In *Proceedings of the 34th International*
660 *Conference on Neural Information Processing Sys-*
661 *tems, NIPS ’20, Red Hook, NY, USA. Curran Asso-*
662 *ciates Inc*.
- 663
664 Arslan Chaudhry, Marcus Rohrbach, Mohamed El-
665 hoseiny, Thalaiyasingam Ajanthan, Puneet Kumar
666 Dokania, Philip H. S. Torr, and Marc Aurelio Ran-
667 zato. 2019. *Continual learning with tiny episodic*
668 *memories*. *CoRR*, abs/1902.10486.
- 669
670 Xi Chen and Min Zeng. 2025. *Prototype conditioned*
671 *generative replay for continual learning in NLP*. In
672 *Proceedings of the 2025 Conference of the Nations*
673 *of the Americas Chapter of the Association for Com-*
674 *putational Linguistics: Human Language Technolo-*
675 *gies (Volume 1: Long Papers)*, pages 12754–12770,
676 Albuquerque, New Mexico. Association for Compu-
677 tational Linguistics.
- 678
679 Zhiyuan Chen and Bing Liu. 2018. *Continual Learning*
680 *and Catastrophic Forgetting*, pages 55–75. Springer
681 International Publishing, Cham.

675	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	Casas, Emma Bou Hanna, Florian Bressand, Gi-	734
676	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	anna Lengyel, Guillaume Bour, Guillaume Lam-	735
677	Akhil Mathur, Alan Schelten, Amy Yang, Angela	ple, L�lio Renard Lavaud, Lucile Saulnier, Marie-	736
678	Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,	Anne Lachaux, Pierre Stock, Sandeep Subramanian,	737
679	Archi Mitra, Archie Sravankumar, Artem Korenev,	Sophia Yang, Szymon Antoniak, Teven Le Scao,	738
680	Arthur Hinsvark, Arun Rao, Aston Zhang, Aur�lien	Th�ophile Gervet, Thibaut Lavril, Thomas Wang,	739
681	Rodr�guez, Austen Gregerson, Ava Spataru, Bap-	Timothe� Lacroix, and William El Sayed. 2024. <i>Mix-</i>	740
682	tiste Rozi�re, Bethany Biron, Binh Tang, Bobbie	<i>tral of experts</i> .	741
683	Chern, Charlotte Caucheteux, Chaya Nayak, Chloe		
684	Bi, Chris Marra, Chris McConnell, Christian Keller,	Zixuan Ke, Hu Xu, and Bing Liu. 2021. <i>Adapting BERT</i>	742
685	Christophe Touret, Chunyang Wu, Corinne Wong,	<i>for continual learning of a sequence of aspect senti-</i>	743
686	Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-	<i>ment classification tasks</i> . In <i>Proceedings of the 2021</i>	744
687	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	<i>Conference of the North American Chapter of the</i>	745
688	David Esiobu, Dhruv Choudhary, Dhruv Mahajan,	<i>Association for Computational Linguistics: Human</i>	746
689	Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,	<i>Language Technologies</i> , pages 4746–4755, Online.	747
690	Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,	Association for Computational Linguistics.	748
691	Emily Dinan, Eric Michael Smith, Filip Radenovic,		
692	Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Geor-	Seungone Kim, Juyoung Suk, Shayne Longpre,	749
693	gia Lewis Anderson, Graeme Nail, Gr�goire Mialon,	Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham	750
694	Guan Pang, Guillem Cucurell, Hailey Nguyen, Han-	Neubig, Moontae Lee, Kyungjae Lee, and Minjoon	751
695	nah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov,	Seo. 2024. <i>Prometheus 2: An open source language</i>	752
696	Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan	<i>model specialized in evaluating other language mod-</i>	753
697	Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan	<i>els</i> . In <i>Proceedings of the 2024 Conference on Empir-</i>	754
698	Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,	<i>ical Methods in Natural Language Processing</i> , pages	755
699	Jeet Shah, Jelmer van der Linde, Jennifer Billock,	4334–4353, Miami, Florida, USA. Association for	756
700	Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,	Computational Linguistics.	757
701	Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu,		
702	Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz,	758
703	Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia,	Joel Veness, Guillaume Desjardins, Andrei A. Rusu,	759
704	Kalyan Vasuden Alwala, Kartikeya Upasani, Kate	Kieran Milan, John Quan, Tiago Ramalho, Ag-	760
705	Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and	gnieszka Grabska-Barwinska, Demis Hassabis, Clau-	761
706	et al. 2024. <i>The llama 3 herd of models</i> . <i>CoRR</i> ,	dia Clopath, Dharshan Kumaran, and Raia Hadsell.	762
707	abs/2407.21783.	2017. <i>Overcoming catastrophic forgetting in neural</i>	763
		<i>networks</i> . <i>Proceedings of the National Academy of</i>	764
		<i>Sciences</i> , 114(13):3521–3526.	765
708	Google. 2022. <i>Learning to prompt for continual learn-</i>		
709	<i>ing</i> .	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	766
		Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gon-	767
710	Jinghan He, Haiyun Guo, Kuan Zhu, Zihan Zhao, Ming	zalez, Hao Zhang, and Ion Stoica. 2023. <i>Efficient</i>	768
711	Tang, and Jinqiao Wang. 2024. <i>SEEKR: Selective</i>	<i>memory management for large language model serv-</i>	769
712	<i>attention-guided knowledge retention for continual</i>	<i>ing with pagedattention</i> . In <i>Proceedings of the 29th</i>	770
713	<i>learning of large language models</i> . In <i>Proceedings</i>	<i>Symposium on Operating Systems Principles, SOSP</i>	771
714	<i>of the 2024 Conference on Empirical Methods in</i>	'23, page 611–626, New York, NY, USA. Association	772
715	<i>Natural Language Processing</i> , pages 3254–3266, Mi-	for Computing Machinery.	773
716	ami, Florida, USA. Association for Computational		
717	Linguistics.	Zhizhong Li and Derek Hoiem. 2018. <i>Learning without</i>	774
		<i>forgetting</i> . <i>IEEE Trans. Pattern Anal. Mach. Intell.</i> ,	775
718	Jianheng Huang, Leyang Cui, Ante Wang, Chengyi	40(12):2935–2947.	776
719	Yang, Xinting Liao, Linfeng Song, Junfeng Yao, and		
720	Jinsong Su. 2024. <i>Mitigating catastrophic forgetting</i>	Chin-Yew Lin. 2004. <i>ROUGE: A package for auto-</i>	777
721	<i>in large language models with self-synthesized re-</i>	<i>matic evaluation of summaries</i> . In <i>Text Summariza-</i>	778
722	<i>hearsal</i> . In <i>Proceedings of the 62nd Annual Meeting</i>	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.	779
723	<i>of the Association for Computational Linguistics (Vol-</i>	Association for Computational Linguistics.	780
724	<i>ume 1: Long Papers)</i> , pages 1416–1428, Bangkok,		
725	Thailand. Association for Computational Linguistics.	David Lopez-Paz and Marc’Aurelio Ranzato. 2017.	781
		Gradient episodic memory for continual learning. In	782
726	Denis Janiak, Jakub Binkowski, Albert Sawczyn, Bog-	<i>Proceedings of the 31st International Conference on</i>	783
727	dan Gabrys, Ravid Schwartz-Ziv, and Tomasz Ka-	<i>Neural Information Processing Systems, NIPS’17</i> ,	784
728	жданович. 2025. The illusion of progress: Re-	page 6470–6479, Red Hook, NY, USA. Curran Asso-	785
729	evaluating hallucination detection in llms. <i>arXiv</i>	ciates Inc.	786
730	<i>preprint arXiv:2508.08285</i> .	Michael McCloskey and Neal J. Cohen. 1989. <i>Catas-</i>	787
		<i>trophic interference in connectionist networks: The</i>	788
731	Albert Q. Jiang, Alexandre Sablayrolles, Antoine	<i>sequential learning problem</i> . volume 24 of <i>Psychol-</i>	789
732	Roux, Arthur Mensch, Blanche Savary, Chris	<i>ogy of Learning and Motivation</i> , pages 109–165. Aca-	790
733	Bamford, Devendra Singh Chaplot, Diego de las	ademic Press.	791

792	Meryem M'hamdi and Jonathan May. 2024. Leitner-guided memory replay for cross-lingual continual learning . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 7808–7821, Mexico City, Mexico. Association for Computational Linguistics.	850	Thomas Scialom, Tuhin Chakrabarty, and Smaranda Muresan. 2022. Fine-tuned language models are continual learners . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 6107–6122, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	851
793		852		853
794		854		855
795		856		
796				
797				
798				
799				
800	Mayank Mishra, Matt Stallone, Gaoyuan Zhang, Yikang Shen, Aditya Prasad, Adriana Meza Soria, Michele Merler, Parameswaran Selvam, Saptha Surendran, Shivdeep Singh, Manish Sethi, Xuan-Hong Dang, Pengyuan Li, Kun-Lung Wu, Syed Zawad, Andrew Coleman, Matthew White, Mark Lewis, Raju Pavuluri, Yan Koyfman, Boris Lublinsky, Maximilien de Bayser, Ibrahim Abdelaziz, Kinjal Basu, Mayank Agarwal, Yi Zhou, Chris Johnson, Aanchal Goyal, Hima Patel, Yousaf Shah, Petros Zerfos, Heiko Ludwig, Asim Munawar, Maxwell Crouse, Pavan Kapanipathi, Shweta Salaria, Bob Calio, Sophia Wen, Seetharami Seelam, Brian Belgodere, Carlos Fonseca, Amith Singhee, Nirmitt Desai, David D. Cox, Ruchir Puri, and Rameswar Panda. 2024. Granite code models: A family of open foundation models for code intelligence .		Yijun Tian, Yikun Han, Xiusi Chen, Wei Wang, and Nitesh V. Chawla. 2025. Beyond answers: Transferring reasoning capabilities to smaller llms using multi-teacher knowledge distillation . In <i>Proceedings of the Eighteenth ACM International Conference on Web Search and Data Mining, WSDM '25</i> , page 251–260, New York, NY, USA. Association for Computing Machinery.	857
801		858		859
802		860		861
803		862		863
804		864		
805				
806				
807				
808				
809				
810				
811				
812				
813				
814				
815				
816				
817	Jisoo Mok, Jaeyoung Do, Sungjin Lee, Tara Taghavi, Seunghak Yu, and Sungroh Yoon. 2023. Large-scale lifelong learning of in-context instructions and how to tackle it . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12573–12589, Toronto, Canada. Association for Computational Linguistics.		Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models .	865
818		866		867
819		868		869
820		870		871
821		872		873
822		874		875
823		876		877
824	NVIDIA. 2020. Nemotron-pretraining-dataset-sample .	878		879
825		880		881
826	German I. Parisi, Ronald Kemker, Jose L. Part, Christopher Kanan, and Stefan Wermter. 2019. Continual lifelong learning with neural networks: A review . <i>Neural Networks</i> , 113:54–71.	882		883
827		884		885
828		886		887
829	Anastasia Razdaibiedina, Yuning Mao, Rui Hou, Madian Khabsa, Mike Lewis, and Amjad Almahairi. 2023. Progressive prompts: Continual learning for language models . In <i>The Eleventh International Conference on Learning Representations</i> .	888		889
830		890		891
831		892		
832				
833				
834	Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. 2017. icarl: Incremental classifier and representation learning . In <i>2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 5533–5542.		Junlin Wang, Jue WANG, Ben Athiwaratkun, Ce Zhang, and James Zou. 2025. Mixture-of-agents enhances large language model capabilities . In <i>The Thirteenth International Conference on Learning Representations</i> .	893
835		894		895
836		896		897
837		898		899
838		900		901
839	David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. 2019. Experience replay for continual learning . In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.	902		903
840		904		905
841		906		907
842		908		
843				
844	Shrey Satapara and P. K. Srijith. 2024. TL-CL: Task and language incremental continual learning . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 12123–12142, Miami, Florida, USA. Association for Computational Linguistics.		Xiao Wang, Tianze Chen, Qiming Ge, Han Xia, Rong Bao, Rui Zheng, Qi Zhang, Tao Gui, and Xuanjing Huang. 2023a. Orthogonal subspace learning for language model continual learning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 10658–10671, Singapore. Association for Computational Linguistics.	902
845		903		904
846		905		906
847		907		908
848				
849				

909 Yifan Wang, Yafei Liu, Chufan Shi, Haoling Li, Chen
 910 Chen, Haonan Lu, and Yujiu Yang. 2024b. *InsCL: A data-efficient continual learning paradigm for fine-tuning large language models with instructions*. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 663–677, Mexico City, Mexico. Association for Computational Linguistics.

919 Yizhong Wang, Swaroop Mishra, Pegah Alipoor-
 920 molabashi, Yeganeh Kordi, Amirreza Mirzaei,
 921 Anjana Arunkumar, Arjun Ashok, Arut Selvan
 922 Dhanasekaran, Atharva Naik, David Stap, et al. 2022. Super-naturalinstructions:generalization via declarative instructions on 1600+ tasks. In *EMNLP*.

925 Zhicheng Wang, Yufang Liu, Tao Ji, Xiaoling Wang,
 926 Yuanbin Wu, Congcong Jiang, Ye Chao, Zhencong
 927 Han, Ling Wang, Xu Shao, and Wenqiu Zeng. 2023b. *Rehearsal-free continual language learning via efficient parameter isolation*. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10933–10946, Toronto, Canada. Association for Computational Linguistics.

934 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,
 935 Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei
 936 Lin, and Daxin Jiang. 2024. *WizardLM: Empowering large pre-trained language models to follow complex instructions*. In *The Twelfth International Conference on Learning Representations*.

940 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
 941 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
 942 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in neural information processing systems*, 36:46595–46623.

946 A Appendix

947 A.1 Prompts

948 **LLM-as-a-Judge.** We use the following prompt
 949 template (Listing A.1) for LLM-as-a-Judge.

LLM-as-a-Judge prompt.

You are a helpful and fair evaluator trained to judge the quality of AI-generated answers. Your job is to assess whether the model’s response is correct and appropriate for the given instruction and input.

Sometimes the provided reference answer (ground truth) might be incomplete, imprecise, or incorrect. So do **not assume** it is always correct. Instead, evaluate the model’s answer based on reasoning and factual correctness.

Please follow these steps:

1. Read the instruction and input carefully.
2. Analyze the model’s response: Is it logically valid, factually correct, and aligned with the instruction?
3. Compare it with the reference answer (if useful), but do not blindly trust it.
4. Rate the model’s response on a scale from 0 to 2, where:
 - 2 = **Correct**
 - 1 = **Acceptable but with minor issues**
 - 0 = **Incorrect and imperfect**
5. Give score within `<score>` and `</score>` tags only

You must explain your judgment briefly.

Now evaluate the following:

```

Instruction:{instruction}
Input:{input}
Model’s Response:{predicted response}
Reference Answer:{golden response}
Your Evaluation:
  
```

Scorer prompt for TF-SSD. We use the following prompt template (Listing A.1) for Scorer model when used in TF-SSD.

950

951

952

953

Scorer prompt for TF-SSD.

You are given a question, the corresponding ground-truth answer, and a prediction from a model. Analyse and categorise the given prediction into one of the four labels.

There should be no contradicting statements in a good prediction. The prediction may contain extra information. If the prediction states something as a possibility, treat it as a definitive answer.

A good prediction must contain all the important information presented in the ground truth, but doesn't have to fully match it word by word.

To make your decision, first read the question and ground-truth answer carefully. Then compare the given prediction with the ground-truth answer in the light of the question.

Start with an explanation and reasoning for your evaluation within `<explanation>` and `</explanation>` tags.

Then, within `<category>` and `</category>` tags, output one of the following four labels number only:

- 4: Fully correct, complete, and fluent.
- 3: Understandable, but lacks minor details.
- 2: Understandable, but lacks major details.
- 1: Contains factual errors, hallucinations, or misleading info.

The output format should as follow: `<explanation>` (the explanation) `</explanation>` `<category>` (category the prediction belongs to) `</category>`

Question: {full prompt} Ground-truth answer: {golden response} Prediction: {predicted response}

Output:

Scorer prompt for $DR_{task} + MoT$ (2 T + 1 A) baseline. We use the following prompt template (Listing A.1) for Scorer model when used in $DR_{task} + MoT$ (2 T + 1 A) baseline.

Scorer prompt when used in $DR_{task} + MoT$ (2 T + 1 A) baseline.

You are given a question, the corresponding ground-truth answer and a prediction from a model. Compare the "Ground-truth answer" and the "Prediction" to determine whether the prediction correctly answers the question.

There should be no contradicting statements in a good prediction. The prediction may contain extra information. If the prediction states something as a possibility, treat it as a definitive answer.

A good prediction must contain all the important information presented in the ground truths, but doesn't have to fully match it word by word.

To make your decision, first read the question and Ground-truth answer carefully. Then compare the given Prediction with the Ground-truth answer in the light of the question. Start with an explanation and reasoning for your evaluation within `<explanation>` and `</explanation>` tags.

Then, within `<score>` and `</score>` tokens, generate "1" if the Prediction is correct in the light of Ground-truth and question. Otherwise generate "0" if it is incomplete or incorrect.

Question: {full prompt}

Ground-truth answer: {golden response}

Prediction: {predicted response}

Answer:

Teacher prompt. We use the following prompt template (Listing A.1) for the teacher.

Teacher prompt

Instruction:
{instruction}. Give only final answer within `<response>` and `</response>` tags without any explanation.

Input:

{input}

Response:

Aggregator prompt. We use the following prompt template (Listing A.1) for the aggregator.

954

955

956

957

958

959

960

961

962

963

964

Aggregator prompt

You have been provided with a set of responses from various open-source models along with user task instruction and input for which the response is generated. Your task is to synthesize these responses into a single, high-quality response. It is crucial to critically evaluate the information provided in these responses, recognizing that some of it may be biased or incorrect. Your response should not simply replicate the given answers but should offer a refined, accurate, and comprehensive reply to the instruction. Ensure your response is well-structured, coherent, and adheres to the highest standards of accuracy and reliability.

Response: {teacher1}

Response: {teacher2}

Task Instruction:

{instruction}. Give only final answer within <response> and </response> tags without any explanation.

Input:

{input}

Response:

A.2 Inference and Training Setup Details

Hardware. We use $8 \times A100$ 80GB GPUs in a single node with 32 cpus and 400Gi memory.

Inference software. We use vLLM³ (Kwon et al., 2023) with tensor parallelism and greedy decoding with temperature set to 0.

Training software. We built our codebase on top of the SSR framework (Huang et al., 2024), which in turn extends the LLaMA Factory—a lightweight and efficient framework for training and fine-tuning large language models. For each task in our continual learning setup, we employed parameter-efficient fine-tuning using LoRA, implemented via the peft library. The training flow is carried using bash scripts, enabling sequential learning of tasks in a continual fashion.

³github.com/vllm-project/vllm

Training hyperparameters. We use a learning rate of $2e-4$ and a batch size of 32 with gradient accumulation set to 1 for all experiments. LoRA is applied with a rank of 8 and a dropout rate of 0.1 to the Q and V projection matrices. The model is trained for 3 epochs in each continual learning task.

Comparison of training time overhead with baselines The overhead added at each stage of the continual learning pipeline is ~ 10 minutes which translates to ~ 2 hours for the entire pipeline consisting 16 stages compared to baseline. We use the state-of-the-art open source setup (vLLM) which uses accelerations such as tensor parallelism to keep this overhead least possible.

A.3 Model Details

We use several language models within the paper. All the models are used from their respective HuggingFace source: Llama 2 7B⁴ (Touvron et al., 2023), Granite 3.3 2B⁵ (Mishra et al., 2024), Mixtral 8x22B Instruct (MixIns)⁶ (Jiang et al., 2024), WizardLM 2 8x22B(Wiz)⁷ (Xu et al., 2024), Prometheus 8X7B⁸ (Kim et al., 2024).

Task Name	Size	Task ID
Classification	18k	50, 1712, 65
Program Execution	17k	63, 93, 370
Mathematics	17k	85, 87, 90
Generation	16k	1, 67, 1730
Summarization	16k	589, 668, 1290
Open QA	15k	2, 24, 1731
Sentiment Analysis	15k	195, 293, 843
Rewriting	14k	402, 413, 1340
Text Quality Evaluation	12k	616, 675, 1283
Code	10k	77, 211, 869
Detection	9k	88, 209, 318
Miscellaneous	9k	305, 383, 700
Comprehension	9k	27, 1664, 223
Dialogue	7.5k	362, 766, 1500
Extraction	6.5k	39, 180, 1568
Closed QA	3k	73, 296, 667

Table 5: The subtask IDs within each task taken from the SuperNI dataset.

⁴hf.co/meta-llama/Llama-2-7b

⁵hf.co/ibm-granite/granite-3.3-2b-base

⁶hf.co/mistralai/Mixtral-8x22B-Instruct-v0.1

⁷huggingface.co/alpindale/WizardLM-2-8x22B

⁸hf.co/prometheus-eval/prometheus-8x7b-v2.0

Task	Max Length	Avg Length
Classification	710	47.17
Text Quality Evaluation	60	17.55
Code	216	41.97
Detection	44	11.75
Sentiment Analysis	71	22.10
Comprehension	1809	117.13
Closed QA	88	42.35
Extraction	152	21.22
Dialogue	653	72.11
Program Execution	31	10.18
Rewriting	414	42.50
Open QA	1160	188.84
Misc.	149	26.46
Generation	1143	141.71
Summarization	7329	212.51
Mathematics	17	10.61

Table 6: Maximum and average length (words) per task

A.4 Dataset Details

SuperNI (Wang et al., 2022) is a collection of diverse NLP tasks with NL instructions, released under the Apache-2.0 license. Each task is stored in a separate file, identified by a unique task ID. Table 5 summarizes the data size and subtask IDs under each task category. Table 6 reports the average and maximum input length in words. We randomly hold out 20% of the dataset for evaluation and 200 samples as buffer replay data.

Task Setup. In our experiments, we consider a 16-task setup grouped under broader task categories, as shown in Table 5. Each category includes 3 subtasks, and their corresponding task IDs are listed in the same table. To illustrate the diversity of subtasks, consider the Mathematics category. Below are the NL instructions used for its three subtasks:

Task ordering We adopt the following continual learning orders:

- **Order 1:** Classification → Text Quality Evaluation → Code → Detection → Sentiment Analysis → Comprehension → Closed QA → Extraction → Dialogue → Program Execution → Rewriting → Open QA → Misc. → Generation → Summarization → Mathematics

- **Order 2:** Generation → Mathematics → Extraction → Comprehension → Text Quality Evaluation → Dialogue → Classification → Code → Misc. → Summarization → Program Execution → Rewriting → Closed QA → Detection → Sentiment Analysis → Open QA

Subtask Instructions in Mathematics Task

1. In this task, you will be given an arithmetic operation and you have to find its answer. The symbols of the operators '+' and '-' have been swapped, i.e., you need to perform subtraction when you see a '+' symbol and addition when you see a '-' symbol.
2. In this task, you will be given an arithmetic operation and you have to find its answer. The operators '+' and '-' have been replaced with new symbols. Specifically, '+' has been replaced with the symbol '@' and '-' with the symbol '#'. You need to perform the operations in the given equation and return the answer.
3. A polynomial equation is a sum of terms. Each term is either a constant number or consists of the variable x raised to a power and multiplied by a coefficient (called the weight). For example, in the polynomial $2x^2 + 3x + 4$, the weights are [2, 3, 4]. A polynomial with weights [6, 4] represents the equation $6x + 4$, while [1, 3, 4] represents $1x^2 + 3x + 4$. In this task, you are given the list of weights and a value for x , and your goal is to compute the result of the polynomial expression.

A.5 Metric Details

Adding to the discussion in Section 4.5, recent work by (Janiak et al., 2025) extensively compared ROUGE-L and LaJ and arrived at the same conclusion that ROUGE-L fails to capture semantic understanding which is crucial in comparing responses with paraphrase variations, while LaJ could accurately capture this difference.

A.6 Additional Experimentation

Figures 5 and 6 show the individual LaJ scores for all 16 tasks across the baseline, ablations and MTL for Llama2 7b and Granite3.3 2B.

A.7 10 Tasks Setting

We also evaluate our method, TF-SSD, in a 10-task continual learning setting using both LLaMA2-7B and Granite 3.3B models. Unlike the 16-task setup, which includes a mix of classification and generation tasks, the 10-task configuration includes majorly generation tasks. The continual learning order adopted: QA → QG → SA → Sum. → Trans.

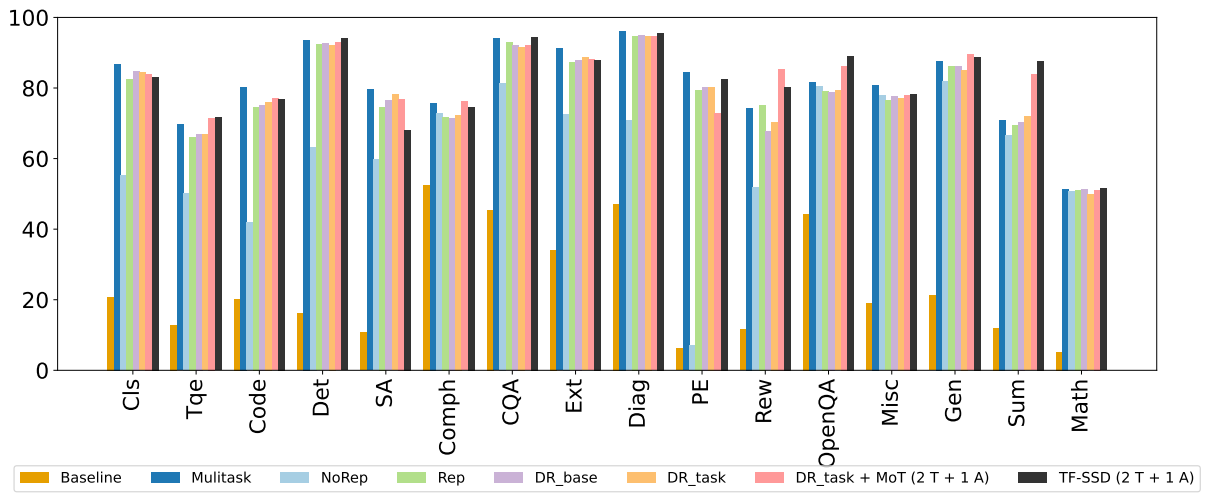


Figure 5: Individual LaJ scores for all 16 tasks for each experiment conducted with LLaMA-2 7B.

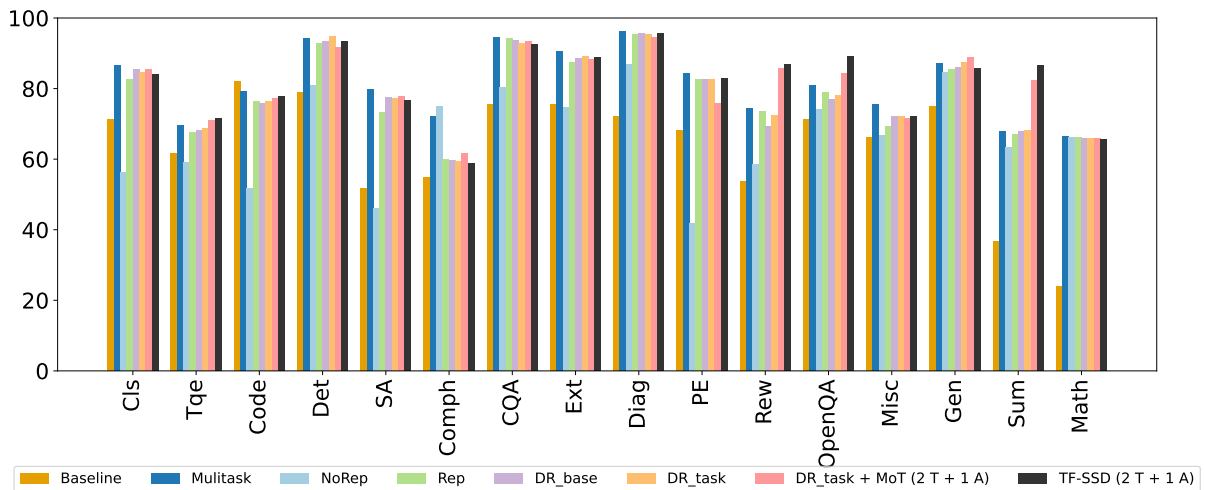


Figure 6: Individual LaJ scores for all 16 tasks for each experiment conducted with Granite 3.3 2B.

→ DSG → Expl. → Para. → PE → POS. Table 12 shows AR, ALaJ, BWT and BWT_{LaJ} for both the models on baselines, ablations and multitask settings.

Task performance As shown in Table 12, our method TF-SSD consistently outperforms all baselines with an improvement of ~ 8 points in ALaJ for Llama2 7B and by ~ 6 points in case of Granite 3.3 2B over best-performing baseline (DR_{task}). Moreover, TF-SSD performs better than MTL by ~ 0.62 points of ALaJ for Llama2 7B and ~ 5.8 points for ALaJ for Granite 3.3 2B. Tables 12 and 4 confirms that our proposed method is robust to diversity in tasks in CFT.

Catastrophic forgetting Similar to 16 tasks results reported in Section 5.2, for 10 tasks our method - TF-SSD outperforms all baselines, with an improvement of ~ 9.54 points of BWT_{LaJ} for Llama2 7B and ~ 8.93 points of BWT_{LaJ} for Granite3.3 2B over best-performing baseline DR_{task} .

A.8 16 Task Setting Additional Order

The two random orders experimented are mentioned in A.4 section. Table 8 shows results for order 2 is consistent with order 1 (see complete results in table 4). Thus confirming our model is robust to task ordering.

A.9 16 Task Setting Additional Replay buffer size

We additionally varied the buffer size for each task instead of keeping it constant, by randomly sampling a buffer size between 100 and 1000. The randomly selected buffer sizes for each task are provided in Table 7. The table 11 shows our method is robust to task buffer data size.

A.10 Concrete Examples

Tables 13 and 14 show example refinement of samples belonging to different scorer category.

- Score 4 task instructions: You are given a short paragraph, a question and two choices to answer from. Choose the correct answer based on the paragraph and write the answer(not the key).
- Score 1,2,3 task instructions: Craft one correct answer to the question given in input. To make it more interesting, try to use non-stereotypical language if possible. Make sure

Task Name	Size
Classification	581
Program Execution	609
Mathematics	741
Generation	742
Summarization	634
Open QA	451
Sentiment Analysis	518
Rewriting	753
Text Quality Evaluation	235
Code	922
Detection	610
Miscellaneous	161
Comprehension	431
Dialogue	837
Extraction	771
Closed QA	716

Table 7: Additional buffer replay size 2 experiment buffer size per task

your correct answer is reasonably long, consistent with the context, and requires common sense (instead of explicit extraction from the context.) In your answer, use as few words as possible from the given context. Use a response that is uncommon/non-stereotypical, so that it is less predictable. To be less repetitive, please vary your language for each question.

1106
1107
1108
1109
1110
1111
1112
1113
1114

	Llama 2 7B				Granite 3.3 2B			
	Order 1		Order 2		Order 1		Order 2	
	DR _{task}	TF-SSD(2 T + 1 A)	DR _{task}	TF-SSD(2 T + 1 A)	DR _{task}	TF-SSD(2 T + 1 A)	DR _{task}	TF-SSD(2 T + 1 A)
AR	70.04	66.97	69.59	66.23	71.01	66.65	69.36	65.80
ALaJ	78.68	81.51	78.63	79.18	79.07	81.74	77.79	79.13
BWT _r	-3.70	-7.16	-4.46	-6.89	-4.32	-8.86	-5.86	-9.55
BWT _{ALaJ}	-2.13	0.46	-2.82	-1.86	-2.73	0.29	-3.77	-2.33

Table 8: Comparative evaluation scores of our method TF-SSD on two random continual learning order with our best-performing baseline DR_{task}. **First-**, best performing method.

	Llama 2 7B					
	Set 1 - (MixIns + Wiz)			Set 2 - (Phi 3.5 MoE + Wiz)		
	TF-SSD(2 T + 1 A)	TF-SSD(1 T - Wiz)	TF-SSD(1 T - MixIns)	TF-SSD(2 T + 1 A)	TF-SSD(1 T - Wiz)	TF-SSD(1 T - Phi 3.5 MoE)
AR	57.23	58.54	59.13	57.41	58.54	57.76
ALaJ	86.48	85.29	83.09	86.28	85.29	85.75
BWT _r	-8.64	-6.68	-6.06	-7.84	-6.68	-7.95
BWT _{ALaJ}	7.47	5.91	3.96	7.46	5.91	6.56

Table 9: Ablation study of our method TF-SSD on two set of MoT models in 10 task setting. **Best-**, number.

	Llama 2 7B		Granite 3.3 2B	
	DR _{task}	TF-SSD(2 T + 1 A)	DR _{task}	TF-SSD(2 T + 1 A)
AR	70.04	66.97	71.01	66.65
ALaJ	68.49	71.76	69.48	71.11
BWT _r	-3.70	-7.16	-4.32	-8.86
BWT _{ALaJ}	-3.27	1.03	-3.49	-2.2

Table 10: Evaluation score of our method TF-SSD using DeepSeek V3.2 model as LLM-as-a-Judge with our best-performing baseline DR_{task}. **First-**, best performing method.

	Llama 2 7B				Granite 3.3 2B			
	Replay size 1		Replay size 2		Replay size 1		Replay size 2	
	DR _{task}	TF-SSD(2 T + 1 A)	DR _{task}	TF-SSD(2 T + 1 A)	DR _{task}	TF-SSD(2 T + 1 A)	DR _{task}	TF-SSD(2 T + 1 A)
AR	70.03	66.97	69.26	66.11	71.01	66.64	71.15	67.52
ALaJ	78.68	81.51	80.14	82.30	79.07	81.74	79.80	82.25
BWT _r	-3.69	-7.16	-2.34	-5.55	-4.31	-8.86	-2.94	-6.86
BWT _{ALaJ}	-2.13	0.46	-1.19	1.16	-2.73	0.29	-1.87	0.89

Table 11: Comparative evaluation scores of our method TF-SSD on two buffer replay data size for each task with our best-performing baseline DR_{task}. **First-**, best performing method.

	Llama 2 7B				Granite 3.3 2B			
	AR	ALaJ	BWT _r	BWT _{LaJ}	AR	ALaJ	BWT _r	BWT _{LaJ}
Baselines								
Base Model	5.59	32.39	—	—	44.79	60.56	—	—
NoRep	10.64	46.08	-61.70	-38.68	45.17	68.38	-23.41	-13.62
Rep	61.19	77.64	-3.90	-1.55	61.66	78.30	-3.89	-1.20
DR _{base}	64.74	78.31	-0.96	-2.04	65.03	79.63	-1.29	-0.52
DR _{task}	64.84	78.40	-0.57	-2.07	65.03	80.54	-0.91	-1.23
Ablations								
DR _{task} + MoT (2 T + 1 A)	59.96	82.95	-5.45	3.14	59.93	83.28	-6.33	3.85
TF-SSD (1 T - MixIns)	59.13	83.09	-6.06	3.96	59.49	85.91	-6.79	6.23
TF-SSD (1 T - Wiz)	58.54	85.29	-6.68	5.91	57.79	86.38	-8.31	6.90
TF-SSD (2 T + 1 A)	57.23(-7.61)	86.48(+8.08)	-8.64(-8.07)	7.47(+9.54)	57.51(-7.52)	86.41(+5.87)	-9.09(-8.18)	7.7(+8.93)
MTL	64.71	85.86	—	—	66.10	80.64	—	—

Table 12: Evaluation scores of various baselines, ablations and MTL for 10 tasks setting are shown. **First-**, **second-**, and **third-** best performing methods are highlighted over all baselines and ablations. MTL is omitted from ranking since it is considered as an ideal target performance reference for a CFT model to achieve. Performance difference of TF-SSD (2T + 1A) with top method (DR_{task}) in the baselines is shown in parentheses.

Input	Stage Model Response	Score	MoT Response	Final Refined Response
<p>Paragraph: The more charged the objects are, the stronger the electrostatic force between them. Question: If Mona is manipulating the charge of a couple of her objects, what happens to the electrostatic force if she increases the charge of one of the objects? Choices: A)increase B)decrease</p>	increase	4	-	increase
<p>Context: What a melodious sound they made . If one of us could n't transport them to an event , the other would . One summer when I was in the doctoral program and we only had one car (the one my husband drove to work) Carol did all the transporting of my kids . Whenever someone commented that they wished they could get their kids involved but they just did n't have time , Carol would say , Its either now or later . We were not stay - at - home moms . Question: Why did they only have one car at the time ?</p>	They didn't have enough money to buy another car.	1	The family's single-vehicle circumstance was a result of the narrator's commitment to a doctoral program and the husband's need for transportation to his workplace, which necessitated the use of their sole car.	The family's single-vehicle circumstance was a result of the narrator's commitment to a doctoral program and the husband's need for transportation to his workplace, which necessitated the use of their sole car.

Table 13: Refinement of samples belonging to different scorer category.

Input	Stage Model Response	Score	MoT Response	Final Refined Response
Context: They asked us to come off rudely , which was needless as we were hardly going to argue . Then a third woman arrived , though she was a real copper and treated us really respectfully and politely . She informed us that she 'd have to take our names because she had stopped us , just for formalities . So we did so , with a laugh and joke , before I said , ' So , common ? ' . Question: What leads the speaker to refer to the third officer as a ' real copper ' ?	She was polite and respectful.	2	The third officer's conduct, characterized by respect and professionalism, stands out to the speaker as embodying the true essence of law enforcement, prompting them to acknowledge her as a 'real copper'.	The third officer's conduct, characterized by respect and professionalism, stands out to the speaker as genuine and authentic.
Context: Since then we have n't done anything like that again but I admit I do sort of miss it . But now it was time for us to play a joke on them and we refused to fail . Tricky , you go set up the camera on the other side of the pool and make sure you can get everything on tape . Panda , give me a few minutes to set up first ok ? Sure . Question: What will happen after they set up the camera ?	They will film their friends.	3	Once the camera is set, they'll initiate their prank, meticulously orchestrating every detail to ensure the unfolding mischief is immortalized in the recording, all for the sake of a good-humored jest.	Once the camera is set, they'll initiate the prank.

Table 14: Refinement of samples belonging to different scorer category. **Highlighted blue part is the prefix added from MoT response**