

Data-Efficient Adaptation to Contextual Shifts in LLM-based Conversational Recommendation

Hyeongjun Yang, Donghyun Kim, Seokju Hwang, Midan Shim, KyuHwan Yeom, Kaehyun Um, Kyong-Ho Lee*

Department of Computer Science, Yonsei University
{edbm95, dhkim92, hsjtjrwn, midans26, tommal121, khyun33, khlee89}@yonsei.ac.kr

Abstract

Large language model (LLM)-based conversational recommender systems (CRSs) have demonstrated strong capabilities in capturing user preferences and providing contextually relevant recommendations. Nevertheless, the recommendation quality of the models frozen after training inevitably degrades under contextual shifts, such as changes in language and social trends. While periodic model updates are essential to maintain alignment with real-world preferences, training on large-scale data incurs substantial costs. This motivates data-efficient adaptation. However, existing data selection methods struggle to distinguish learnable samples under contextual shifts. To address this, we propose **Contextual Shift-Adaptive Data Pruning and Training (CAPT)**, a framework agnostic to underlying LLM-based CRSs. Specifically, we conceptualize a three-class data taxonomy comprising familiar, valuable, and outlier samples to formalize data behavior under contextual shifts. Based on this taxonomy, we design an importance score estimation scheme that quantifies a sample’s relative learnability for shift adaptation. Leveraging these importance scores, CAPT prioritizes highly learnable samples and further guides shift-adaptive training to actively steer the model toward current contexts. Experiments on three CRS benchmarks with real-world temporal splits demonstrate that CAPT outperforms baselines, matching or surpassing full-data fine-tuning performance using only 10–50% of the training data.

1 Introduction

Conversational recommender systems (CRSs) (Jan-nach et al., 2021) capture user preferences through free-form conversations and recommend relevant items. Recent studies on CRSs (He et al., 2023; Yang and Chen, 2024; He et al., 2025) have focused on leveraging large language models (LLMs),

*Corresponding author.

Necessity of Periodic CRS Updates

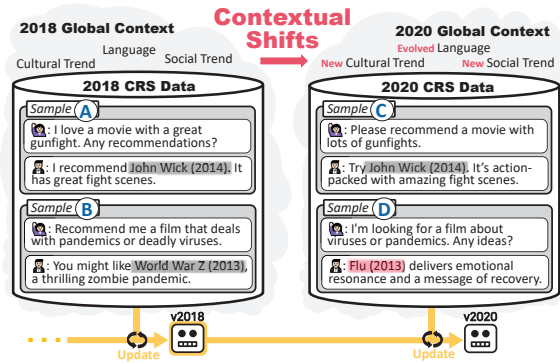


Figure 1: Contextual shifts (*i.e.*, evolving language & changing social/cultural trends) alter optimal recommendations, necessitating periodic CRS updates. Under new contexts, the recommendation in *Sample D* becomes preferable to that in *Sample B*.

which possess human-like reasoning capabilities and extensive world knowledge, to provide contextually appropriate and accurate recommendations.

Despite these advantages, a model frozen after training gradually diverges from real-world contexts, resulting in suboptimal recommendations (Koren, 2009). In other words, the optimal recommendation even for a similar request may vary across contextual shifts, which refer to the evolution of language and changes in social/cultural trends (Adomavicius et al., 2005; Hamilton et al., 2016). As illustrated in Figure 1, in some cases, the same item remains an appropriate recommendation (*e.g.*, *Sample A* → *C*), whereas in others, the suitable recommendation changes (*e.g.*, *Sample B* → *D*). To elaborate on the latter case (*i.e.*, *Sample B* → *D*), before the COVID-19 pandemic, a CRS might have recommended ‘*World War Z*’, which offers a thrilling depiction of fictional scenarios. However, general users after the pandemic may no longer perceive pandemic-themed content purely as entertainment. Instead, such topics can now evoke

Data Learnability in CRS Updates

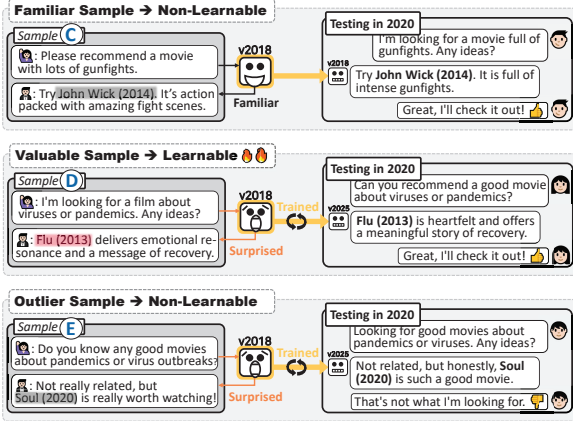


Figure 2: Data learnability from the perspective of the previously frozen model, where the familiar sample \textcircled{C} is non-learnable, while the high-surprisal samples \textcircled{D} and \textcircled{E} can either be valuable for CRS updates or outliers to be excluded, as revealed in user feedback from testing.

emotional resonance and messages of recovery, reflecting a shift in public perception shaped by the hardships of the pandemic. Consequently, a film like ‘*Flu*’, which embodies such characteristics, might become more appealing to general users. Accordingly, it is essential for CRSs to be periodically updated to reflect these contextual shifts.

However, the significant cost of updating CRSs, driven by LLM parameters and extensive data, necessitates cost-efficient fine-tuning techniques. The techniques comprise parameter-efficient fine-tuning (PEFT) (Dettmers et al., 2023; Hu et al., 2022; Li and Liang, 2021), which modifies only a subset of parameters (e.g., Adapters, LoRA), and data-efficient fine-tuning (DEFT) (Wang et al., 2024a; Mindermann et al., 2022; Killamsetty et al., 2021), which selects informative samples. Recognizing that these methods are often complementary, our study focuses on DEFT¹.

To address data-efficient adaptation of CRSs under contextual shifts, we conceptualize a three-class taxonomy of data behavior: **Familiar**, **Valuable**, and **Outlier**. From the model’s perspective on data learnability, as shown in Figure 2, data samples (e.g., Sample \textcircled{C}) where the same item remains appropriate are already familiar to the previously frozen model and thus non-learnable. In contrast, the model exhibits high surprisal for Sample \textcircled{D} , which is essential for adapting to the new contexts.

¹We further examine the compatibility of our DEFT approach with the PEFT approach (i.e., LoRA) in Figure 6.

While outliers (e.g., Sample \textcircled{E}) also induce high surprisal, they should be excluded because training on them yields negative feedback at test time.

Based on this understanding of data behavior under contextual shifts, existing DEFT methods exhibit the following limitations. Gradient-based approaches (Wang et al., 2024a; Xia et al., 2024; Killamsetty et al., 2021) estimate data importance by the gradient similarity with validation samples, which makes them effective at excluding outliers. However, these criteria are inherently limited in estimating data learnability of familiar samples, as they do not explicitly assess whether a sample is already well predicted by the model. In addition, they incur per-sample gradient storage overhead. In contrast, loss-based approaches (Zheng et al., 2023; Mindermann et al., 2022; Paul et al., 2021; Jiang et al., 2019) deprioritize familiar samples and avoid the cost of gradient storage. Nevertheless, they are vulnerable to outliers, as samples with large absolute losses can dominate selection. This issue becomes particularly severe under shift adaptation, where the loss distribution spans a wide scale.

Motivated by these observations, we propose a loss-based data importance estimation method that explicitly incorporates the three-class data taxonomy and is robust to the loss-scale disparities induced by the shift. Using the importance scores, we prioritize valuable samples for shift adaptation.

However, when fine-tuning on a small subset of selected samples, the standard CRS training objective often leads to insufficient adaptation to the shift. Unlike factual tasks (e.g., QA), where outdated facts require binary correction, CRS supervision under contextual shifts entails adjusting the relative relevance of items given a user’s request. Consequently, effective adaptation requires reshaping the recommendation distribution to prioritize currently preferred items while demoting those with faded appeal. To actively steer the model’s relative likelihoods without counterfactual supervision data (e.g., “Recommend B instead of A”), we propose utilizing our importance scores for loss weighting.

Furthermore, random negative sampling in conventional CRS training often yields trivial negatives that are easily distinguishable from the ground-truth item. Thus, we incorporate samples with high importance scores into the negative set, thereby enhancing the model’s discriminative ability.

To this end, we propose **Contextual Shift-Adaptive Data Pruning and Training (CAPT)**, a framework agnostic to underlying LLM-based

CRSs. Our contributions are summarized as follows: (1) This is the first work to formulate the problem of data-efficient adaptation of CRSs under contextual shifts. (2) We propose a robust data importance estimation scheme grounded in the data taxonomy for shift adaptation. (3) We develop a shift-adaptive training strategy incorporating loss weighting and negative sampling that actively drives adaptation to contextual shifts. (4) On three CRS benchmarks with temporal splits, CAPT consistently outperforms state-of-the-art data selection methods. Notably, CAPT matches or surpasses full-data fine-tuning using only 10–50% of the data.

2 Related Work

2.1 Conversational Recommender System

Knowledge graphs (KGs) have been widely employed to incorporate domain-specific knowledge (Bizer et al., 2009; Yang et al., 2024; Won et al., 2023). A CRS aims to infer user preferences from conversations and match them to relevant items. Early CRS approaches (Yang et al., 2023; Zhou et al., 2020; Chen et al., 2019; Yang et al., 2025) primarily infer user preferences from conversations and match them to items in KGs. However, these KG-based CRSs rely on the assumption that entity recognition and linking from utterances are already in place, which limits their practical deployment.

To overcome this limitation, PECS (Ravaut et al., 2024) aligns items, item-attributes, and dialogue context within a single language model. This approach facilitates end-to-end training and inference, seamlessly integrating the entire process from dialogue understanding to recommendation. More recently, LLMs with human-like reasoning capabilities and extensive world knowledge have significantly advanced CRS research (He et al., 2023). ReFICR (Yang and Chen, 2024) adopts LLM instruction tuning, providing more accurate and nuanced recommendations. RTA (He et al., 2025) incorporates an aggregator for better control of LLM-based recommendations.

Despite these advances, LLM-based CRSs may generate suboptimal recommendations when the knowledge learned during training no longer aligns with current real-world contexts, highlighting the need for efficient shift adaptation strategies.

2.2 Adapting Language Models to Shifts

Language models require adaptation not only to new tasks, but also to evolving language usage and

social trends within the same task (Zheng et al., 2025; Lazaridou et al., 2021). Continual learning (CL) (Wu et al., 2024) is a principal approach for adapting to new tasks while retaining past knowledge. Replay-based methods (M’hamdi and May, 2024; Wang et al., 2024b) in CL reuse past representative samples to preserve prior knowledge during current training. However, these approaches are not universally beneficial. This is because past data may reflect outdated language use and trends, which can hinder the model’s adaptation to the current context (Jin et al., 2022). In such scenarios, sequential fine-tuning is effective, as this approach incrementally trains on newly emerging data without replaying prior data (Zhao et al., 2024; Loureiro et al., 2022; Luu et al., 2022; Röttger and Pierrehumbert, 2021). While previous works (Mitchell et al., 2022; Meng et al., 2023; Wang et al., 2025) have investigated parameter efficiency, data efficiency remains underexplored. In this study, we focus on data efficiency in sequential fine-tuning.

2.3 Data Selection

Data selection aims to identify informative samples from large datasets to enhance performance or training efficiency using heuristic and optimization-based criteria. Recent heuristic methods (Wettig et al., 2024; Chen et al., 2024) employ LLMs as proxies for human evaluators to assess data quality. However, they are limited in identifying samples that directly contribute to model optimization.

Optimization-based approaches can be categorized by their reliance on the Hessian matrix, gradient, or loss. Hessian-based (Koh and Liang, 2017; Park et al., 2023; San Joaquin et al., 2024) and gradient-based (Killamsetty et al., 2021; Garima et al., 2020; Mirzasoleiman et al., 2020) methods require computing Hessians and full-parameter gradients, respectively. Although these approaches offer theoretical performance guarantees, they entail higher computational costs than full-dataset fine-tuning. Therefore, they are not suitable for achieving cost-efficient updates.

To mitigate this, recent works restrict the parameter space for importance estimation. DEALRec (Lin et al., 2024) approximates influence scores by computing Hessians from a smaller model. GREATS (Wang et al., 2024a) and LESS (Xia et al., 2024) leverage gradients from the last layer and the LoRA (Hu et al., 2022) modules, respectively. Despite their efficiency gains, these methods still require gradient storage overhead. Furthermore,

their criteria do not account for whether a sample is already predicted correctly by the model.

In comparison, loss-based methods (Paul et al., 2021; Zheng et al., 2023; Jiang et al., 2019; Mindermann et al., 2022) prioritize samples that the model still needs to learn, even without storage overhead. However, they are susceptible to outliers with large loss values. This motivates the need for a selection criterion that robustly prioritizes valuable samples under varying loss scales.

3 Method: CAPT

3.1 CRS Task under Contextual Shifts

Given a conversation x , a CRS parameterized by θ aims to recommend the most appropriate item y . In real-world scenarios, however, CRSs face a non-stationary environment where data accumulates sequentially, accompanied by emerging contexts and shifting general preferences. Accordingly, we formulate the task as a sequence of learning phases over time to mitigate performance degradation caused by the shift. Let D_t and D_{t+1} represent datasets collected in consecutive temporal phases. Our objective is to effectively adapt the parameters θ pre-trained on the historical dataset D_t to the context of the subsequent dataset D_{t+1} .

3.2 Revisiting Data Selection for Robust Shift Adaptation

For data-efficient shift adaptation, samples can be categorized into familiar, valuable, and outlier. Familiar samples are correctly predicted by the model and are unnecessary for training. In contrast, although both valuable and outlier samples induce high surprisal in the model, valuable samples are beneficial for adapting to new contexts.

Gradient-based Criteria. To mitigate the high computational cost of full-parameter gradient computation (Killamsetty et al., 2021; Garima et al., 2020), recent works such as GREATS (Wang et al., 2024a) and LESS (Xia et al., 2024) approximate importance using gradients from a subset of parameters $\mathbf{W} \subset \theta$, such as the last layer or LoRA modules (Hu et al., 2022). Concretely, they compute gradients of \mathbf{W} for a training sample $(x, y) \in D_{t+1}$ and a validation sample $(x', y') \in D_{\text{val}}$, and measure their alignment:

$$\sum_{(x', y') \in D_{\text{val}}} \frac{\partial \mathcal{L}(y | x; \theta_{D_t})}{\partial \mathbf{W}} \cdot \frac{\partial \mathcal{L}(y' | x'; \theta_{D_t})}{\partial \mathbf{W}}, \quad (1)$$

where \cdot denotes cosine similarity. A high similarity indicates that updating the model with the candidate sample (x, y) is likely to improve validation performance, rendering this approach effective at excluding outliers. However, familiar samples may receive high importance scores when similar samples are present in the validation set.

Loss-based Criteria. These methods (Zheng et al., 2023; Paul et al., 2021; Jiang et al., 2019) assume that samples with high training loss \mathcal{L}_t provide stronger learnability signals, deprioritizing familiar samples. However, high-loss samples inevitably include outliers. RHO (Mindermann et al., 2022) incorporates a validation set into importance estimation to exclude outliers. Specifically, the importance of a candidate sample (x, y) is estimated by fine-tuning a model pre-trained on D_t with (x, y) added, and then measuring how much this update reduces the loss on a validation set D_{val} . Using Bayes' rule, the objective can be approximated as:

$$\begin{aligned} & \sum_{(x', y') \in D_{\text{val}}} \log p_{\theta}(y' | x'; D_t \cup \{(x, y)\}) \\ & \propto \underbrace{\mathcal{L}(y | x; \theta_{D_t})}_{\mathcal{L}_t} - \underbrace{\mathcal{L}(y | x; \theta_{D_t \cup D_{\text{val}}})}_{\mathcal{L}_p} \end{aligned} \quad (2)$$

Detailed derivations are provided in Appendix A. While the proxy loss \mathcal{L}_p helps exclude outliers in in-distribution settings, its effectiveness degrades under shift adaptation where loss scales vary widely. Thus, outliers are often prioritized due to their absolute loss magnitude, as they tend to exhibit both high training loss and high proxy loss.

Diagnostic Analysis. To empirically verify these limitations, we conduct a pilot study using the LLM-based CRS (He et al., 2025) on the Reddit dataset (He et al., 2023), considering the transition from D_t to D_{t+1} . We identify 300 familiar samples representing redundant context and artificially inject 300 outliers featuring mismatched user requirements. Detailed experimental setups are provided in Appendix B. As shown in Table 1, gradient-based criteria (*i.e.*, GREATS and LESS) fail to deprioritize familiar samples, including over 38.6% such samples within the top 20% selection. In contrast, loss-based RHO selection is highly susceptible to noise, with 63.6% outliers being prioritized due to their high loss magnitudes. A more comprehensive comparison across varying data budgets, including our proposed method, is

Criterion	Familiar (%)	Outlier (%)
GREATS	38.6	7.3
LESS	45.3	4.3
RHO	3	63.6

Table 1: Inclusion percentages of familiar and outlier samples within the top 20% priority pool. Lower values indicate better importance weighting performance.

provided in Figure 3. Consequently, effective shift adaptation requires a data selection mechanism that prioritizes high-surprisal samples while robustly excluding outliers, regardless of loss scale. Appendix C presents an empirical analysis of the training loss scales and the prevalence of familiar samples across varying temporal granularities.

3.3 Importance Score Estimation

Recognizing that prior absolute-loss metrics are sensitive to the loss-scale variations induced by contextual shifts, we formalize a three-class taxonomy of data behavior by comparing the training loss \mathcal{L}_t and the proxy loss \mathcal{L}_p . Specifically, \mathcal{L}_t denotes the loss incurred by the target model on the training sample $(x, y) \in D_{t+1}$. \mathcal{L}_p refers to the loss on the same training sample, computed by a version of the target model that has undergone additional training² on the validation set D_{val} at temporal phase $t + 1$. The taxonomy is as follows:

- **Familiar:** $\mathcal{L}_t \approx \mathcal{L}_p \approx \epsilon$
- **Valuable:** $\mathcal{L}_t \gg \mathcal{L}_p, \mathcal{L}_p \approx \epsilon$
- **Outlier:** $\mathcal{L}_t \approx \mathcal{L}_p, \mathcal{L}_p \gg \epsilon$

Here, ϵ denotes a reference scalar representing the expected loss magnitude of a well-fitted sample. We use \approx to indicate that two values are of comparable magnitude within the same scale.

Grounded in this taxonomy, we design an importance score $\mathcal{I}_{(x,y)}$ that prioritizes valuable samples over familiar and outlier samples in a scale-invariant manner:

$$\mathcal{I}_{(x,y)} = \frac{\mathcal{L}_t}{\mathcal{L}_p}. \quad (3)$$

Intuitively, this criterion functions as a signal-to-noise ratio, quantifying the learnability relative to the sample’s inherent difficulty. The score is maximized for valuable samples, where $\mathcal{L}_t \gg \mathcal{L}_p$, while

²All training procedures and loss computations share the same objective function, which follows the original implementation of each CRS model.

naturally remaining low for familiar and outlier samples regardless of absolute loss magnitude.

Data Pruning. Leveraging this criterion, we prune the full dataset D_{t+1} by selecting the top- k :

$$D_p = \text{Top-}k \left(\{(x, y) \in D_{t+1}\}, \text{ based on } \mathcal{I}_{(x,y)} \right). \quad (4)$$

Fine-tuning on the pruned dataset D_p focuses the model on samples most critical for shift adaptation.

3.4 Shift-Adaptive Training

Although the implementation details of the loss function vary slightly across CRS studies, LLM-based CRSs (He et al., 2025; Yang and Chen, 2024; Ravaut et al., 2024) are typically trained with InfoNCE (Gutmann and Hyvärinen, 2010), a standard contrastive objective that aligns the conversation with the appropriate item while pushing it away from negative items:

$$p_\theta(y|x) = \frac{e^{s(x,y)}}{e^{s(x,y)} + \sum_{y_n \in \mathcal{N}} e^{s(x,y_n)}}, \quad (5)$$

where $e^{s(\cdot)}$ denotes an exponentiated similarity score, and \mathcal{N} is a set of negative items y_n obtained through random sampling.

Loss Weighting. The standard objective is insufficient to actively steer the model’s relative likelihoods. Thus, we regard the importance score $\mathcal{I}_{(x,y)}$ as a measure of the ‘knowledge gap’ between the data sample (x, y) and the target model. To update parameters in proportion to this gap, we incorporate the score as a weight for each sample’s loss:

$$\mathcal{L}_{\text{weight}} = -\mathbb{E}_{(x,y) \sim D_p} [\mathcal{I}_{(x,y)} \log p_\theta(y|x)]. \quad (6)$$

Notably, we do not incorporate additional parameter regularization terms (e.g., KL divergence) and strictly adhere to the fine-tuning iteration counts and other hyperparameters specified in the original LLM-based CRS implementations. Empirical observations confirm that training remains stable under these settings without exhibiting excessive parameter drift (Detailed discussion in Appendix D). **Negative Sampling.** Conventional CRS studies rely on random negative sampling, which is insufficient for improving the model’s ability to make fine-grained distinctions among candidates for shift adaptation. Thus, we include the items of high-importance samples in the negative set \mathcal{N} , controlled by a sampling ratio τ . The remaining $(1 - \tau)$ portion of negatives is randomly sampled from the

entire item list to preserve generalizability. The optimal value for the mixing ratio τ is validated in the experimental results shown in Figure 7.

4 Experimental Setting

We evaluate recommendation performance using $\text{Recall}@N$ for hit rate and $\text{NDCG}@N$ and $\text{MRR}@N$ for ranking accuracy, with $N \in \{1, 10, 50\}$. Implementation details of CAPT are provided in Appendix E.

4.1 CRS Benchmark Datasets

The **Reddit** dataset (He et al., 2023) is constructed from **Reddit**³, where users engage in discussions about movie recommendations. Unlike other CRS benchmarks, the Reddit dataset provides timestamps for each dialogue, covering a time span from June 2018 to December 2022. To simulate periodic model updates to temporal shifts, we sorted the data chronologically and divided it into three equally-sized subsets, each containing 250,214 utterances. The first subset, **Reddit_{pretrain}**, is used for pre-training the model to acquire basic CRS capabilities. The remaining two, **Reddit_{phase1}** and **Reddit_{phase2}**, are used to evaluate how effectively the model adapts to subsequent splits. Other standard CRS benchmarks, such as **REDIAL** (Li et al., 2018) and **INSPIRED** (Hayati et al., 2020), do not provide temporal information. Therefore, we used them to assess ‘task-level adaptation efficiency’, focusing on how effectively a general-purpose LLM can be adapted to the CRS task (see Appendix F for more detail).

4.2 CRS Models

We evaluate our framework on three representative CRS models. **PECRS** (Ravaut et al., 2024): While utilizing GPT-2 as its backbone, PECRS was the first work to integrate the entire CRS task into a single language model. **ReFICR** (Yang and Chen, 2024): It leverages GRITLM-7B (Muennighoff et al.) and adopts instruction tuning. **RTA** (He et al., 2025): Built upon Llama-7B, RTA combines the LLM with an aggregator.

4.3 Data Selection Baselines

(1) **Full** refers to fine-tuning a model on the entire training set without data pruning. It serves as the primary reference for performance. (2) **Random** chooses a random subset of the training data, equal

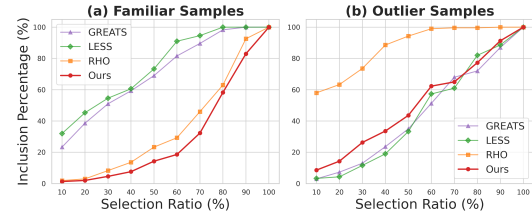


Figure 3: Accumulated inclusion percentages of (a) familiar samples and (b) outliers across varying selection ratios in the pilot study.

in size to the other methods’ pruned set. It helps verify the benefit of intelligent data selection. (3) **GREATS** (Wang et al., 2024a) reduces computational overhead by using only the gradients from the last layer. It selects data samples that are expected to yield improvements on a validation set. (4) **LESS** (Xia et al., 2024) computes gradients by leveraging the PEFT module (i.e., LoRA), thereby reducing the cost of gradient-based selection. (5) **MaxLoss** (Loshchilov and Hutter, 2016) prioritizes samples with the highest training loss, focusing on the most surprising examples the model has not yet learned. (6) **RHO** (Mindermann et al., 2022) selects samples based on the absolute difference between the training loss and the loss produced by a model further trained on a validation set.

To examine compatibility with PEFT, we also apply LoRA to our method and other baselines, and compare these variants with LESS in Figure 6.

5 Experimental Results

Extended Pilot Study on Familiar and Outlier Samples. To further investigate our pilot study (Section 3.2), we analyzed the accumulated inclusion percentages of familiar samples and outliers. As illustrated in Figure 3, the results indicate that gradient-based criteria exhibit significant limitations in deprioritizing familiar samples, while the loss-based RHO method is highly susceptible to outliers. Importantly, our proposed CAPT framework remains robust to both types of samples.

Data-Efficient Adaptation under Contextual Shifts across Different CRS Models. Figure 4 and Table 2 present the performance of three CRS models with all data-selection methods and the detailed results of the main methods on **Reddit_{phase1}**, respectively. The first row of Figure 5 presents the performance of the RTA model on **Reddit_{phase2}**. Random selection generally shows that performance improves as more data is

³<https://www.reddit.com/>

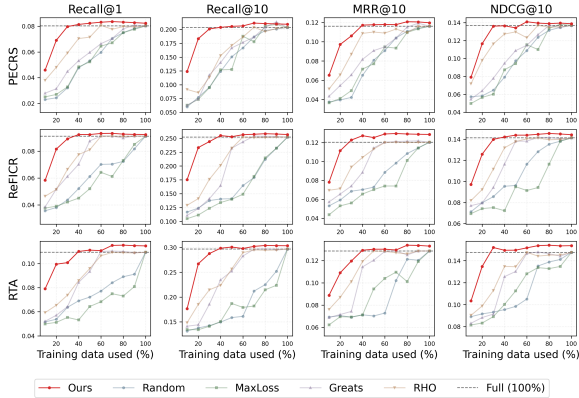


Figure 4: Performance comparison of data selection methods on the $\text{Reddit}_{\text{Phase1}}$ dataset. Rows denote CRS models (PECRS, ReFICR, RTA), and columns denote evaluation metrics. The x-axis shows the percentage of training data used.

	Method	Recall@10		MRR@10		NDCG@10	
		$\geq \text{Full}$	Best	$\geq \text{Full}$	Best	$\geq \text{Full}$	Best
PECRS	Full	-	0.2037	-	0.1161	-	0.1369
	GREATS	80%	0.2042	90%	0.1193	80%	0.1370
	RHO	70%	0.2041	100%	0.1161	70%	0.1369
	Ours	40%	0.2116	40%	0.1208	50%	0.1410
ReFICR	Full	-	0.2521	-	0.1204	-	0.1412
	GREATS	70%	0.2539	60%	0.1218	70%	0.1417
	RHO	60%	0.2521	70%	0.1204	70%	0.1412
	Ours	40%	0.2583	30%	0.1296	40%	0.1454
RTA	Full	-	0.2969	-	0.1287	-	0.1476
	GREATS	80%	0.2971	70%	0.1287	70%	0.1479
	RHO	100%	0.2969	50%	0.1287	60%	0.1476
	Ours	40%	0.3037	40%	0.1339	30%	0.1541

Table 2: Performance comparison on the $\text{Reddit}_{\text{Phase1}}$ dataset. ‘ $\geq \text{Full}$ ’ indicates the percentage of data required to match/exceed the performance of full dataset. ‘Best’ denotes the best performance of each method.

used. MaxLoss occasionally performs worse than Random, showing that focusing only on samples with high training loss can harm performance. This highlights the use of a validation set in RHO for excluding outliers. The gradient-based method, GREATS, consistently outperforms Random selection, showing the benefit of validation gradients. CAPT consistently outperforms all baselines and achieves full-dataset fine-tuning performance using only 30–50% of the data. Interestingly, performance often declines slightly when more than 70% of the data is used. This indicates that selecting informative data is more effective than using the entire dataset. Additional results on @50 metrics and a different temporal granularity are provided in Appendix G.

CRS Task Adaptation on Small Datasets. The second and third rows of Figure 5 show the results on the REDIAL and INSPIRED datasets,

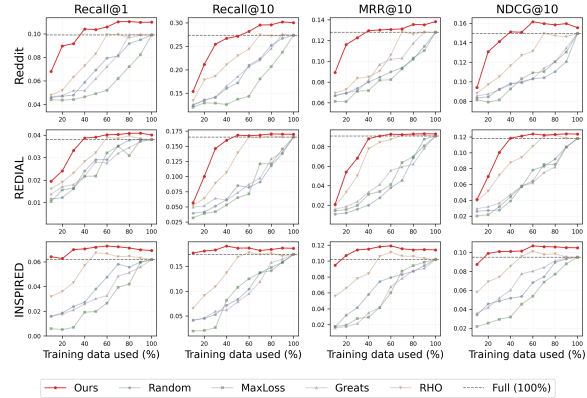


Figure 5: Performance comparison using the RTA model. Each row corresponds to a different dataset: $\text{Reddit}_{\text{Phase2}}$ for temporal shifts adaptation, and REDIAL and INSPIRED for CRS task adaptation.

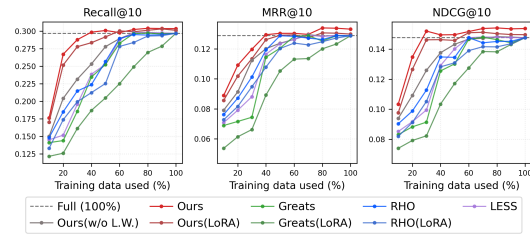


Figure 6: Performance comparison of data selection methods with PEFT (*i.e.*, LoRA) on the RTA model, using the $\text{Reddit}_{\text{Phase1}}$ dataset, with an ablation on importance-guided loss weighting (*i.e.*, w/o L.W.).

respectively. Despite the smaller dataset sizes compared to Reddit, our method achieves strong performance using pruned data. Notably, on INSPIRED, CAPT surpasses the performance of full-dataset fine-tuning with only 10% of the data. As the data ratio increases to 50%, CAPT continues to show clear performance gains. Beyond this point, the improvements become smaller, and performance may slightly decline. This observation can be attributed to the remaining samples contributing little to adaptation.

Compatibility with PEFT. Figure 6 shows the performance comparison of our method and the baselines when augmented with LoRA. For this experiment, the gradient-based methods (*i.e.*, GREATS and LESS) utilize gradients from the LoRA modules for data selection, while the loss-based methods (*i.e.*, RHO and CAPT) estimate proxy loss by fine-tuning on a validation set using LoRA. GREATS shows a significant performance drop in this setting. Since GREATS already relies on a limited signal from only the

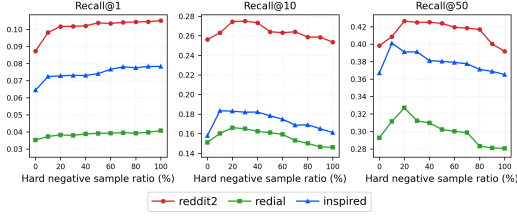


Figure 7: Analysis of the negative sampling ratio (τ) on the RTA model using the three datasets.

final layer’s gradients, the additional constraints of LoRA further reduce its representational capacity for effective data selection. LESS performs slightly better than the original GREATS since it computes gradients for the full model parameters within the LoRA modules, providing a richer signal. Despite a slight reduction in the proxy model’s representational capability due to LoRA, CAPT with LoRA still outperforms all baselines.

Effect of Loss Weighting and Negative Sampling.

Figure 6 also includes an ablation study (*i.e.*, w/o L.W.) on the contribution of our importance-guided loss weighting. The results show that assigning larger weights to high-importance samples leads to better performance. Notably, even without the loss weighting strategy, our data selection criterion outperforms the baselines, reaching full-dataset performance earlier. Figure 7 presents the performance variation according to the negative sampling ratio τ . Across all metrics, there is a significant drop in performance when $\tau=0$, which corresponds to removing the importance-guided negative sampling. This result clearly demonstrates the effectiveness of this strategy. Interestingly, different metrics show different trends as τ increases. For Recall@1, performance improves steadily as the ratio increases from 20% to 100%. However, for Recall@10 and Recall@50, performance peaks at a ratio of 20–30% before gradually declining. This suggests that while a high ratio of important negatives forces the model to develop an excellent ability to distinguish the single best item, it may slightly hinder its ability to generalize for a broader set of top-N recommendations.

Computational Complexity. Table 3 presents a detailed comparison of the computational and storage costs. The reported runtimes are measured under a LoRA-based setting, and they scale linearly with the corresponding computation factors. While

Method	Importance Score Computation	Computation Time (hrs)	Storage for Gradients
Greats	$O((N + M)T + MNd)$	3.26	$O(Md)$
LESS	$O((N + M)T + MNd)$	6.85	$O(Md)$
RHO	$O(MT + 2NF)$	0.97	-
Ours	$O(MT + 2NF)$	0.97	-

Table 3: Computation and storage costs for importance score estimation in each data-selection method. N =|training set|, M =|validation set|, F =forward pass cost, T =training step cost (*e.g.*, full-parameter, last layer, or LoRA), and d =gradient dimension. Time is measured in wall-clock hours under the same settings.

conventional gradient-based methods are computationally prohibitive for LLMs, GREATS and LESS leverage PEFT to reduce this overhead. As a result, their computational costs have become broadly comparable to those of RHO and CAPT. Specifically, for gradient-based methods, the cost involves calculating gradients for each sample in the training and validation sets $O((N + M)T)$ and then computing the influence of each training sample on the validation samples $O(MNd)$. On the other hand, loss-based methods require training a proxy model $O(MT)$ and then obtaining losses from both the proxy and target models, $O(2NF)$. Excluding the common $O(MT)$ term shared by all methods, the computational cost of $O(NT + MNd)$ is comparable to the $O(NF)$ of RHO and Ours. However, GREATS and LESS require $O(Md)$ for gradient storage, which loss-based methods avoid.

6 Conclusion

In this work, we introduced the challenge of data-efficient adaptation for LLM-based CRSs under contextual shifts. By conceptualizing a data taxonomy, we identified the fundamental limitations of existing data selection methods for shift adaptation in the CRS task. To address this, we proposed CAPT, a framework that incorporates a scale-invariant importance scoring scheme grounded in the data taxonomy and shift-adaptive training strategies. Experimental results demonstrate that CAPT surpasses full-data performance using 10–50% of the training data, reducing the computational burden of periodic model updates. These results suggest that more data is not always better for shift adaptation. Instead, prioritizing valuable samples with adaptive weighting to bridge the gap between prior knowledge and changing trends is effective in maintaining recommendation relevance.

Limitations

Focus on Global Preference. Our framework specifically addresses shifts in global preferences rather than turn-level preference changes within a single dialogue session. While modeling fine-grained changes in a dialogue is a core component of conversational recommendation, a CRS’s ability to serve cold-start users or navigate initial dialogue turns depends heavily on its alignment with contemporary global preferences. We treat global adaptation as a foundational step that stabilizes the backbone CRS model, ensuring that the system is better prepared to interpret and incorporate personalized signals as they accumulate over future interactions. A natural extension of this work is to develop a multi-granular adaptation framework that jointly optimizes both global and turn-level preference shifts.

Long-tail Items. Our error analysis of CAPT reveals that a frequent failure case occurs when the ground-truth item is a long-tail item that is rarely observed during training. Because such items receive limited exposure, their representations are often insufficiently trained, making them difficult to rank correctly even after the application of shift-adaptive training. Consequently, achieving robust performance on long-tail recommendations remains an open challenge for data-efficient adaptation.

Applicability to RAG-augmented LLM-based CRSs. Recent studies on LLM-based CRSs (Li et al., 2025; Xi et al., 2024) have explored augmenting the model with Retrieval-Augmented Generation (RAG) to incorporate factual information and user history. Such RAG-augmented designs improve the system’s ability to access up-to-date knowledge and contextual evidence beyond the model’s parametric knowledge. However, the LLM’s ability to reason over and rank those retrieved items still hinges on its fine-tuned parameters. Therefore, our work focuses on the backbone LLM, which serves as the fundamental decision-maker in a CRS. Extending CAPT to jointly account for the RAG components for retrieval distribution shifts or user memory changes is a promising future work.

Dependence on Validation Set Quality. Our importance scoring relies on a validation set to estimate proxy loss. If the validation data is noisy or

misaligned with the target environment, the estimated scores may be less reliable. However, this limitation is not unique to CAPT. Most data selection techniques and even standard machine learning pipelines inherently depend on the representativeness of the validation data. Moreover, our experiments across multiple benchmarks suggest that even a small validation set is sufficient for CAPT to provide meaningful gains.

Ethical Considerations

All authors of this paper acknowledge the ACL code of ethics. For our experiments, we utilized publicly available datasets, REDIAL, INSPIRED and Reddit. The REDIAL and INSPIRED datasets were constructed through crowdsourcing on Amazon Mechanical Turk, ensuring that all data was collected anonymously. The Reddit dataset was crawled from the Reddit community and processed to exclude any personal information, maintaining the privacy of individual users. Although we have not observed any inappropriate or harmful content in the responses generated during our study, we inform that our framework leverages pre-trained LLMs, which may inherently produce biased or inappropriate responses under certain conditions.

References

- Gediminas Adomavicius, Ramesh Sankaranarayanan, Shahana Sen, and Alexander Tuzhilin. 2005. [Incorporating contextual information in recommender systems using a multidimensional approach](#). *ACM Trans. Inf. Syst.*, 23(1):103–145.
- Christian Bizer, Jens Lehmann, Georgi Kobilarov, Sören Auer, Christian Becker, Richard Cyganiak, and Sebastian Hellmann. 2009. [Dbpedia - a crystallization point for the web of data](#). *Web Semant.*, 7(3):154–165.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Sriniwasan, Tianyi Zhou, Heng Huang, and Hongxia Jin. 2024. [Alpagasus: Training a better alpaca with fewer data](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.
- Qibin Chen, Junyang Lin, Yichang Zhang, Ming Ding, Yukuo Cen, Hongxia Yang, and Jie Tang. 2019. [Towards knowledge-based recommender dialog system](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1803–1813.

- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *Advances in neural information processing systems*, 36:10088–10115.
- Garima, Frederick Liu, Satyen Kale, and Mukund Sundararajan. 2020. Estimating training data influence by tracing gradient descent. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20*, Red Hook, NY, USA. Curran Associates Inc.
- Michael Gutmann and Aapo Hyvärinen. 2010. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 297–304. JMLR Workshop and Conference Proceedings.
- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Cultural shift or linguistic drift? comparing two computational measures of semantic change. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2116–2121, Austin, Texas. Association for Computational Linguistics.
- Shirley Anugrah Hayati, Dongyeop Kang, Qingxi-aoyang Zhu, Weiyang Shi, and Zhou Yu. 2020. Inspired: Toward sociable recommendation dialog systems. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8142–8152.
- Zhankui He, Zhouhang Xie, Rahul Jha, Harald Steck, Dawen Liang, Yesu Feng, Bodhisattwa Prasad Majumder, Nathan Kallus, and Julian McAuley. 2023. Large language models as zero-shot conversational recommenders. In *Proceedings of the 32nd ACM international conference on information and knowledge management*, pages 720–730.
- Zhankui He, Zhouhang Xie, Harald Steck, Dawen Liang, Rahul Jha, Nathan Kallus, and Julian McAuley. 2025. Reindex-then-adapt: Improving large language models for conversational recommendation. In *Proceedings of the Eighteenth ACM International Conference on Web Search and Data Mining, WSDM '25*, page 866–875, New York, NY, USA. Association for Computing Machinery.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*.
- Dietmar Jannach, Ahtsham Manzoor, Wanling Cai, and Li Chen. 2021. A survey on conversational recommender systems. *ACM Comput. Surv.*, 54(5).
- Angela H Jiang, Daniel L-K Wong, Giulio Zhou, David G Andersen, Jeffrey Dean, Gregory R Ganger, Gauri Joshi, Michael Kaminsky, Michael Kozuch, Zachary C Lipton, and 1 others. 2019. Accelerating deep learning by focusing on the biggest losers. *arXiv preprint arXiv:1910.00762*.
- Xisen Jin, Dejiao Zhang, Henghui Zhu, Wei Xiao, Shang-Wen Li, Xiaokai Wei, Andrew Arnold, and Xiang Ren. 2022. Lifelong pretraining: Continually adapting language models to emerging corpora. In *Proceedings of BigScience Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models*, pages 1–16, virtual+Dublin. Association for Computational Linguistics.
- Krishnateja Killamsetty, Durga S, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. 2021. Grad-match: Gradient matching based data subset selection for efficient deep model training. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 5464–5474. PMLR.
- Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML'17*, page 1885–1894. JMLR.org.
- Yehuda Koren. 2009. Collaborative filtering with temporal dynamics. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '09*, page 447–456, New York, NY, USA. Association for Computing Machinery.
- Angeliki Lazaridou, Adhiguna Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liška, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d’Autume, Tomas Kocisky, Sebastian Ruder, Dani Yogatama, Kris Cao, Susannah Young, and Phil Blunsom. 2021. Mind the gap: assessing temporal generalization in neural language models. In *Proceedings of the 35th International Conference on Neural Information Processing Systems, NIPS '21*, Red Hook, NY, USA. Curran Associates Inc.
- Chuang Li, Yang Deng, Hengchang Hu, Min-Yen Kan, and Haizhou Li. 2025. ChatCRS: Incorporating external knowledge and goal guidance for LLM-based conversational recommender systems. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 295–312, Albuquerque, New Mexico. Association for Computational Linguistics.
- Raymond Li, Samira Ebrahimi Kahou, Hannes Schulz, Vincent Michalski, Laurent Charlin, and Chris Pal. 2018. Towards deep conversational recommendations. *Advances in neural information processing systems*, 31.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language*

- Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics.
- Dongding Lin, Jian Wang, Wenjie Li, abc, and abc. 2023. Cola: Improving conversational recommender systems by collaborative augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 4462–4470.
- Xinyu Lin, Wenjie Wang, Yongqi Li, Shuo Yang, Fuli Feng, Yinwei Wei, and Tat-Seng Chua. 2024. [Data-efficient fine-tuning for llm-based recommendation](#). In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '24*, page 365–374, New York, NY, USA. Association for Computing Machinery.
- Ilya Loshchilov and Frank Hutter. 2016. Online batch selection for faster training of neural networks. *International Conference on Learning Representations Workshop*.
- Daniel Loureiro, Francesco Barbieri, Leonardo Neves, Luis Espinosa Anke, and Jose Camacho-collados. 2022. [TimeLMs: Diachronic language models from Twitter](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 251–260, Dublin, Ireland. Association for Computational Linguistics.
- Kelvin Luu, Daniel Khashabi, Suchin Gururangan, Karishma Mandyam, and Noah A. Smith. 2022. [Time waits for no one! analysis and challenges of temporal misalignment](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5944–5958, Seattle, United States. Association for Computational Linguistics.
- Wenchang Ma, Ryuichi Takanobu, Minlie Huang, and abc. 2021. [Cr-walker: Tree-structured graph reasoning and dialog acts for conversational recommendation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1839–1851.
- Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. [Mass-editing memory in a transformer](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.
- Meryem M'hamdi and Jonathan May. 2024. [Leitner-guided memory replay for cross-lingual continual learning](#). 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 7808–7821, Mexico City, Mexico. Association for Computational Linguistics.
- Sören Mindermann, Jan M Brauner, Muhammed T Razak, Mrinank Sharma, Andreas Kirsch, Winnie Xu, Benedikt Höltgen, Aidan N Gomez, Adrien Morisot, Sebastian Farquhar, and Yarin Gal. 2022. [Prioritized training on points that are learnable, worth learning, and not yet learnt](#). In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 15630–15649. PMLR.
- Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. 2020. Coresets for data-efficient training of machine learning models. In *Proceedings of the 37th International Conference on Machine Learning, ICML'20*. JMLR.org.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022. [Fast model editing at scale](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*.
- Niklas Muennighoff, SU Hongjin, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, and Douwe Kiela. Generative representational instruction tuning. In *The Thirteenth International Conference on Learning Representations*.
- Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Mądry. 2023. [Trak: attributing model behavior at scale](#). In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*. JMLR.org.
- Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. 2021. Deep learning on a data diet: finding important examples early in training. In *Proceedings of the 35th International Conference on Neural Information Processing Systems, NIPS '21*, Red Hook, NY, USA. Curran Associates Inc.
- Mathieu Ravaut, Hao Zhang, Lu Xu, Aixin Sun, and Yong Liu. 2024. Parameter-efficient conversational recommender system as a language processing task. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 152–165.
- Paul Röttger and Janet Pierrehumbert. 2021. [Temporal adaptation of BERT and performance on downstream document classification: Insights from social media](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2400–2412, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ayrton San Joaquin, Bin Wang, Zhengyuan Liu, Nicholas Asher, Brian Lim, Philippe Muller, and Nancy F. Chen. 2024. [In2Core: Leveraging influence functions for coreset selection in instruction finetuning of large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 10324–10335, Miami, Florida, USA. Association for Computational Linguistics.
- Jiachen Tianhao Wang, Tong Wu, Dawn Song, Prateek Mittal, and Ruoxi Jia. 2024a. [Greats: Online selection of high-quality data for llm training in every](#)

- iteration. *Advances in Neural Information Processing Systems*, 37:131197–131223.
- Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. 2025. Wise: rethinking the knowledge memory for lifelong model editing of large language models. In *Proceedings of the 38th International Conference on Neural Information Processing Systems, NIPS '24*, Red Hook, NY, USA. Curran Associates Inc.
- Yifan Wang, Yafei Liu, Chufan Shi, Haoling Li, Chen Chen, Haonan Lu, and Yujia 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.
- Alexander Wettig, Aatmik Gupta, Saumya Malik, and Danqi Chen. 2024. Qurating: selecting high-quality data for training language models. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org.
- Heesoo Won, Byungkook Oh, Hyeongjun Yang, and Kyong-Ho Lee. 2023. Cross-modal contrastive learning for aspect-based recommendation. *Information Fusion*, 99:101858.
- Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari. 2024. [Continual learning for large language models: A survey](#). Preprint, arXiv:2402.01364.
- Yunjia Xi, Weiwen Liu, Jianghao Lin, Bo Chen, Ruiming Tang, Weinan Zhang, and Yong Yu. 2024. [Memocrs: Memory-enhanced sequential conversational recommender systems with large language models](#). In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM '24*, page 2585–2595, New York, NY, USA. Association for Computing Machinery.
- Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. Less: selecting influential data for targeted instruction tuning. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org.
- Hyeongjun Yang, Donghyun Kim, Gayeon Park, KyuHwan Yeom, and Kyong-Ho Lee. 2025. [Coresense: Social commonsense knowledge-aware context refinement for conversational recommender system](#). *IEEE Transactions on Knowledge and Data Engineering*, 37(4):1702–1713.
- Hyeongjun Yang, Yerim Lee, Gayeon Park, Taeyoung Kim, Heesun Kim, Kyong-Ho Lee, and Byungkook Oh. 2024. Granular intents learning via mutual information maximization for knowledge-aware recommendation. *Knowledge-Based Systems*, 306:112705.
- Hyeongjun Yang, Heesoo Won, Youbin Ahn, and Kyong-Ho Lee. 2023. Click: Contrastive learning for injecting contextual knowledge to conversational recommender system. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1867–1877.
- Ting Yang and Li Chen. 2024. [Unleashing the retrieval potential of large language models in conversational recommender systems](#). In *Proceedings of the 18th ACM Conference on Recommender Systems, RecSys '24*, page 43–52, New York, NY, USA. Association for Computing Machinery.
- Bowen Zhao, Zander Brumbaugh, Yizhong Wang, Hananeh Hajishirzi, and Noah Smith. 2024. [Set the clock: Temporal alignment of pretrained language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 15015–15040, Bangkok, Thailand. Association for Computational Linguistics.
- Haizhong Zheng, Rui Liu, Fan Lai, and Atul Prakash. 2023. [Coverage-centric coreset selection for high pruning rates](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*.
- Junhao Zheng, Shengjie Qiu, Chengming Shi, and Qianli Ma. 2025. [Towards lifelong learning of large language models: A survey](#). *ACM Comput. Surv.*, 57(8).
- Kun Zhou, Wayne Xin Zhao, Shuqing Bian, Yuanhang Zhou, Ji-Rong Wen, and Jingsong Yu. 2020. [Improving conversational recommender systems via knowledge graph based semantic fusion](#). In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1006–1014.
- Yuanhang Zhou, Kun Zhou, Wayne Xin Zhao, Cheng Wang, Peng Jiang, and He Hu. 2022. [C²-crs: Coarse-to-fine contrastive learning for conversational recommender system](#). In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, pages 1488–1496.

A Detailed Derivation of RHO Importance Score

The value of a candidate sample (x, y) can be quantified by how much fine-tuning a model pre-trained on D_t with (x, y) reduces the loss on the validation set D_{val} . This process is formalized as selecting the sample that solves the following optimization problem:

$$\arg \min_{(x,y) \in D_{\text{cand}}} \sum_{(x',y') \in D_{\text{val}}} -\log p(y' | x'; D_t \cup \{(x, y)\}). \quad (7)$$

However, this calculation is prohibitively expensive, as it requires re-training the model with each individual sample and re-evaluating the loss on the validation set repeatedly. Thus, the above objective can be approximated by using Bayes' rule and conditional independence for a more cost-efficient estimation of the sample. The derivation proceeds as follows:

$$\begin{aligned} & \log p(y' | x'; D_t \cup \{(x, y)\}) \\ &= \log \frac{p(y | x; x', y', D_t) p(y' | x', x; D_t)}{p(y | x, x'; D_t)} \quad \text{Bayes rule} \\ &= \log \frac{p(y | x; y', x', D_t) p(y' | x'; D_t)}{p(y | x; D_t)} \quad \text{conditional independence} \\ &\propto L(y | x; \theta_{D_t}) - L(y | x; \theta_{D_t \cup D_{\text{val}}}). \end{aligned}$$

This approximation defines the importance score of the sample (x, y) as the difference between the two loss values. $L(y | x; \theta_{D_t})$ is the training loss of the sample (x, y) on the current model pre-trained on D_t , reflecting how surprising the sample is to the model. $L(y | x; \theta_{D_t \cup D_{\text{val}}})$ is a proxy loss computed using a model that further trained on the validation set D_{val} . This term estimates the sample's inherent learnability and helps filter out outliers.

B Pilot Study: Vulnerability of Existing Data Selections to Familiar Samples and Outliers

To empirically illustrate the failure mode of existing data selections, we construct a controlled pilot study using the LLM-based CRS (He et al., 2025) on a Reddit conversational recommendation dataset (He et al., 2023). We consider two consecutive datasets, D_t and D_{t+1} , where D_{t+1} represents the shifted environment. We curate two groups within D_{t+1} :

- **Familiar.** We sample 300 instances from D_{t+1} that share the same recommended items with instances in D_t , while also preserving similar user intents and conversational contexts. These samples represent already-learned and context-consistent patterns that are less informative for shift adaptation.
- **Outlier.** From D_{t+1} , we construct 300 perturbed instances by replacing the originally relevant recommended item with an intentionally irrelevant one that is semantically and entity-wise disjoint from the original preference. For example, a request aligned with U.S. blockbuster Marvel movies is paired with a low-budget independent film that critiques social issues of South Korea, ensuring no overlap in production company, director, or cast. These samples mimic contextual mismatch that yields high losses but minimal useful learnability.

C Impact of Temporal Adaptation Granularity on Training Loss Distribution

To analyze how the training loss distribution varies with temporal granularity, we measure the distribution of the training loss (\mathcal{L}_t) using the LLM-based CRS (He et al., 2025) on the Reddit dataset (He et al., 2023). We partitioned the Reddit dataset into several temporal granularities, specifically using 3, 9, and 25 splits. For each setting, the initial split was utilized for pre-training. The loss distributions are then evaluated over the remaining 2, 8, and 24 phases, respectively. Figure 8 presents the resulting loss distributions using boxplots across the three temporal split settings. Our observations are as follows:

- **Coarse-grained Adaptation (2 Split).** In scenarios where updates are infrequent, the model encounters substantial contextual shifts within a single phase. Consequently, the loss distribution covers an extensive range, with the upper whisker extending significantly to reach 96.8. At the same time, a high density of low-loss samples (*i.e.*, familiar samples) is observed, indicating that many items from the previous context remain valid or relevant despite the large contextual transition.
- **Intermediate Adaptation (8 Split).** As the update frequency increases, the spread of the

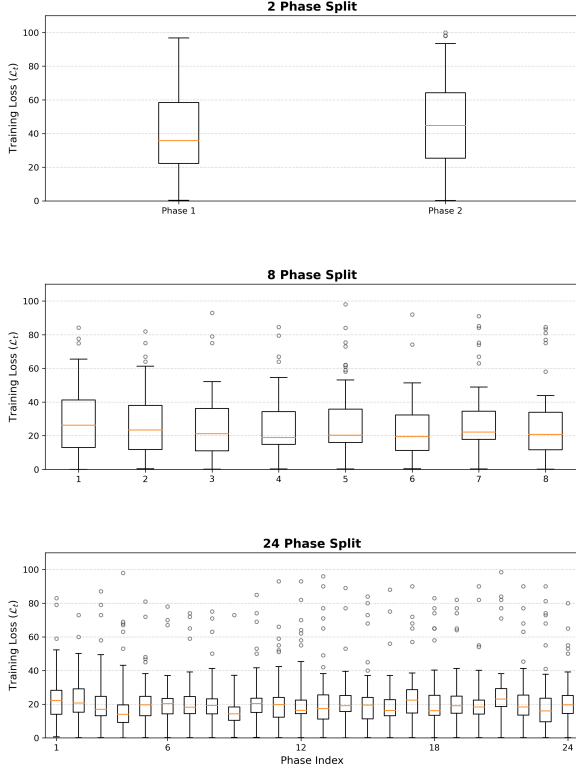


Figure 8: Distribution of training loss (\mathcal{L}_t) for the current model π_θ on the subsequent phase data D_{t+1} across different temporal split settings (2, 8, and 24).

loss distribution becomes narrower compared to the 2-phase setting, suggesting reduced loss-scale disparities under more frequent adaptation.

- **Fine-grained Adaptation (24 Split).** Under the most frequent update setting, the loss distribution exhibits the narrowest range. Similar to the other settings, a high volume of samples maintains low loss values, further highlighting the persistent presence of familiar samples.

Implications. Our analysis reveals that large contextual transitions, such as those in the 2-phase setup, induce severe scale volatility in training loss. In such environments, the absolute magnitude of the loss value ceases to be a reliable indicator of a sample’s actual learnability. Meanwhile, The consistent presence of high-density low-loss samples across all granularities confirms that loss-based approaches are inherently more suitable than gradient-based methods for CRS tasks, as they can more efficiently deprioritize these non-informative samples. The nature of CRS task dictates that while some optimal recommendations change due to evolving trends, many others remain effective over time.

D Impact of KL-Divergence Regularization for Shift Adaptation in LLM-based CRSs

In the context of continual learning, regularization techniques are often employed to prevent the model from deviating excessively from its previous state (*i.e.*, catastrophic forgetting). To investigate the necessity of such constraints in our setting, we experimented with a KL-divergence regularization term:

$$\mathcal{L}_{reg} = \lambda \cdot KL(p_{\theta_{t+1}}(y|x) || p_{\theta_t}(y|x)), \quad (8)$$

where θ_t represents the parameters of the previous phase, and θ_{t+1} is the updated parameters for shift adaptation. However, our empirical results indicated that this regularization did not yield performance gains. Instead, adding regularization terms increases computational complexity and requires additional training iterations. We assumed that since our proposed method selects high-quality shift-adaptive samples (D_p), the optimization landscape is naturally guided toward meaningful updates without requiring explicit constraints. Furthermore, consistent with standard LLM-based CRS practices, limiting the number of training iterations proved to be a sufficient and more efficient strategy for maintaining stability.

E Implementation Details

CAPT. To ensure numerical stability (*e.g.*, preventing division-by-zero or numerical explosion) in Equation 4 of the importance score estimation, we employ a stability constant $\delta = 10^{-2}$. The importance scores are calculated as $\mathcal{I}_{(x,y)} = \mathcal{L}_t / (\mathcal{L}_p + \delta)$. The overall procedure for the CAPT adaptation across temporal phases is detailed in Algorithm 1.

CRS Models. The original implementations of PECRS and RefiCR without modification are used. For RTA, we adopted the ‘w/ Bias’ variant and fully trained both the additional aggregator and the LLM. To ensure a fair comparison, we adhered to the default hyperparameters such as batch size, learning rate, and epochs as specified in the original implementation of each CRS.

Evaluation Settings. Reported results are averaged over 10 runs with random seeds, and the performance gains of CAPT are statistically significant relative to the strongest baseline, as verified by a t-test ($p < 0.05$). Figure 7 provides a sensitivity

Algorithm 1: The CAPT Framework for Adapting to Temporal Phase $t + 1$

Input: Training dataset D_{t+1} ; validation dataset D_{val} at phase $t+1$; initial target model π_θ pre-trained on D_t ; selection budget k ; hard negative sampling ratio τ ; stabilization constant δ

Output: Adapted target model π_θ

- 1 **Step 1: Proxy Model Preparation**
 - 2 Initialize proxy model $\pi_\phi \leftarrow \pi_\theta$
 - 3 Fine-tune π_ϕ on D_{val}
 - 4 **Step 2: Importance Score Estimation**
 - 5 **for** each training sample $(x, y) \in D_{t+1}$ **do**
 - 6 Compute proxy loss $\mathcal{L}_p(x, y)$ using π_ϕ
 - 7 Compute training loss $\mathcal{L}_t(x, y)$ using π_θ
 - 8 Calculate importance score:
$$\mathcal{I}(x, y) = \frac{\mathcal{L}_t(x, y)}{\mathcal{L}_p(x, y) + \delta}$$
 - 9 Select top- k samples from D_{t+1} based on $\mathcal{I}(x, y)$ to form a subset D_p
 - 10 **Step 3: Importance-Weighted Tuning**
 - 11 Sample hard negatives from high-importance samples with ratio τ
 - 12 **for** each mini-batch from D_p **do**
 - 13 Update π_θ using the weighted loss function: $\mathcal{L}_{\text{total}} = \sum \mathcal{I}(x, y) \cdot \mathcal{L}(x, y)$
 - 14 **return** π_θ
-

analysis of the negative sampling ratio τ ranges from 0% to 100%. We reported the experimental results in Figures 4-6 using a sampling ratio of $\tau = 20\%$. We conducted experiments on a NVIDIA RTX A6000 GPU (48GB memory) The computation time reported in Table 3 is measured on the RTA model with LoRA using the Reddit Phase-2 dataset with a batch size of 16. To evaluate compatibility with PEFT in Figure 6, we applied LoRA with rank $r = 16$.

F CRS Datasets Details

For the Reddit dataset, each split consists of training, validation, and test sets in an 8:1:1 ratio. The REDIAL (Li et al., 2018) and INSPIRED (Hayati et al., 2020) datasets are both constructed using the Amazon Mechanical Turk (AMT) platform, where workers role-played as either users seeking movie recommendations or recommenders providing them. REDIAL consists of 182,150 utterances, while INSPIRED contains 35,811 utterances. Both datasets were split 8:1:1 into training, validation,

	Method	Recall@10		MRR@10		NDCG@10	
		$\geq Full$	Best	$\geq Full$	Best	$\geq Full$	Best
2 Splits	Full	–	0.2969	–	0.1287	–	0.1476
	GREATS	80%	0.2971	70%	0.1287	70%	0.1479
	RHO	100%	0.2969	50%	0.1287	60%	0.1476
	Ours	40%	0.3037	40%	0.1339	30%	0.1541
24 Splits	Full	–	0.3003	–	0.1351	–	0.1463
	GREATS	67.3%	0.3011	59.2%	0.1361	53.9%	0.1465
	RHO	69.4%	0.3018	47.6%	0.1358	56.5%	0.1471
	Ours	33.7%	0.3053	39.4%	0.1402	28.3%	0.1510

Table 4: Performance comparison on the Reddit_{Phase1} dataset. Results for 24 Splits are averaged across all 24 sequential adaptation phases. ‘ $\geq Full$ ’ indicates the percentage of data required to match/exceed the performance of full dataset. ‘Best’ denotes the best performance of each method.

and test sets. In summary, Reddit is used to evaluate ‘temporal shifts adaptation’, while REDIAL and INSPIRED evaluate ‘task-level adaptation’.

G Additional Experimental Results

Evaluation on @50 Metrics. While Figures 4 and 5 exclude the @50 results due to space constraints, Figures 9 and 10 respectively report performance on additional evaluation metrics.

Robustness to Temporal Granularity (24-Phase Adaptation). To evaluate the consistency of CAPT in more frequent update scenarios, we conduct experiments on the Reddit dataset partitioned into 24 temporal phases. This setup simulates a real-world environment where the CRS is updated at a higher frequency (e.g., monthly). Table 4 reports the average performance across all 24 adaptation phases. We observe that the performance gap between CAPT and other methods is narrower in this 24-phase setting than in the 2-phase experiments. This is because more frequent updates naturally reduce the magnitude of contextual shift per phase, making the data selection task less challenging as most samples remain relatively “familiar” to the model. Nevertheless, CAPT still achieves the highest accuracy demonstrates its consistent ability to identify learnable samples regardless of the update frequency.

Replay Ratio	Random	GREATS	CAPT
0% (Phase 1 only)	0.296	0.296	0.296
10%	0.281	0.287	0.264
20%	0.264	0.258	0.243

Table 5: Recall@10 of the RTA model on Reddit_{Phase1} when replaying past data from Reddit_{Pretrain}.

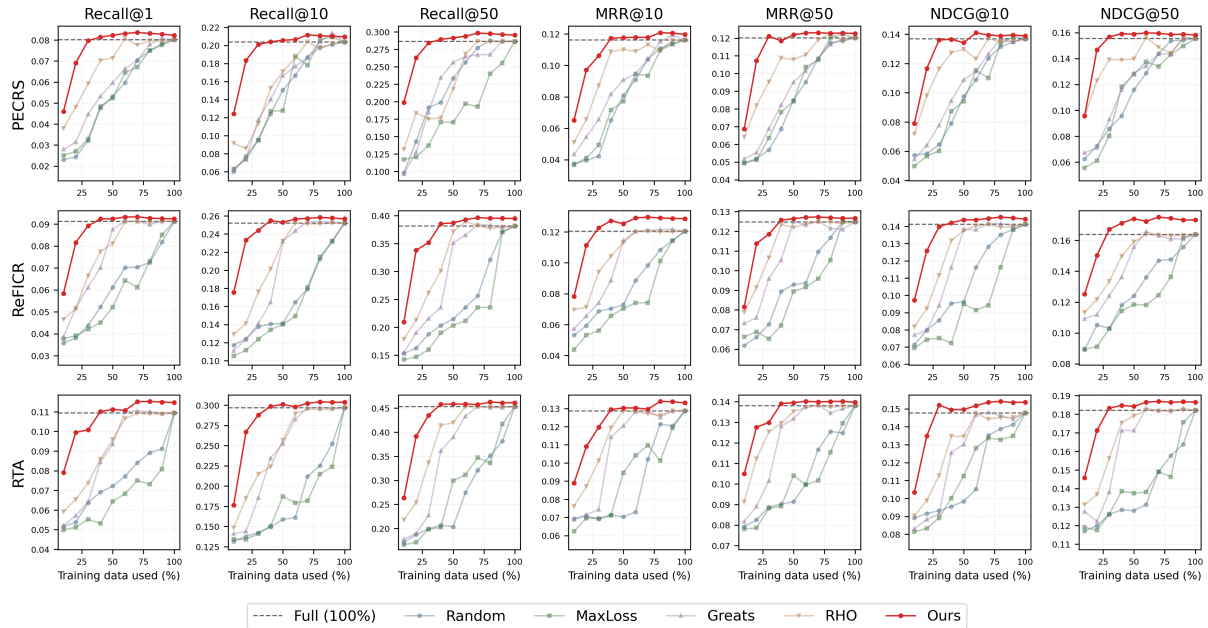


Figure 9: Additional evaluation metrics for data selection methods on the $\text{Reddit}_{\text{Phase1}}$ dataset across various CRS models.

H Replaying Past Data for CRS Adaptation

As discussed in the ‘Adapting Language Models to Shifts’ subsection of the Related Work, we further examine whether replaying past data (*i.e.*, Experience Replay) benefits temporal adaptation in CRSs, which is a common technique in continual learning. Initially, we trained the RTA model on the full combined datasets of $\text{Reddit}_{\text{Pretrain}}$ and $\text{Reddit}_{\text{Phase1}}$ and evaluated its performance on $\text{Reddit}_{\text{Phase1}}$. Contrary to the typical success of Experience Replay in other domains, we observed a significant performance degradation: the Recall@10 score on $\text{Reddit}_{\text{Phase1}}$ dropped from 0.296 to 0.211.

To explore if intelligent selection could mitigate this, we merged $\text{Reddit}_{\text{Phase1}}$ with subsets of $\text{Reddit}_{\text{Pretrain}}$ sampled using different strategies (Random, GREATS, and CAPT) at varying ratios (10% and 20%). The results are summarized in Table 5. Across all baselines, replaying past data consistently led to inferior performance compared to training on the most recent data alone. We attribute this to the shift in global preferences, where historical data conflict with current contexts. Consequently, these results support our emphasis on data pruning of the current phase rather than replaying past phases for efficient CRS adaptation.

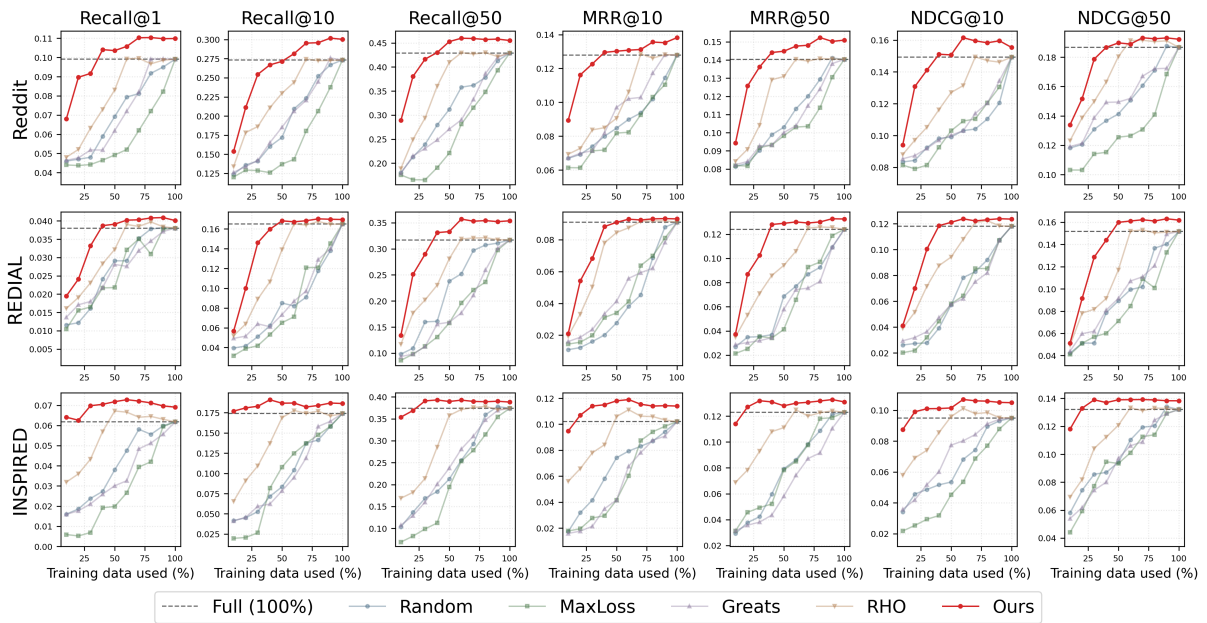


Figure 10: Additional evaluation metrics for the RTA model across different datasets. Results are reported on $\text{Reddit}_{\text{Phase2}}$ (temporal shift adaptation) and REDIAL / INSPIRED (CRS task adaptation).