

# **Knowledge-augmented Large Language Model**

- **Check Your Facts and Try Again: Improving Large Language Models with External Knowledge and Automated Feedback**
- **Towards Continual Knowledge Learning of Language Models**
- **Plug-and-Play Knowledge Injection for Pre-trained Language Models**

# **Check Your Facts and Try Again: Improving Large Language Models with External Knowledge and Automated Feedback\***

**Baolin Peng<sup>†</sup> Michel Galley<sup>†</sup> Pengcheng He<sup>†</sup> Hao Cheng<sup>†</sup> Yujia Xie<sup>†</sup>  
Yu Hu<sup>†</sup> Qiuyuan Huang<sup>†</sup> Lars Liden<sup>†</sup> Zhou Yu<sup>‡</sup> Weizhu Chen<sup>†</sup> Jianfeng Gao<sup>†</sup>  
<sup>†</sup> Microsoft Research   <sup>‡</sup> Columbia University**

# Introduction

- LLM에서 external knowledge 활용의 한계
  - LLM의 knowledge encoding은 손실이 많으며, knowledge generalization은 memory distortion으로 이어짐  
→ **Hallucination**
  - 증가하는 모델 크기에도 불구하고 LLM은 많은 애플리케이션에 필요한 **모든 정보를 인코딩할 수 없음**
  - 이전에 제안된 external knowledge을 활용하는 거의 모든 방법은 LLM의 parameter를 finetuning해야 하며 LLM의 크기가 기하급수적으로 증가함에 따라 엄청나게 **비용**이 많이 들 수 있음

# Introduction

- 본 논문은 plug-and-play (PnP) module로 black-box LLM을 보강한 **LLM-AUGMENTER**를 제안함
- (1) LLM이 **External Knowledge**에 근거한 응답을 생성하도록 함
- (2) **Automated Feedback**을 사용하여 LLM의 응답 수정

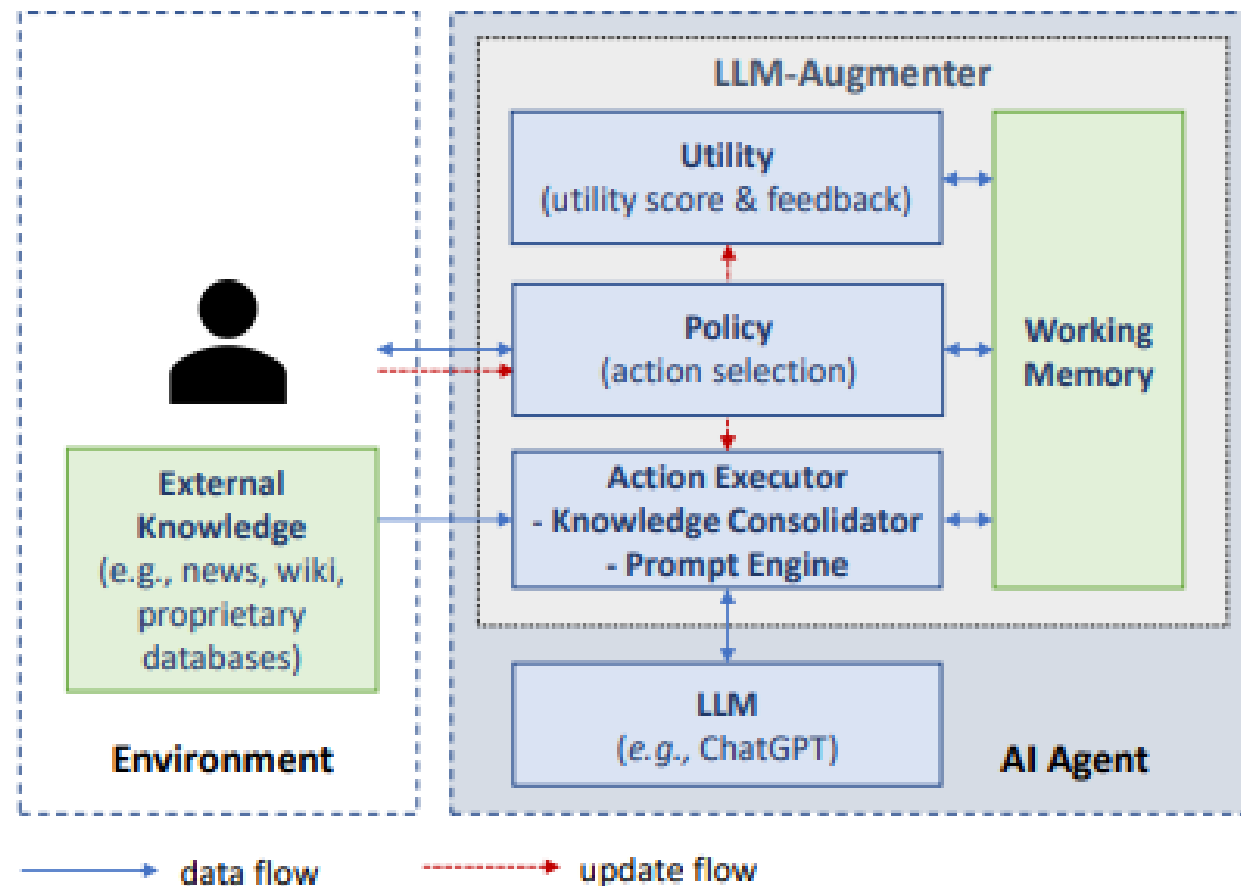


# Introduction

- LLM-AUGMENTER의 효과는 (1) **task-oriented dialog**와 (2) **open-domain question answering**에서 검증함
- LLM-AUGMENTER는 응답의 fluency와 informativeness를 희생하지 않으면서 ChatGPT의 **hallucination**을 크게 줄임

# LLM-AUGMENTER

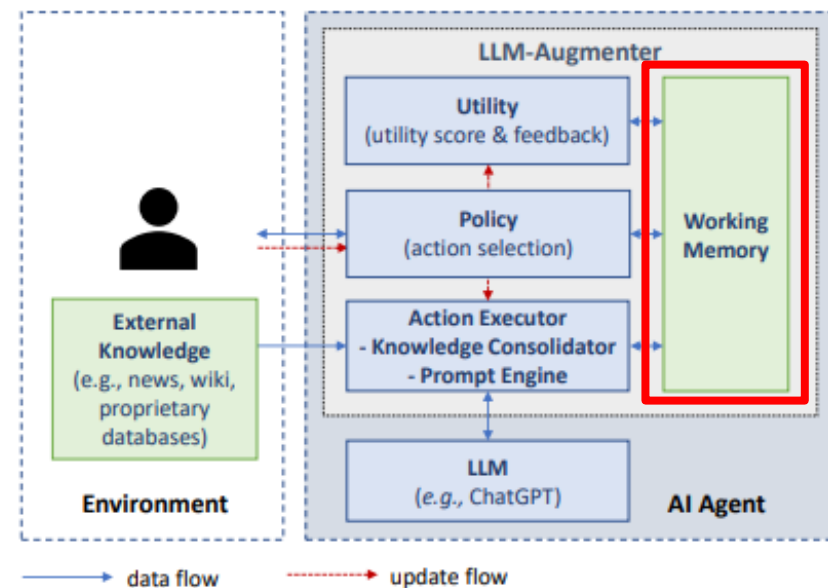
- Fixed LLM을 향상시키기 위한 PnP module의 구성 요소
  - (1) Working Memory
  - (2) Policy
  - (3) Action Executor
  - (4) Utility



# LLM-AUGMENTER

## (1) Working Memory

- 대화에 필수적인 모든 정보를 가짐
- **q**: current user query
- **e**: evidence for q, consolidated from external knowledge by Knowledge Consolidator
- **o**: a set of the LLM-generated candidate responses for q
- **u**: a score assessing the utility of each element of o
- **f**: a verbalized feedback to guide the LLM to improve its utility
- **hq**: dialog history before q

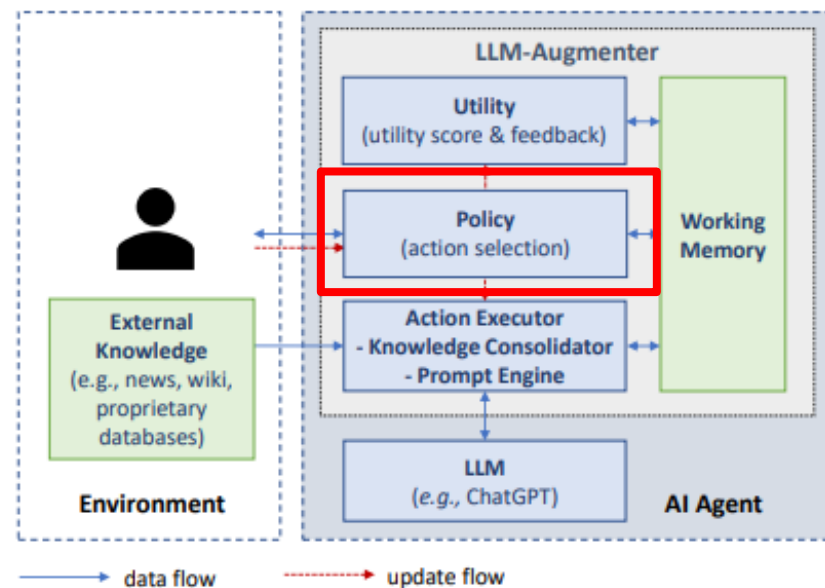




# LLM-AUGMENTER

## (2) Policy

- 이 모듈은 best expected reward R에 따라 다음 시스템의 action을 선택함
  - (1) external knowledge에서 q에 대한 evidence e 획득
  - (2) LLM을 호출하여 후보 응답 생성
  - (3) Utility module의 확인을 통과한 경우 사용자에게 응답 보내기
- 학습 가능한 정책  $\pi$ 를  $\theta$ 로 parameterize된 PLM (e.g., T5)로 구현
- 예상 보상을 최대화하기 위해 REINFORCE를 사용하여 optimized

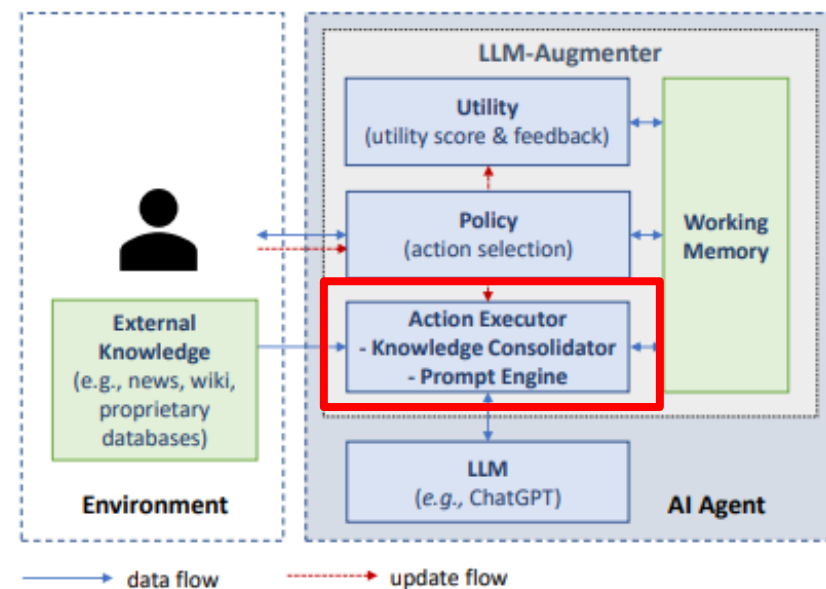


# LLM-AUGMENTER

## (3) Action Executor

- Knowledge Consolidator

- Hallucination을 완화하기 위해 external knowledge에 대한 응답을 기반으로 LLM을 강화함
- (Ma et al., 2022)에 따라 knowledge retriever, entity linker, evidence chainer로 구성됨

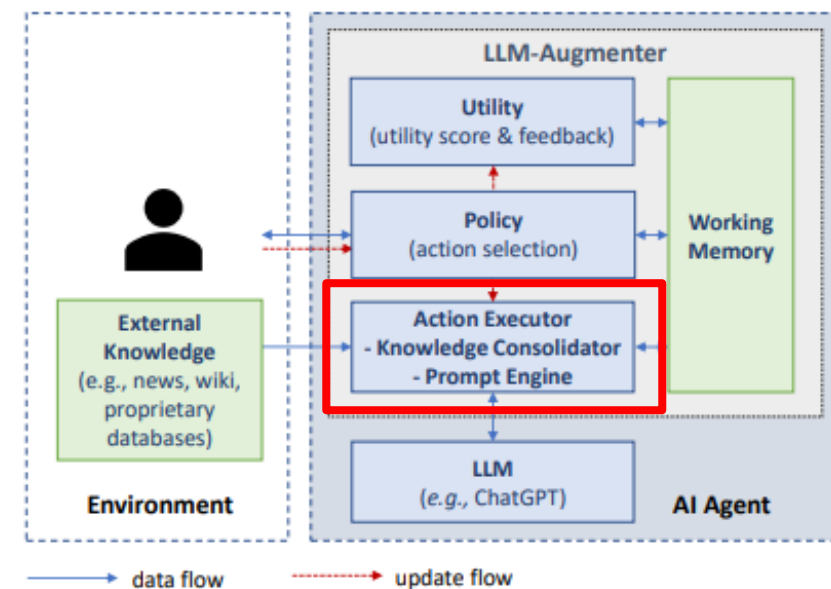


# LLM-AUGMENTER

## (3) Action Executor

- Prompt Engine

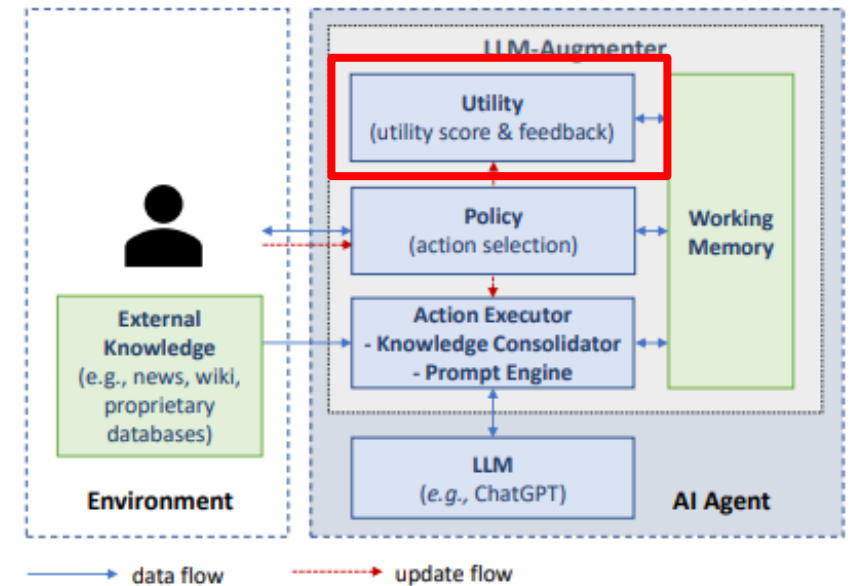
- 프롬프트는 task instruction, user query  $q$ , dialog history  $h_q$ , (Knowledge Consolidator에서 사용할 수 있는 경우) evidence  $e$ , (Utility module에서 사용할 수 있는 경우) feedback  $f$ 로 구성된 텍스트



# LLM-AUGMENTER

## (4) Utility

- **Utility Score  $u$** 와 **Corresponding Feedback  $f$**  를 생성함
- Model-based utility function or Rule-based utility function
- 텍스트 생성 모델이나 LLM, 규칙 기반 생성기를 피드백 생성에 활용



# LLM-AUGMENTER

Instruction	I want you to act as a chatbot AI for travel planning. You need to answer customer's questions nicely.
Context	<p><b>User:</b> Hello, I am looking for a moderately priced restaurant on the north side of town.</p> <p><b>Assistant:</b> I recommend golden wok. It is in the north part of town with a moderate price range.</p> <p><b>User:</b> Yes, I am looking for a moderately priced restaurant in the north part of town.</p> <p><b>Assistant:</b> I recommend golden wok. It is in the north part of town with a moderate price range.</p> <p><b>User:</b> I don't want Golden Wok, is there anything else?,</p> <p><b>Assistant:</b> Yes, your other option is The Nirala. It serves Indian food and is also located in the north part of town.</p> <p><b>User:</b> Is their beer a good value for the money?</p>
ChatGPT Response	I am sorry, I do not have information on the beer prices at The Nirala. It is best to contact the restaurant directly for more information.
Working Memory	<p>Review: They also have a modest drinks menu with some affordable prices, but I was more interested in the beautiful view from my window seat, which allowed me to watch the sun setting as I dined.</p> <p>Q: Is alcohol served at this restaurant? A: Yes, alcohol is served at this restaurant.</p> <p>Review: I found them to be overpriced and mediocre.</p> <p>Review: The drinks are priced well.</p> <p>Review: Something that I disliked, however was the prices for their drinks.</p>
LLM-AUGMENTER Response	I am sorry, I do not have information on the value of their beer. I only have information on the location and price range of the restaurant.
Rule-based Feedback	The response is inconsistent with the knowledge. Please generate again?
LLM-AUGMENTER Response	Some reviews mention that the drinks at the restaurant are priced well, while others mention that they found the prices to be overpriced and mediocre.
Human Response	Unfortunately the reviews are split, half say they are overpriced and have say they are affordable.

# Information Seeking Dialog

- News Chat

Model	K.C.	Feedback	KF1 ↑	BLEU ↑	ROUGE ↑	chrF ↑	METEOR ↑	BERTScore ↑	BARTScore ↑	BLEURT ↑	Avg. length
CHATGPT	-	-	26.71	1.01	16.78	23.80	7.34	82.14	0.25	26.98	58.94
LLM-AUGMENTER	BM25	✗	34.96	6.71	22.25	27.02	9.35	<b>83.46</b>	0.34	26.89	46.74
LLM-AUGMENTER	BM25	✓	<b>36.41</b>	<b>7.63</b>	<b>22.80</b>	<b>28.66</b>	<b>10.17</b>	83.33	<b>0.35</b>	<b>27.71</b>	54.24
LLM-AUGMENTER	gold	✗	57.44	19.24	38.89	40.02	17.21	86.65	0.82	40.55	44.35
LLM-AUGMENTER	gold	✓	60.76	21.49	40.56	42.14	18.50	86.89	0.93	42.15	47.19

- Customer Service

Model	K.C.	Feedback	KF1 ↑	BLEU ↑	ROUGE ↑	chrF ↑	METEOR ↑	BERTScore ↑	BARTScore ↑	BLEURT ↑	Avg. length
CHATGPT	-	-	31.33	4.70	24.02	27.14	12.83	87.88	1.53	<b>47.99</b>	28.81
LLM-AUGMENTER	BM25	✗	34.07	<b>4.78</b>	<b>24.52</b>	28.95	13.61	<b>87.96</b>	1.78	47.21	32.65
LLM-AUGMENTER	BM25	✓	<b>37.41</b>	3.86	24.20	<b>30.90</b>	<b>14.74</b>	87.58	<b>2.09</b>	44.71	45.07
LLM-AUGMENTER	gold	✗	45.63	6.54	29.77	33.32	16.93	89.35	2.59	54.38	33.04
LLM-AUGMENTER	gold	✓	52.83	5.63	29.65	35.68	18.66	89.01	3.14	52.49	45.09

# Wiki QA

- OTT-QA

Model	Knowledge Consolidator	Feedback	Wiki QA		
			P ↑	R ↑	F1 ↑
CHATGPT	-	-	0.48	1.52	0.59
LLM-AUGMENTER	DPR	✗	2.08	4.31	2.38
LLM-AUGMENTER	CORE	✗	7.06	14.77	8.08
LLM-AUGMENTER	CORE	✓	8.93	33.87	11.80

# Conclusions

- External knowledge과 Automated feedback으로 black-box LLM을 보강하기 위한 프레임워크인 **LLM-AUGMENTER**를 제안함
- LLM 프롬프트의 일부로 제공되는 **external knowledge**은 현재 대화와 관련된 knowledge에 더욱 기반을 둔 응답을 생성하는 데 도움이 됨
- **Automated feedback**은 ChatGPT와 같은 모델의 follow-up correction 능력을 이끌어내어 주어진 Utility function에 따라 rank가 더 높은 수정된 응답을 생성함



# TOWARDS CONTINUAL KNOWLEDGE LEARNING OF LANGUAGE MODELS

**Joel Jang<sup>1</sup> Seonghyeon Ye<sup>1</sup> Sohee Yang<sup>1</sup> Joongbo Shin<sup>2</sup>  
Janghoon Han<sup>2</sup> Gyeonghun Kim<sup>2</sup> Stanley Jungkyu Choi<sup>2</sup> Minjoon Seo<sup>1</sup>**

<sup>1</sup>KAIST AI <sup>2</sup>LG AI Research

{joeljang, vano1205, sohee.yang, minjoon}@kaist.ac.kr

{jb.shin, janghoon.han, ghkayne.kim, stanleyjk.choi}@lgresearch.ai

# INTRODUCTION

- **Real-world 시나리오에서 LLM이 word knowledge를 학습할 시 어려운 점**
  - (1) Catastrophic forgetting을 피하고
  - (2) Invariant knowledge을 보존하면서
  - (3) 새로운 knowledge을 안정적으로 획득하는 것

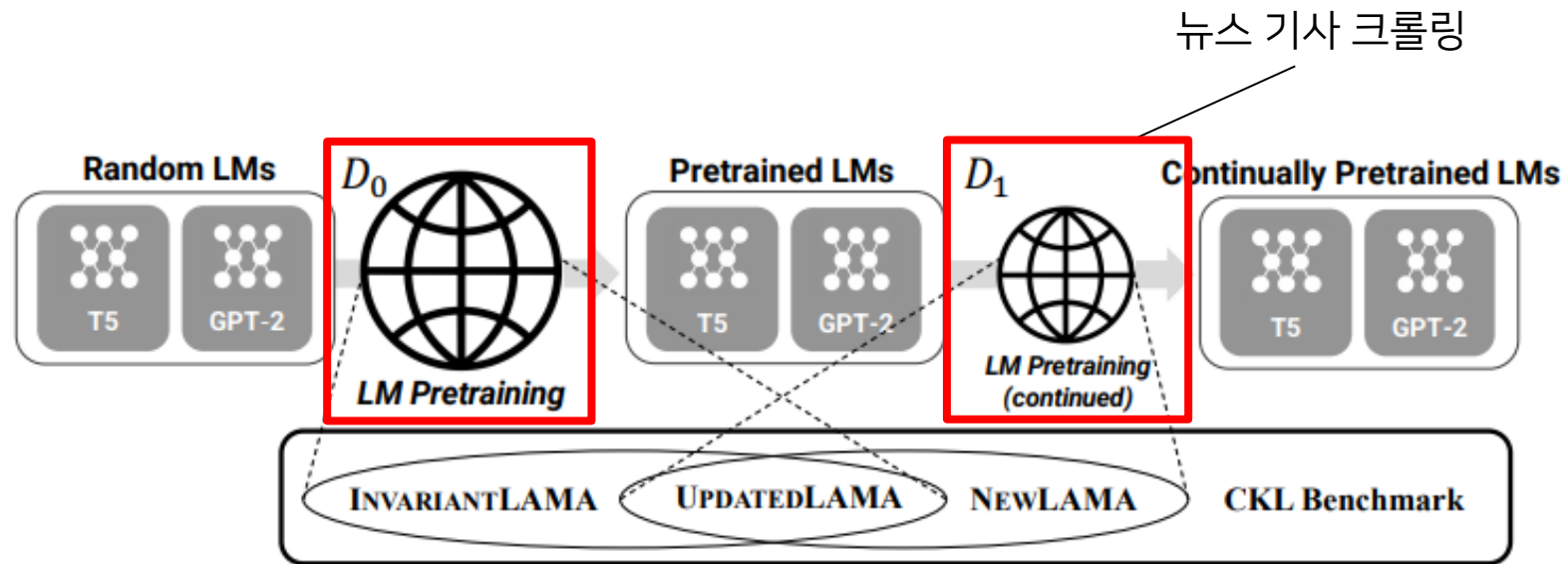
# INTRODUCTION

- 끊임없이 변화하는 LM의 더 나은 유지 관리를 위해

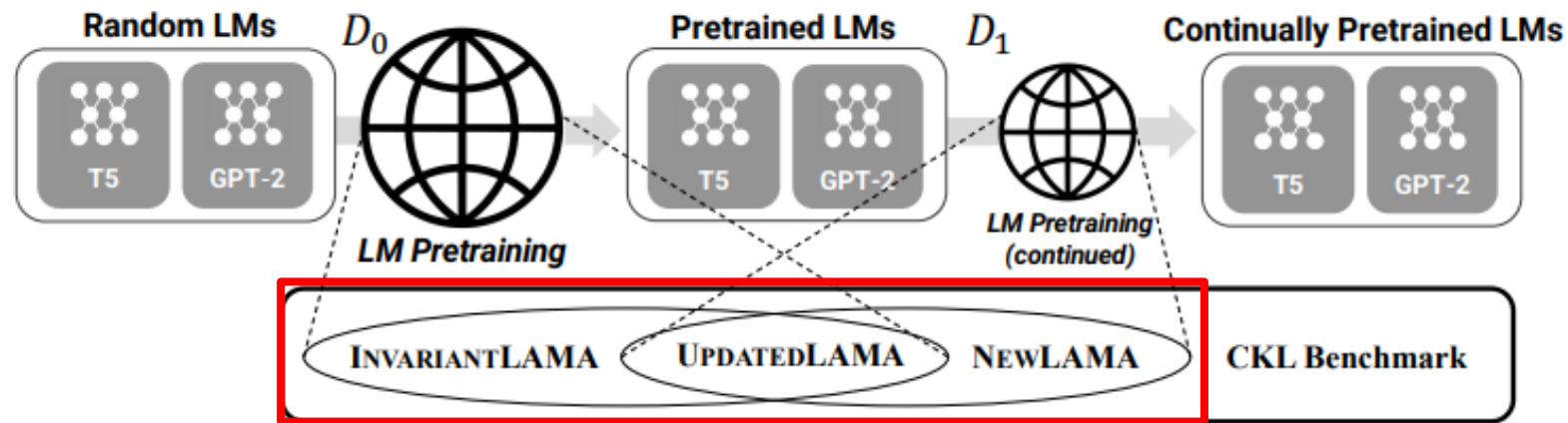
New problem formulation : **Continual Knowledge Learning (CKL)**

- (1) **time-invariant** world knowledge의 retention, (2) **update** of outdated knowledge, (3) acquisition of **new** knowledge을 정량화하기 위해 새로운 벤치마크와 메트릭을 구성함
- 적용 가능한 최근 방법을 채택하여 강력한 베이스라인을 구성함
- 실험을 통해 CKL이 지식을 안정적으로 유지하고 동시에 학습하기 위해서는 파라미터 확장이 필요함을 발견함

# CONTINUAL KNOWLEDGE LEARNING (CKL)



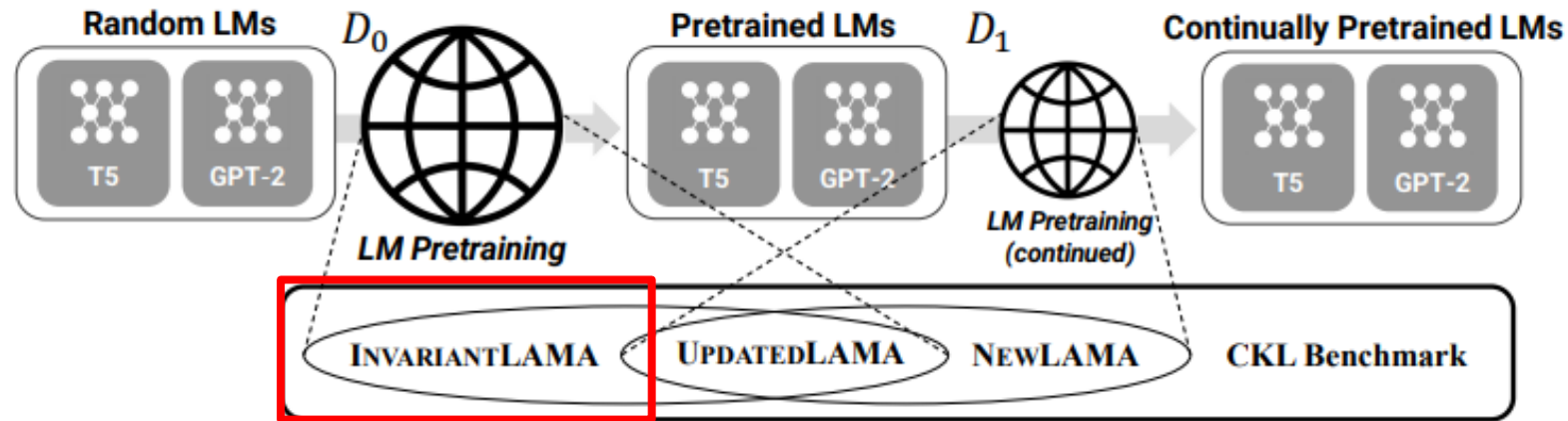
# CONTINUAL KNOWLEDGE LEARNING (CKL)



## LAMA

- Cloze sentence 안 masked entity를 제로 샷으로 맞추면 그 지식을 알고 있다고 가정함  
Dante was born in **[MASK]** → Florence

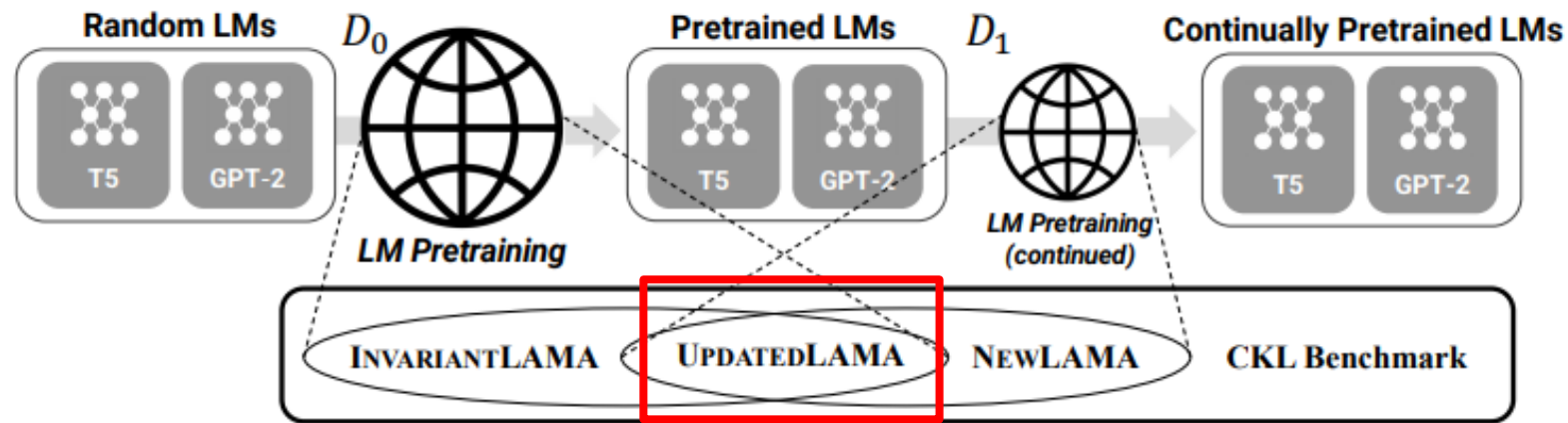
# CONTINUAL KNOWLEDGE LEARNING (CKL)



## INVARIANTLAMA

- Measuring Retention of Time-invariant World Knowledge
- Time-invariant knowledge를 가진 문장들로만 구성됨
  - iPod Touch is produced by **[MASK]** → Apple

# CONTINUAL KNOWLEDGE LEARNING (CKL)

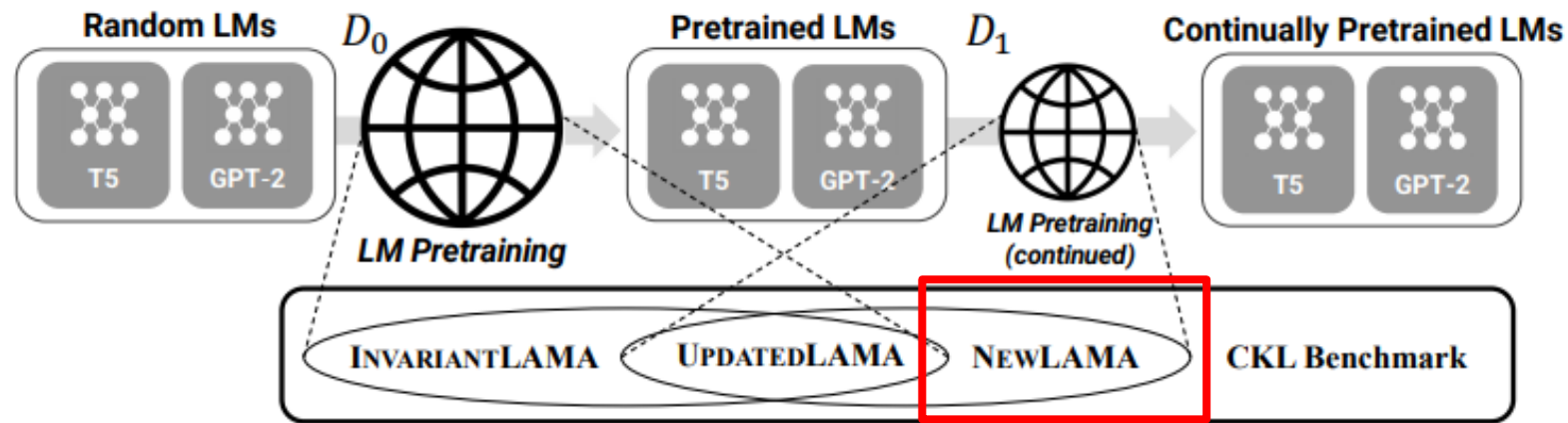


## UPDATEDLAMA

- Measuring Update of Outdated World Knowledge
- **Outdated world knowledge** : D<sub>0</sub>과 D<sub>1</sub> 사이에서 충돌하는 정보

[MASK] is the prime minister of England. → (Theresa May → Boris Johnson)

# CONTINUAL KNOWLEDGE LEARNING (CKL)



## NEWLAMA

- Measuring Acquisition of New World Knowledge
- **NEWLAMA** : 각 인스턴스에 대한 답이  $D_0$ 에는 없고  $D_1$ 에만 존재하는 것으로 확인됨  
    [MASK] owns the rights to the Falcon and the Winter Soldier. → Disney
- **NEWLAMA-EASY** : 각 인스턴스가 기준을 완벽하게 준수하지는 않음  
    Allen Lazard is officially listed as questionable with a nuclear injury after missing the last [MASK] games. → six



# CONTINUAL KNOWLEDGE LEARNING (CKL)

- COMBINED METRIC FOR CKL

- FUAR (FORGOTTEN / (UPDATED + ACQUIRED) RATIO)
- 하나의 새롭거나 업데이트된 knowledge 인스턴스를 학습하기 위해 상대적으로 얼마나 많은 time-invariant knowledge 인스턴스를 잊어버렸는지를 나타냄

$$\text{FUAR}(T^F, T_n^U, T_n^A) = \begin{cases} \frac{\sum_{i=0}^{n-1} \max(0, \text{Gap}(T_i^F, D_i, D_n)) \mathbb{1}_{\{T_i^F \neq n.d.\}}}{\sum_{i=0}^{n-1} \{ \max(0, \text{Gap}(T_n^U, D_n, D_i)) \mathbb{1}_{\{T_i^F \neq n.d.\}} + \max(0, \text{Gap}(T_n^A, D_n, D_i)) \mathbb{1}_{\{T_i^F \neq n.d.\}} \}}, \\ \text{if denominator} > 0, \\ \text{no gain, otherwise.} \end{cases}$$

# EXPERIMENTAL RESULTS

Method	# of Params (Trainable / Total)	IL	UL	NL	NLE	FUAR
		<u>EM</u>	<u>EM</u>	<u>EM</u>	<u>EM</u>	((IL), UL, NL) ↓
T5-Initial	0M / 737M	<b>24.17</b>	1.62	1.88	10.32	-
T5-Vanilla	737M / 737M	12.89	10.17	3.77	17.75	1.08
T5-RecAdam	737M / 737M	13.20	12.55	4.02	17.85	0.84
T5-MixReview	737M / 737M	13.92	6.49	2.89	14.86	1.74
T5-LoRA	403M / 738M	16.58	<b>12.77</b>	4.52	<b>19.56</b>	0.55
T5-Kadapters (k=2)	427M / 762M	19.59	12.34	<b>5.03</b>	18.75	<u>0.33</u>
T5-Kadapters (k=3)	440M / 775M	19.76	<u>12.66</u>	4.02	19.00	<u>0.33</u>
T5-Modular	438M / 773M	<u>20.29</u>	<u>12.66</u>	<u>4.65</u>	<u>19.24</u>	<b>0.28</b>

# Conclusions

- 벤치마크 데이터셋 및 메트릭을 설정하고 끊임없이 변화하는 LM의 지속적인 knowledge 학습을 위한 방법론을 탐색하는 **CONTINUAL KNOWLEDGE LEARNING (CKL)**을 제안함
- 파라미터 확장 방법이 모든 실험 설정에서 가장 강력한 성능을 보여줌

# Plug-and-Play Knowledge Injection for Pre-trained Language Models

**Zhengyan Zhang<sup>1\*</sup>, Zhiyuan Zeng<sup>1\*</sup>, Yankai Lin<sup>2,3</sup>, Huadong Wang<sup>1</sup>, Deming Ye<sup>1</sup>  
Chaojun Xiao<sup>1</sup>, Xu Han<sup>1†</sup>, Zhiyuan Liu<sup>1,4,5†</sup>, Peng Li<sup>6</sup>, Maosong Sun<sup>1,4</sup>, Jie Zhou<sup>7</sup>**

<sup>1</sup>NLP Group, DCST, IAI, BNRIST, Tsinghua University, Beijing

<sup>2</sup>Gaoling School of Artificial Intelligence, Renmin University of China, Beijing

<sup>3</sup> Beijing Key Laboratory of Big Data Management and Analysis Methods

<sup>4</sup>International Innovation Center of Tsinghua University, Shanghai <sup>5</sup> Quan Cheng Laboratory

<sup>6</sup> Institute for AI Industry Research (AIR), Tsinghua University, China

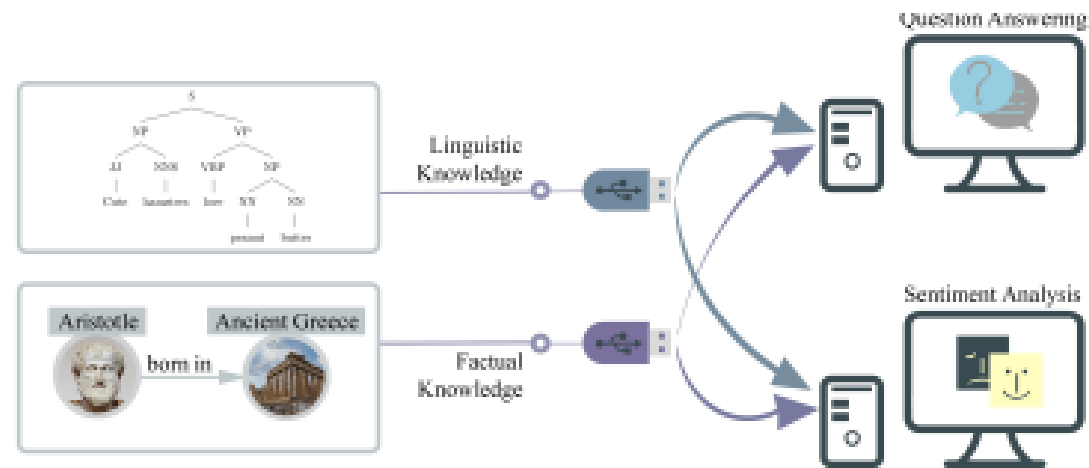
<sup>7</sup> Pattern Recognition Center, WeChat AI, Tencent Inc

{zy-z19, zengzy20}@mails.tsinghua.edu.cn {hanxu2022, liuzy}@tsinghua.edu.cn

# Introduction

- Downstream 작업을 위한 새로운 external knowledge injection 방법은 massive retraining이 필요함
- 이 논문에서는 **기존 downstream model을 reuse**하여 knowledge의 유연성과 효율성을 개선하는 방법을 처음으로 연구함

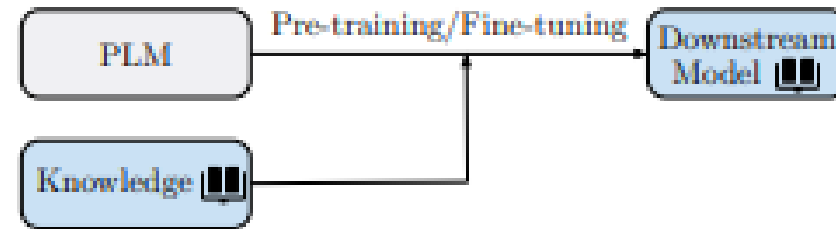
# Introduction



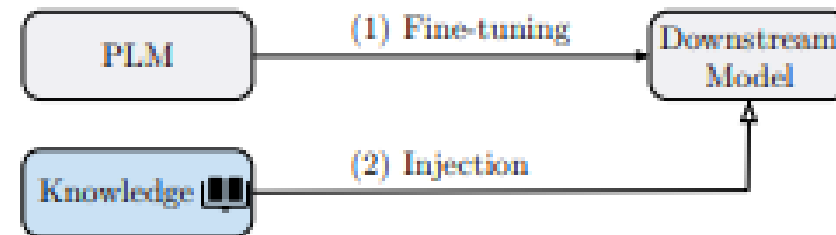
- Knowledge plugin에 의해 frozen된 기존 downstream 모델에 knowledge base가 주입되는 **plug-and-play knowledge injection**을 탐색함
- 그에 상응하여 모델 파라미터를 frozen 상태로 유지하면서 매핑된 임베딩으로 모델 입력을 풍부하게 하기 위해 knowledge 임베딩의 매핑을 훈련시키는 **plug-and-play injection method map-tuning**을 제안함
- Three knowledge-driven NLP task에 대한 실험 결과 기존 injection 방법은 새로운 패러다임에 적합하지 않은 반면 map-tuning은 downstream 모델의 성능을 효과적으로 향상시킴

# Plug-and-Play Knowledge Injection

## Previous Knowledge Injection (a)



## Plug-and-play Knowledge Injection



# Plug-and-Play Knowledge Injection

- **Paradigm Description**

- Extra knowledge base B를 통합하고 Downstream model D의 파라미터를 frozen하여 성능을 개선하려고 함
- knowledge plugin M는 훈련해야 함



# Plug-and-Play Knowledge Injection

- **Two Injection Settings**
  - **General plug-and-play knowledge injection** : M is obtained based on only P and B, and then it is directly plugged into all downstream models,  $D_1, D_2, \dots$ , without any additional training
  - **Task-specific plug-and-play knowledge injection** : It is allowed to train  $M_1, M_2, \dots$  for  $D_1, D_2, \dots$  respectively while keeping  $D_1, D_2, \dots$  frozen

# Plug-and-Play Knowledge Injection

- **Potentiality of Using Existing Methods**

- General plug-and-play knowledge injection에 사용할 수 있는 injection 방법들을 탐색함
- **(1) Embedding-based methods:** E-BERT 및 PELT는 토큰 임베딩의 representation space에 entity 임베딩 lookup table을 구축하고 entity 임베딩과 토큰 임베딩을 결합하여 입력 임베딩을 구성함
- **(2) Retrieval-based methods:** RAG는 지식 베이스에서 일반 텍스트를 검색하고 주입된 지식으로 원래 입력 텍스트를 보강함
- **(3) Adapter-based methods:** K-Adapter는 knowledge가 있는 어댑터와 함께 downstream 모델의 출력을 기반으로 knowledge이 있는 표현을 계산하며, 고정된 PLM으로 교육되고 모든 downstream 모델에 연결됨

# Map-Tuning

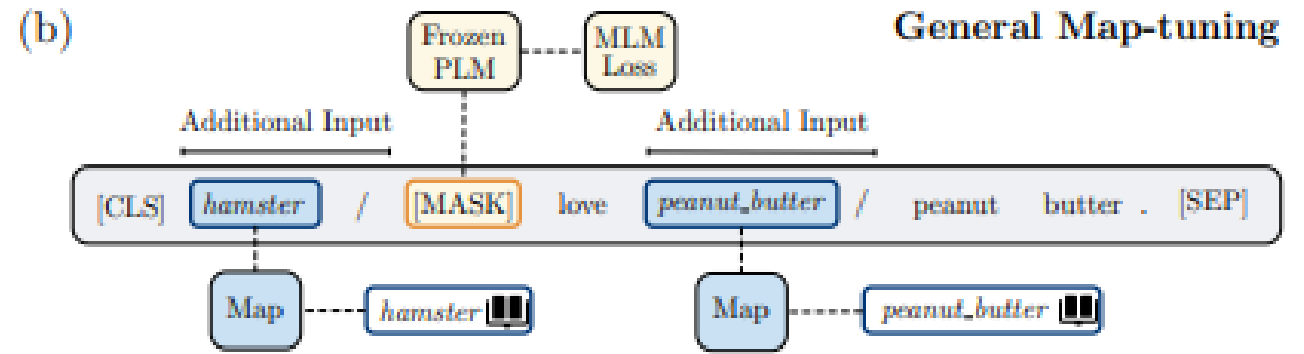
- **Overall Framework**

- Knowledge representation을 토큰 임베딩 공간에 매핑하고 매핑된 표현을 추가 입력으로 사용하여 knowledge를 injection함

# Map-Tuning

- **General Map-tuning**

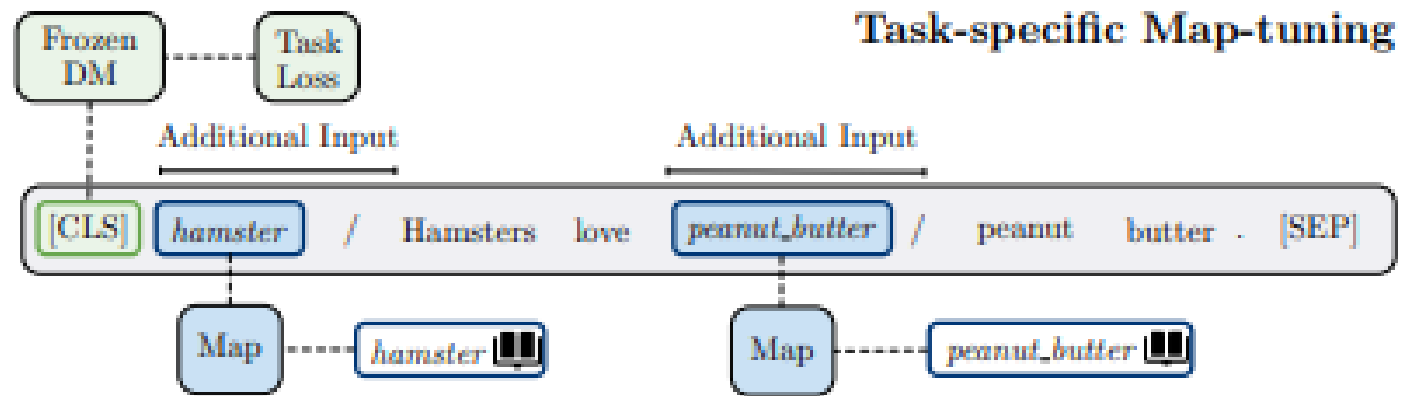
- PLM P와 knowledge representation model K를 기반으로 **매핑 네트워크 M**을 훈련하는 것을 목표로 함
- P의 파라미터를 동결하고 plug-and-play knowledge injection의 요구 사항을 충족하기 위해 매핑 네트워크 M만 함
- 매핑 네트워크가 충분히 훈련되고 매핑된 임베딩이 PLM에서 잘 활용되도록 입력 텍스트에서 **entity mention**만 랜덤 마스킹
- 이러한 방식으로 마스킹된 entity mention을 예측하는 PLM의 기능은 마스킹된 entity와 컨텍스트의 다른 entity 모두의 매핑된 임베딩을 통해 향상됨
- General map-tuning 후 **M**을 **general plug-and-play knowledge injection**에 사용할 수 있음



# Map-Tuning

- **Task-specific Map-tuning**

- 주어진 downstream 모델 D에 대해 매핑 네트워크 M을 조정하는 것을 목표로 함



# Experiments

- General Plug-and-Play Injection

Method	Injection	FewRel 1.0				Wiki80	Wiki-ET	EntityQuestions
		5-1	5-5	10-1	10-5			
Fine-tuning	–	91.0	95.1	85.4	90.8	86.1	77.5	41.7
	E-BERT	91.0 (+0.0)	95.0 (−0.1)	86.5 (+1.1)	90.5 (−0.3)	85.4 (−0.7)	77.0 (−0.5)	42.9 (+1.2)
	PELT	90.5 (−0.5)	94.8 (−0.3)	85.3 (−0.1)	89.8 (−1.0)	85.0 (−1.1)	76.8 (−0.7)	46.8 (+5.1)
	RA	91.5 (+0.5)	95.5 (+0.4)	85.8 (+0.4)	<b>91.7 (+0.9)</b>	85.9 (−0.2)	76.7 (−0.8)	<b>69.5 (+27.8)</b>
	K-Adapter	88.6 (−2.4)	94.5 (−0.6)	82.3 (−3.1)	89.9 (−0.9)	86.0 (−0.1)	<b>77.8 (+0.3)</b>	39.2 (−2.5)
	Map-tuning	<b>92.6 (+1.6)</b>	<b>95.6 (+0.5)</b>	<b>88.1 (+2.7)</b>	91.2 (+0.4)	<b>86.7 (+0.6)</b>	76.6 (−0.9)	49.0 (+7.3)
LoRA	–	90.7	95.1	84.9	91.2	85.3	77.5	42.4
	E-BERT	90.7 (+0.0)	95.2 (+0.1)	85.4 (+0.5)	90.4 (−0.8)	83.7 (−1.6)	77.6 (+0.1)	44.0 (+1.6)
	PELT	89.9 (−0.8)	94.8 (−0.3)	84.6 (−0.3)	89.8 (−1.4)	83.1 (−2.2)	77.5 (+0.0)	47.7 (+5.3)
	RA	91.3 (+0.6)	95.8 (+0.7)	85.0 (+0.1)	<b>92.5 (+1.3)</b>	83.8 (−1.5)	76.8 (−0.7)	47.7 (+5.3)
	K-Adapter	90.0 (−0.7)	94.8 (−0.3)	83.4 (−1.5)	89.1 (−2.1)	85.0 (−0.3)	77.3 (−0.2)	41.1 (−1.3)
	Map-tuning	<b>92.3 (+1.6)</b>	<b>96.0 (+0.9)</b>	<b>87.4 (+2.5)</b>	91.9 (+0.7)	<b>85.8 (+0.5)</b>	<b>78.3 (+0.8)</b>	<b>49.6 (+7.2)</b>
Adapter	–	91.2	95.2	86.2	91.1	85.7	77.5	43.6
	E-BERT	91.3 (+0.1)	95.4 (+0.2)	86.9 (+0.7)	91.6 (+0.5)	84.4 (−1.3)	78.4 (+0.9)	45.1 (+1.5)
	PELT	91.0 (−0.2)	95.4 (+0.2)	86.3 (+0.1)	91.3 (+0.2)	84.3 (−1.4)	77.9 (+0.4)	48.4 (+4.8)
	RA	91.7 (+0.5)	95.5 (+0.3)	85.8 (−0.4)	<b>92.3 (+1.2)</b>	85.0 (−0.7)	76.8 (−0.7)	42.9 (−0.7)
	K-Adapter	89.9 (−1.3)	94.7 (−0.5)	83.6 (−2.6)	90.0 (−1.1)	<b>85.9 (+0.2)</b>	77.7 (+0.2)	41.5 (−2.1)
	Map-tuning	<b>92.6 (+1.4)</b>	<b>95.8 (+0.6)</b>	<b>88.2 (+2.0)</b>	91.8 (+0.7)	<b>85.9 (+0.2)</b>	<b>79.2 (+1.7)</b>	<b>50.8 (+7.2)</b>
BitFit	–	89.2	94.8	83.0	90.0	82.7	77.1	41.3
	E-BERT	88.7 (−0.5)	94.5 (−0.3)	83.5 (+0.5)	89.6 (−0.4)	81.3 (−1.4)	77.2 (+0.1)	42.3 (+1.0)
	PELT	88.2 (−1.0)	94.3 (−0.5)	80.9 (−2.1)	88.3 (−1.7)	80.3 (−2.4)	77.6 (+0.5)	46.7 (+5.4)
	RA	89.5 (+0.3)	95.2 (+0.4)	82.7 (−0.3)	<b>91.1 (+1.1)</b>	81.8 (−0.9)	74.0 (−3.1)	33.9 (−7.4)
	K-Adapter	86.4 (−2.8)	93.7 (−1.1)	78.8 (−4.2)	87.5 (−2.5)	81.5 (−1.2)	77.2 (+0.1)	40.7 (−0.6)
	Map-tuning	<b>90.4 (+1.2)</b>	<b>95.5 (+0.7)</b>	<b>85.2 (+2.2)</b>	90.8 (+0.8)	<b>83.7 (+1.0)</b>	<b>78.0 (+0.9)</b>	<b>48.4 (+7.1)</b>

# Experiments

- Task-specific Plug-and-Play Injection

Method	Wiki80	Wiki-ET	EntityQuestions
Fine-tuning	86.1	77.5	41.7
+ General Map-tuning	86.7	76.6	49.0
+ Task-specific Map-tuning			
Train from Scratch	87.2	78.8	57.7
Train from the General Map	<b>87.8</b>	<b>78.9</b>	<b>58.9</b>

# Conclusion

- 본 논문에서는 flexible and efficient knowledge injection의 새로운 paradigm을 제안함
- 먼저 기존 방법을 평가하고 **plug-and-play injection**에 적합하지 않음을 찾음
- 그 다음 이 paradigm에 대한 **map-tuning**을 제안하여 downstream 모델에 knowledge를 효과적으로 injection하여 향상시킴



**Thank You**

---