

# Multi-Stage LLM Fine-Tuning with a Continual Learning Setting

Changhao Guan<sup>1</sup>, Chao Huang<sup>1</sup>, Hongliang Li<sup>1</sup>

, You Li<sup>1</sup>, Ning Cheng<sup>1</sup>, Ziheng Liu<sup>1</sup>

, Jinan Xu<sup>1</sup>, Yufeng Chen<sup>\*1</sup>, Jian Liu<sup>2</sup>

<sup>1</sup>Beijing Jiaotong University, Beijing, China

<sup>2</sup>University of Science and Technology Beijing, Beijing, China

{guanchanghao, huangchao, hongliangli, youlee}@bjtu.edu.cn

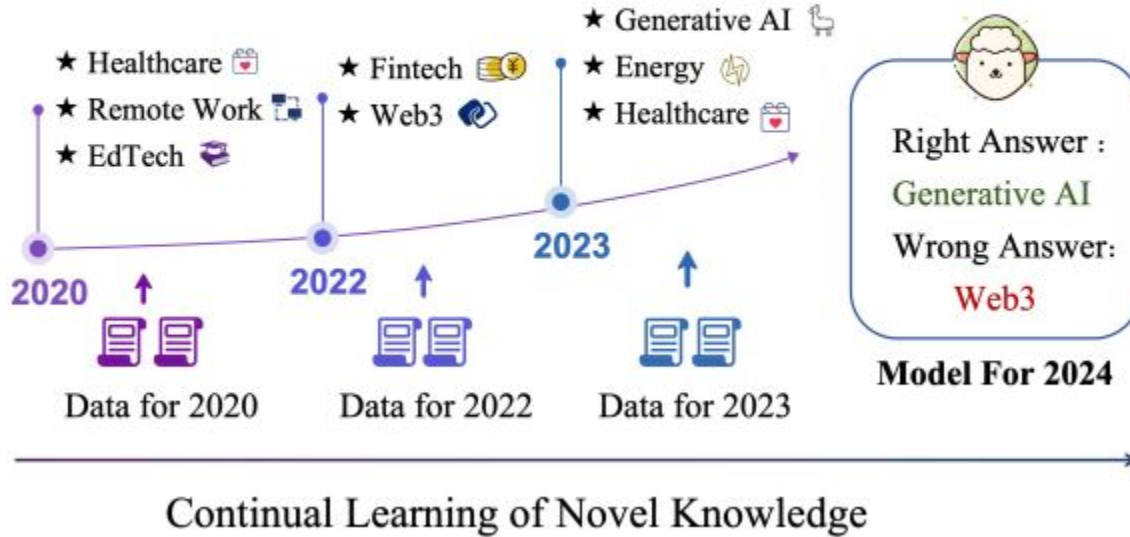
{ningcheng, 23120386, jaxu, chenylf}@bjtu.edu.cn, [jian.liu@ustb.edu.cn](mailto:jian.liu@ustb.edu.cn)

Year of Publication :- 2025

Number of Citations :- 2

# Introduction

**What are the most prominent technology sectors for global venture capital investment in today's society?**



# Introduction

- Large Language Models (LLMs) are considered complex knowledge repositories as they have ability to represent diverse general information [1,2].
- However, it is necessary to fine-tune them on customized datasets when applying them to specific domains [3,4].
- In addition, there is a continual learning requirement, especially when domain knowledge rapidly changes[5,6].
- Pilot experiments show that employing the standard fine-tuning methods for LLMs significantly degrades their performance. However, this problem has not received much research attention.

[1] Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., ... & Agarwal, S. (2020). Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 1(3), 3.

[2] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35, 27730-27744.

[3] Xu, R., Luo, F., Zhang, Z., Tan, C., Chang, B., Huang, S., & Huang, F. (2021). Raise a child in large language model: Towards effective and generalizable fine-tuning. arXiv preprint arXiv:2109.05687.

[4] Xie, T., Wan, Y., Huang, W., Yin, Z., Liu, Y., Wang, S., ... & Hoex, B. (2023). Darwin series: Domain specific large language models for natural science. arXiv preprint arXiv:2308.13565.

[5] McCann, B., Keskar, N. S., Xiong, C., & Socher, R. (2018). The natural language decathlon: Multitask learning as question answering. arXiv preprint arXiv:1806.08730.

[6] Gururangan, S., Marasović, A., Swayamdipta, S., Lo, K., Beltagy, I., Downey, D., & Smith, N. A. (2020). Don't stop pretraining: Adapt language models to domains and tasks. arXiv preprint arXiv:2004.10964.

# Introduction - Problems in Continual Learning

- Potential knowledge conflict [7,8]. When a domain undergoes rapid changes, potential conflicts between new and old knowledge may arise, potentially leading to “hallucinations” in LLMs.
- Incomparable amount of fine-tuning data compared to pre-training data [9,10]. Compared to the extensive data leveraged during pre-training, the domain specific data available for fine-tuning is typically scarce, making it challenging to adapt the model’s parameters to fit for fine-tuning.
- Proposed a new approach for fine-tuning LLMs in the multi-stage continual learning settings using a preference-based forgetting strategy and self-distillation based data augmentation.

[7] Longpre, S., Perisetla, K., Chen, A., Ramesh, N., DuBois, C., & Singh, S. (2021). Entity-based knowledge conflicts in question answering. arXiv preprint arXiv:2109.05052.

[8] Liu, Y., Yao, Z., Lv, X., Fan, Y., Cao, S., Yu, J., ... & Li, J. (2024). Untangle the knot: Interweaving conflicting knowledge and reasoning skills in large language models. arXiv preprint arXiv:2404.03577.

[9] Jiang, G., Jiang, C., Xue, S., Zhang, J. Y., Zhou, J., Lian, D., & Wei, Y. (2023). Towards anytime fine-tuning: Continually pre-trained language models with hypernetwork prompt. arXiv preprint arXiv:2310.13024.

[10] Dong, G., Yuan, H., Lu, K., Li, C., Xue, M., Liu, D., ... & Zhou, J. (2023). How abilities in large language models are affected by supervised fine-tuning data composition. arXiv preprint arXiv:2310.05492.

# Related Works - Fine tuning LLMs

- Fine-tuning is a widely adopted approach to adopt LLMs to new domains and tasks using domain specific data [11,12].
  - Fine-tuning for complex instructions.
  - Fine-tuning for specific domains.
- Solely relying on fine-tuning often struggles to acquire new knowledge when facing significant domain shifts[13].
- Existing works have 2 stage approach [14].
  - First Stage - Fine-tuning to acquire domain knowledge.
  - Second Stage - Fine-tuning to enhance task specific capabilities.

[11] Ding, R., Han, X., & Wang, L. (2022). A unified knowledge graph augmentation service for boosting domain-specific NLP tasks. arXiv preprint arXiv:2212.05251.

[12] Zheng, J., Hong, H., Liu, F., Wang, X., Su, J., Liang, Y., & Wu, S. (2024). Fine-tuning large language models for domain-specific machine translation. arXiv preprint arXiv:2402.15061.

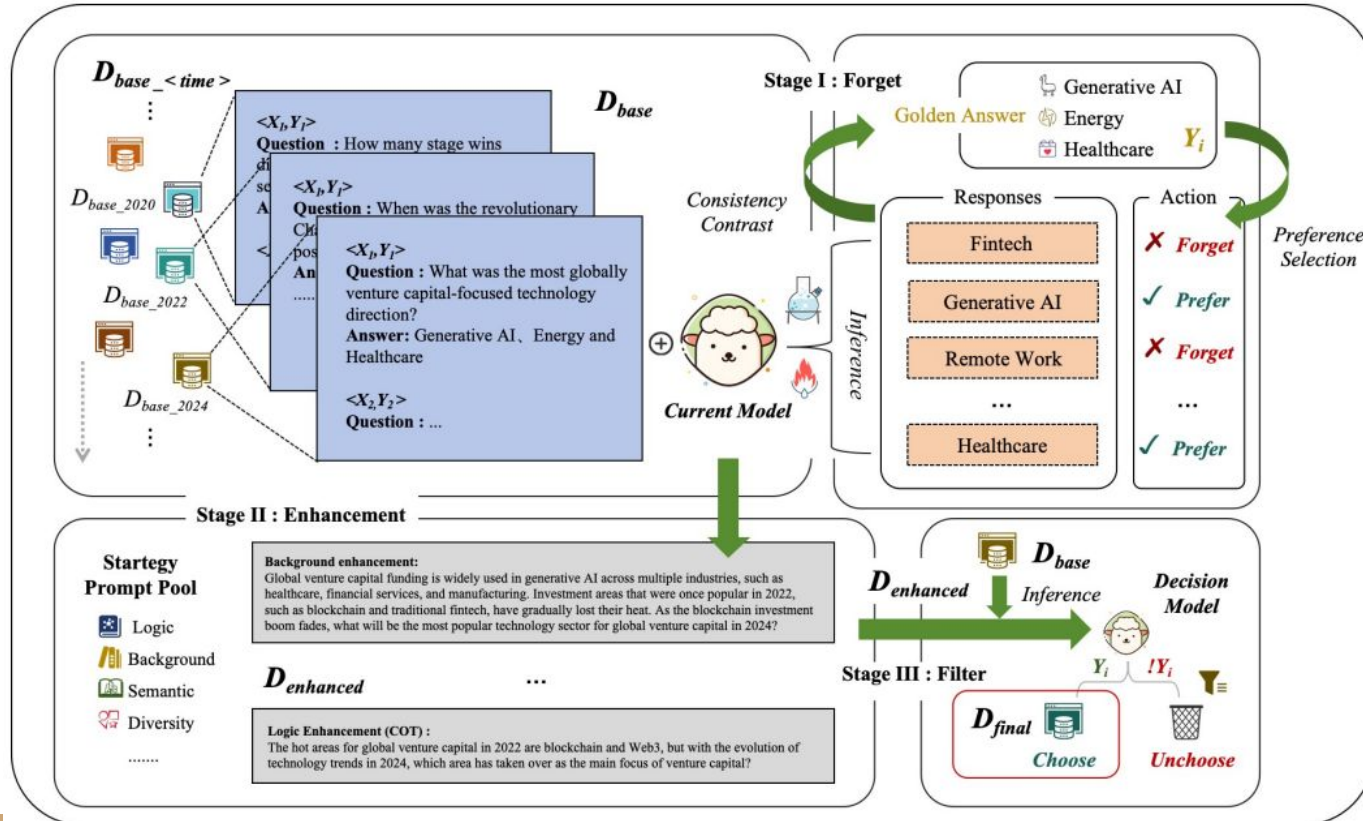
[13] Emelin, D., Bonadiman, D., Alqahtani, S., Zhang, Y., & Mansour, S. (2022). Injecting domain knowledge in language models for task-oriented dialogue systems. arXiv preprint arXiv:2212.08120.

[14] Han, R., Ren, X., & Peng, N. (2020). ECONET: Effective continual pretraining of language models for event temporal reasoning. arXiv preprint arXiv:2012.15283.

# Related Works - Continual Learning with LLMs

- Traditional CL methods [15].
  - Regularization based.
  - Replay Based.
  - Architecture based strategies.
- Continuous fine tuning strategy [16] .
- Modular continual learning [17].
- Forget-before-learn [18].

# Approach - Overview



# Preference Based Learning Bias

- Let  $(x, y)$  be a training example.
- $x$  is utilized as an input and apply the model  $K$  times to get a prediction set  $Y$  for each input  $x$ .
- Measure compatibility between each element of the set  $Y$  and the desired  $y$ .
- Divide  $Y$  into subsets  $Y_{\text{align}}$  and  $Y_{\text{conflict}}$ .
- Main motivation is to bias the model to generate responses similar to those in  $Y_{\text{align}}$  and avoid those in  $Y_{\text{conf}}$ .



# Preference Losses

- Positive preference loss.

$$\mathcal{L}_{\text{PP}} = - \sum_{y' \in Y_{\text{align}}} \log \left( \frac{\pi_{\theta}(y' | x)}{\pi_{\text{ref}}(y' | x)} \right)$$

- Negative preference loss.

$$\mathcal{L}_{\text{NP}} = \sum_{y' \in Y_{\text{conf}}} \log \left( \frac{\pi_{\theta}(y' | x)}{\pi_{\text{ref}}(y' | x)} \right)$$

- Total loss

$$\mathcal{L}_{\text{total}} = \alpha \cdot \mathcal{L}_{\text{PP}} + \beta \cdot \mathcal{L}_{\text{NP}}$$

# Data Augmentation with Self Distillation

- Using the LLMs themselves for augmentation.
- Augmentation strategies.
  - Background knowledge integration - LLM is asked to provide more background knowledge in order for the input to contain more context related information.
  - Logic-Compatible Expansion - LLM is asked to incorporate the logic-related information to expand the semantic complexity of the input.
  - Paraphrase augmentation - This method involves rewriting and rephrasing the original example to write more similar examples with various structures and expressions.
- Based on the above strategies, a new pair  $(x', y)$  can be created for any training example  $(x, y)$ .

# Dynamic Data Selection Strategy

- To evaluate the effectiveness and validity if the augmented data for model training.
- Based on heuristic criteria.

- Mutual Information.

$$MI(x'; \hat{x}') = \sum_{w, \hat{w}'} N(w, \hat{w}') \log \left( \frac{N(w, \hat{w}')}{N(w)N(\hat{w}')} \right)$$

- Indication from LLMs.

By this criterion, it is measured whether  $x'$  produces same result as  $x$ .

- The filtered set of high-quality samples will be used as input for subsequent augmentation and fine-tuning processes.

# Experimental Setup - Datasets

- Considered rapidly evolving domains: Natural sciences, medicine, technology, transportation, tourism, finance and social sciences.
- Total 21,000 question- answer pairs with 3000 question-answer pairs for each domain.
- In this, there are 1000 examples shared by any 2 domains to evaluate cross domain conflict situations.
- 6000 additional samples as general purpose samples to indicate model's performance in domain agnostic setting.

# Experimental Setup - Evaluation Settings

- Domain independent continual learning.
  - Use the domain independent dataset (6000 samples) for fine-tuning in the initial stage.
  - Manually edit the answers and use the revised dataset for fine-tuning.
  - Generate the same number of examples compatible with the fine-tuning data as evaluation set.
- Cross domain scenarios.
  - Conducting continual learning using cross domain data by gradually adding domains one by one.
  - Manual verification to ensure that a minimum 1000 examples are shared between 2 domains.
  - Before fine-tuning, answers were edited to be different from previous domain to mimic domain dispute.
  - Each fine-tuning stage contains, 2000 domain independent examples and 1000 cross domain conflict examples.

# Experimental Setup - Evaluation Metric

- Knowledge Gain Ratio (KGR) - Assessing model's improvement in learning dynamically evolving knowledge.
- Post-injection Accuracy - To measure overall accuracy improvement in a given test set.

# Experimental Setup - Baselines and backbone LLMs

- Baselines

- Continual instruction fine-tuning (CIF).
- Modular continual learning (MoCL).
- Forgetting before learning (F-Learning).

- Backbone LLMs.

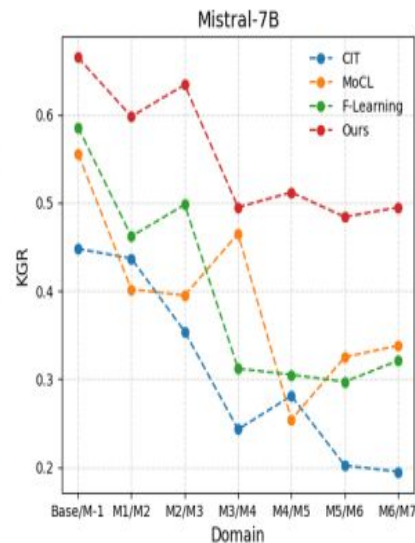
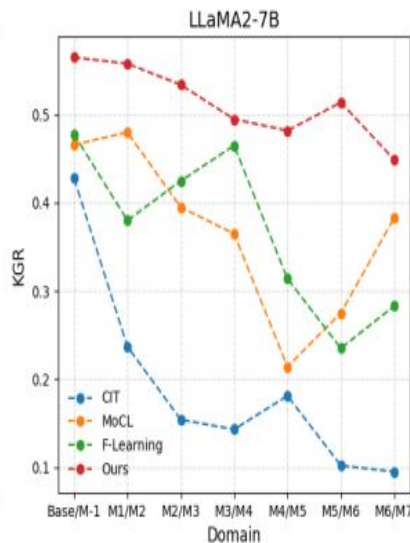
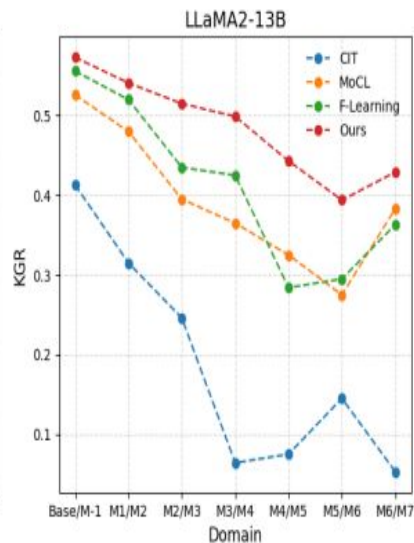
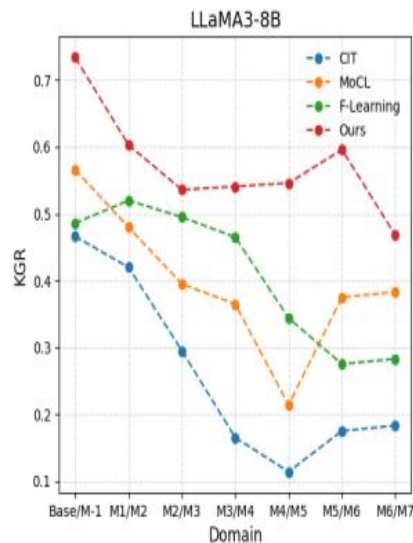
- Llama2.
- Llama3.
- Misral7B.

# Results - Domain independent continual learning

Eval	Stage1 (Initial)		Stage2		Stage3		Stage4		Stage5	
	KGR	ACC	KGR	ACC	KGR	ACC	KGR	ACC	KGR	ACC
<b>Llama2-7B</b>										
<b>CIT (2023b)</b>	50.50	66.55	44.30 <sub>↓6.20</sub>	46.10 <sub>↓20.45</sub>	12.50 <sub>↓38.00</sub>	25.39 <sub>↓41.16</sub>	12.20 <sub>↓38.30</sub>	24.23 <sub>↓42.32</sub>	11.70 <sub>↓38.80</sub>	17.60 <sub>↓48.95</sub>
<b>MoCL (2024a)</b>	—	—	48.20 <sub>↑3.90</sub>	49.60 <sub>↑3.50</sub>	45.50 <sub>↑33.00</sub>	47.70 <sub>↑22.31</sub>	46.80 <sub>↑34.60</sub>	47.10 <sub>↑22.87</sub>	26.20 <sub>↑14.50</sub>	27.35 <sub>↑9.75</sub>
<b>F-Learning (2023)</b>	—	—	54.40 <sub>↑10.10</sub>	57.50 <sub>↑11.40</sub>	49.40 <sub>↑36.90</sub>	52.25 <sub>↑26.86</sub>	49.20 <sub>↑37.00</sub>	49.75 <sub>↑25.52</sub>	21.80 <sub>↑10.10</sub>	20.90 <sub>↑3.30</sub>
<b>Ours</b>	—	—	60.20 <sub>↑15.90</sub>	62.55 <sub>↑16.45</sub>	69.40 <sub>↑56.90</sub>	69.65 <sub>↑44.26</sub>	68.60 <sub>↑56.40</sub>	67.53 <sub>↑43.30</sub>	75.80 <sub>↑64.10</sub>	75.49 <sub>↑57.89</sub>
<b>Llama2-13B</b>										
<b>CIT (2023b)</b>	68.90	81.95	32.20 <sub>↓36.70</sub>	34.20 <sub>↓47.75</sub>	26.80 <sub>↓42.10</sub>	33.55 <sub>↓48.40</sub>	24.60 <sub>↓44.30</sub>	32.85 <sub>↓49.10</sub>	11.20 <sub>↓57.70</sub>	17.85 <sub>↓64.10</sub>
<b>MoCL (2024a)</b>	—	—	41.60 <sub>↑9.40</sub>	41.95 <sub>↑7.75</sub>	50.40 <sub>↑23.60</sub>	51.25 <sub>↑17.70</sub>	48.70 <sub>↑24.10</sub>	50.00 <sub>↑17.15</sub>	25.80 <sub>↑14.60</sub>	26.50 <sub>↑8.65</sub>
<b>F-Learning (2023)</b>	—	—	43.30 <sub>↑11.10</sub>	43.80 <sub>↑9.60</sub>	59.30 <sub>↑32.50</sub>	60.70 <sub>↑27.15</sub>	53.90 <sub>↑29.30</sub>	54.90 <sub>↑22.05</sub>	33.60 <sub>↑22.40</sub>	33.90 <sub>↑16.05</sub>
<b>Ours</b>	—	—	66.30 <sub>↑34.10</sub>	67.25 <sub>↑33.05</sub>	76.50 <sub>↑49.70</sub>	77.40 <sub>↑43.85</sub>	66.20 <sub>↑41.60</sub>	65.45 <sub>↑32.60</sub>	76.50 <sub>↑65.30</sub>	77.65 <sub>↑59.80</sub>
<b>Llama3-8B</b>										
<b>CIT (2023b)</b>	65.90	81.55	48.80 <sub>↓17.10</sub>	49.90 <sub>↓31.65</sub>	31.90 <sub>↓34.00</sub>	34.05 <sub>↓47.50</sub>	30.30 <sub>↓35.60</sub>	35.30 <sub>↓46.25</sub>	23.40 <sub>↓42.50</sub>	27.70 <sub>↓53.85</sub>
<b>MoCL (2024a)</b>	—	—	70.40 <sub>↑21.6</sub>	70.65 <sub>↑20.75</sub>	52.50 <sub>↑20.60</sub>	52.65 <sub>↑18.60</sub>	63.30 <sub>↑33.00</sub>	63.65 <sub>↑28.35</sub>	31.30 <sub>↑7.90</sub>	30.95 <sub>↑3.25</sub>
<b>F-Learning (2023)</b>	—	—	67.00 <sub>↑18.20</sub>	67.45 <sub>↑17.55</sub>	58.40 <sub>↑26.50</sub>	57.95 <sub>↑23.90</sub>	61.40 <sub>↑31.10</sub>	61.20 <sub>↑25.90</sub>	57.70 <sub>↑34.30</sub>	57.90 <sub>↑30.20</sub>
<b>Ours</b>	—	—	83.60 <sub>↑34.80</sub>	83.90 <sub>↑34.00</sub>	69.40 <sub>↑37.50</sub>	69.55 <sub>↑35.50</sub>	69.20 <sub>↑33.90</sub>	69.20 <sub>↑39.03</sub>	74.80 <sub>↑51.40</sub>	74.60 <sub>↑46.90</sub>
<b>Mistral-7B</b>										
<b>CIT (2023b))</b>	61.00	74.20	41.20 <sub>↓19.80</sub>	44.65 <sub>↓29.55</sub>	38.50 <sub>↓22.50</sub>	38.90 <sub>↓35.30</sub>	32.80 <sub>↓28.20</sub>	36.40 <sub>↓37.80</sub>	21.60 <sub>↓39.40</sub>	23.50 <sub>↓50.70</sub>
<b>MoCL (2024a)</b>	—	—	54.40 <sub>↑13.20</sub>	51.25 <sub>↑6.60</sub>	62.50 <sub>↑24.00</sub>	59.10 <sub>↑20.20</sub>	58.79 <sub>↑25.99</sub>	57.25 <sub>↑20.85</sub>	47.10 <sub>↑25.50</sub>	46.00 <sub>↑22.50</sub>
<b>F-Learning (2023)</b>	—	—	48.30 <sub>↑7.10</sub>	51.50 <sub>↑6.85</sub>	58.10 <sub>↑19.60</sub>	57.65 <sub>↑18.75</sub>	54.20 <sub>↑21.40</sub>	55.30 <sub>↑18.90</sub>	33.80 <sub>↑12.20</sub>	34.90 <sub>↑11.40</sub>
<b>Ours</b>	—	—	64.60 <sub>↑23.40</sub>	65.72 <sub>↑21.07</sub>	72.60 <sub>↑34.10</sub>	72.80 <sub>↑33.90</sub>	71.20 <sub>↑38.40</sub>	71.58 <sub>↑35.18</sub>	77.50 <sub>↑55.90</sub>	77.34 <sub>↑53.84</sub>



# Results - Cross Domain Setting



# Discussion - Ablation study

- Impact of Preference based Learning Bias.

Method	ACC(%)	KGR(%)
CIT	48.80	49.90
<b>PBL (Ours)</b>	<b>59.00</b>	<b>59.30</b>

- Impact of Data Augmentation with Self-Distillation.

Method	ACC (%)	KGR (%)
No Argument	49.90	48.80
+ BKI	<u>58.85</u>	58.60
+ LCE	52.80	52.50
+ PA	56.40	<u>64.10</u>
<b>CA (Ours)</b>	<b>69.90</b>	<b>69.30</b>

# Discussion - Ablation Study

- Impact of Dynamic Data Selection Strategy.

Method	ACC (%)	KGR (%)
No Argument	49.90	48.80
Data Argument	69.90	69.30
+ RS (50%)	<u>79.85</u>	<u>79.60</u>
+ RS (25%)	69.85	63.50
+ RS (12.5%)	70.40	66.70
<b>DS (Ours)</b>	<b>81.20</b>	<b>82.50</b>

# Discussion - Analysis of Knowledge Retention

- 3000 data points were added in the first round of the experiment.
- Model's knowledge retention and forgetting of original domain data were monitored in every round of the experiment.
- Knowledge Retention Rate (KRR) is introduced as an additional evaluation metric.
- The traditional CIT method results in a significant decline in both ACC and KRR after multiple training rounds, indicating a severe catastrophic forgetting phenomenon and a sharp deterioration in the model's ability to recall original information.

# Conclusion

- This work introduced a novel approach that incorporates conflict-based learning to address knowledge conflicts.
- This work introduces a self-distillation based data augmentation approach to enhance training data.
- Extensive experiments were done through which the method demonstrated significant improvements in both knowledge acquisition efficiency and long term retention of previously learned information.

# References

1. Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., ... & Agarwal, S. (2020). Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 1(3), 3.
2. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35, 27730-27744.
3. Xu, R., Luo, F., Zhang, Z., Tan, C., Chang, B., Huang, S., & Huang, F. (2021). Raise a child in large language model: Towards effective and generalizable fine-tuning. arXiv preprint arXiv:2109.05687.
4. Xie, T., Wan, Y., Huang, W., Yin, Z., Liu, Y., Wang, S., ... & Hoex, B. (2023). Darwin series: Domain specific large language models for natural science. arXiv preprint arXiv:2308.13565.
5. McCann, B., Kesar, N. S., Xiong, C., & Socher, R. (2018). The natural language decathlon: Multitask learning as question answering. arXiv preprint arXiv:1806.08730.
6. Gururangan, S., Marasović, A., Swayamdipta, S., Lo, K., Beltagy, I., Downey, D., & Smith, N. A. (2020). Don't stop pretraining: Adapt language models to domains and tasks. arXiv preprint arXiv:2004.10964.
7. Longpre, S., Perisetla, K., Chen, A., Ramesh, N., DuBois, C., & Singh, S. (2021). Entity-based knowledge conflicts in question answering. arXiv preprint arXiv:2109.05052.
8. Liu, Y., Yao, Z., Lv, X., Fan, Y., Cao, S., Yu, J., ... & Li, J. (2024). Untangle the knot: Interweaving conflicting knowledge and reasoning skills in large language models. arXiv preprint arXiv:2404.03577.
9. Jiang, G., Jiang, C., Xue, S., Zhang, J. Y., Zhou, J., Lian, D., & Wei, Y. (2023). Towards anytime fine-tuning: Continually pre-trained language models with hypernetwork prompt. arXiv preprint arXiv:2310.13024.
10. Dong, G., Yuan, H., Lu, K., Li, C., Xue, M., Liu, D., ... & Zhou, J. (2023). How abilities in large language models are affected by supervised fine-tuning data composition. arXiv preprint arXiv:2310.05492.
11. Ding, R., Han, X., & Wang, L. (2022). A unified knowledge graph augmentation service for boosting domain-specific NLP tasks. arXiv preprint arXiv:2212.05251.
12. Zheng, J., Hong, H., Liu, F., Wang, X., Su, J., Liang, Y., & Wu, S. (2024). Fine-tuning large language models for domain-specific machine translation. arXiv preprint arXiv:2402.15061.
13. Emelin, D., Bonadiman, D., Alqahtani, S., Zhang, Y., & Mansour, S. (2022). Injecting domain knowledge in language models for task-oriented dialogue systems. arXiv preprint arXiv:2212.08120.
14. Han, R., Ren, X., & Peng, N. (2020). ECONET: Effective continual pretraining of language models for event temporal reasoning. arXiv preprint arXiv:2012.15283.
15. Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., ... & Hadsell, R. (2017). Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13), 3521-3526.
16. Xin, C., Lu, Y., Lin, H., Zhou, S., Zhu, H., Wang, W., ... & Sun, L. (2024, May). Beyond full fine-tuning: Harnessing the power of LoRA for multi-task instruction tuning. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)* (pp. 2307-2317).