

University of Stuttgart Institute of Industrial Automation

and Software Engineering

Customization of LLMs with Fine-tuning





**Domain Knowledge Document** 

Base LLM



synthetic data creation for training LLMs on domain-specific knowledge

Presenter: Juntao Lin Supervisor: Yuchen Xia Examiner: Prof. Dr. Ing. Michael Weyrich



# Contents

- Introduction
- Basics
- Training experiment design
- Testing and Evaluation
- Overall Results
- Findings and Outlook

## Introduction

## Current solution to achieve LLM customization



## Scenario: Pharmaceutical Manufacturing



## Proof the feasibility of Fine-tuning



Actual pharmaceutical industry processes

- Academic papers
- Patent Specification
- Authoritative databases
- etc.

## A production process of a medicine(ibuprofen)

ono-Alsobutybenzose: ndergans latonination reaction with latonine in sarkon total

Reaction Formulas

nula-ID 12 1

- Description
- Function

Fur

- · Conditions of reaction
- Reactant properties
- Influential factors



#### **Customization of LLM**

#### **QA** Chatbot

# **Basics**

- Incontext Learning(Prompting)
- Retrieval-Augmented Generation (RAG)
- Fine Tuning

## incontext Learning(Prompting)



Characteristic of incontext learning in Customized LLMs

- No need to change model parameters
- Fast and easy to deploy
- High requirements on example design



## **Retrieval-Augmented Generation (RAG)**



Characteristic of RAG in Customized LLMs

- No need to change model parameters
- Suitable for knowledge-intensive tasks
- High requirements for the retrieval system



## **Fine-tuning Approach**

## Literature Review









# **Training experiment design**

- Base model
- Training resources
- Training Data
- Experiments execution

## **Base Model**

## GPT-4o-mini vs. Qwen2.5-7B-Instruct

index	GPT4o-mini	Qwen2.5-7B-Instruct	
Model parameter size	about 8B	7B	
Developer	OpenAl(2024)	Alibaba(2024)	
Call resource consumption	\$0.26 /1M tokens	\$0.03/1M tokens	
Max. Context length	128k	32k	
Specialize in task types	Question answering programming tasks reasoning tasks	Question answering dialogue generation	
Pre-Trained data source	Wide range Covering many fields	Wide range Covering many fields Focus on Chinese corpus	

## **Training resources**

- For Qwen2.5-7B-Instruct
- Cloud Computing Platform: Runpod
- Fine-tuning Platform:LLaMa-Factory
- Training environment configuration:

GPU	2 x A100 SXM
vRAM	160GB
pytorch version	2.4.0
python version	3.11
CUDA	12.4.1
Ubuntu	22.04

			••••	
RunPod Pytorch 2.4.0 🖍 ID: c8ja104671mazk	<b>2 x A100 SXM</b> 64 vCPU 251 GB RAM	runpod/pytorch:2.4.0-py3.11-cuda12.4.1-devel-ubuntu22. On-Demand - Secure Cloud	04	Running
	- <del>1</del>		GPU Memory Used	
30 GB Disk 130 GB Pod Volume			Avg	97%
Volume Path: /workspace			0	97%
			1	96%
🕤 EUR-IS-1 🔹 4738 Mbps 🔹 8558 Mbps 📑	4337 MBps			

- > 70GB disk volume for each check-point step of 16-bit quantized model
- For GPT-4o-mini
  - Fine-tuning Platform: OpenAl Platform



\$3.78/hr





## **Training Data**

## Dataset comparison

index	PDF-synthetic dataset	Knowledge-distillation dataset	
Sample number	155 trainset(101) validation set(30) testset(25)	1,730 Trainset(1,440) validation set(180) testset(110)	
Dataset size	12,764	172,113	
Dataset Format	JSONL sharegpt	JSONL sharegpt	

## **Experiments execution**

• GPT4o-mini with PDF-synthetic dataset







- Qwen2.5-7B with Knowledge-distillation dataset
  - Best performance of 6 experiments:
    - Exp1:Epoch=1,Batch\_size=1,Learning\_rate=1.8e-5
    - Exp2:Epoch=3,Batch\_size=1,Learning\_rate=1.8e-5
    - Exp3:Epoch=5,Batch\_size=1,Learning\_rate=1.8e-5
    - Exp4:Epoch=3,Batch\_size=2,Learning\_rate=1.8e-5
    - Exp5:Epoch=3,Batch\_size=4,Learning\_rate=1.8e-5

Exp6:Epoch=1,Batch\_size=1,Learning\_rate=1.0e-5
University of Stuttgart, IAS

	hyperpar ameters	GPT4o- mini +PDF synthetic datase	GPT4o- mini+Kno wledge- distillatio n dataset	Qwen2.5- 7B+Knowledg e-distillation dataset
Best performance of 6 expe	epoch	3	1	1
• Exp1:Epoch=1,Batch s	batch_size	1	1	1
• Exp2:Epoch=3,Batch	learning_rat e	1.8e-5	1.8e-5	1.0e-5
<ul> <li>Exp3:Epoch=5,Batch s</li> </ul>	ize=1,Learnin	g rate=1.8e-5		

- Exp4:Epoch=3,Batch size=2,Learning rate=1.8e-5
- Exp5:Epoch=3,Batch\_size=4,Learning\_rate=1.8e-5
- Exp6:Epoch=3,Batch\_size=1,Learning\_rate=1.0e-5

Best performance of 6 experiments:

٠

- Exp1:Epoch=1,Batch\_size=1,Learning\_rate=1.8e-5
- Exp2:Epoch=3,Batch\_size=1,Learning\_rate=1.8e-5
- Exp3:Epoch=5,Batch\_size=1,Learning\_rate=1.8e-5
- Exp4:Epoch=3,Batch\_size=2,Learning\_rate=1.8e-5
- Exp5:Epoch=3,Batch\_size=4,Learning\_rate=1.8e-5

Exp6:Epoch=3,Batch\_size=1,Learning\_rate=1.0e-5



# **Testing and Evaluation**

- Question answering quality (LLM-as-a-judge)
- Task solving (success rate)

## Test & evaluation method

- Question answering quality
  - LLM as a Judge[7]

"question": "How are ibuprofen esters prepared?"

"reference answer": "(...)The process involves mixing nickel chloride hexahydrate, 2,2'-bipyridine, sodium chloride, manganese powder, and triethylamine hydrochloride in tetrahydrofuran under nitrogen at 40°C. This mixture is then combined with 1-bromo-4-isobutylbenzene and methyl 2chloropropionate, and the reaction proceeds at 40°C for 12 hours. (...)



"answer before SFT": "(...)This process usually involves the use of a catalyst, such as sulfuric acid, under reflux conditions to promote the reaction. The reaction can be conducted by combining ibuprofen with a suitable alcohol (like methanol, ethanol, or isopropanol), leading to the formation of ibuprofen esters. Post-reaction, purification methods like distillation or chromatography may be employed to isolate the desired ester."

Final score:4

GPT-40

"answer after SFT": "(...)Ibuprofen esters are prepared through a nickel-catalyzed coupling reaction. 1-bromo-4-isobutylbenzene is coupled with methyl 2-chloropropionate using nickel(0) catalyst, ligands, metal salts, and additives. The reaction occurs in a solvent mix of tetrahydrofuran and water under an argon atmosphere at 40°C for 12 hours. This forms ibuprofen acid methyl ester, with by-products including manganese(0) and 1,3dioxide."

#### Final score:8

Metrics	Description
Factual Accuracy(0-10)	Is the answer correct and consistent with the reference?
Completeness(0-10)	Does the answer cover the main points?
Clarity(0-10)	Is the answer clear and well-structured?

## Test result by training set

LLMs	Dataset	Improved performance
		training set question
GPT4o-mini	PDF_synthetic	5.86>7.29 +14.3%
	Knowledge_ distillation	5.86>7.18 +13.2%
Qwen2.5-7B	Knowledge_ distillation	4.40>6.18 +17.8%

## Response before SFT

#### User

What is the solvent used in the preparation of 1-bromo-4-isobutylbenzene?

#### Assistant

In the preparation of 1-bromo-4-isobutylbenzene, a common solvent used is acetone or dichloromethane (DCM). These solvents facilitate the bromination reaction by dissolving the reactants and enabling better interaction between them. Depending on the specific reaction conditions, other solvents such as toluene may also be used. Always consider the reaction scale and conditions when choosing a solvent.

## Response after SFT

# User Abisticate Ansioner Image: Construction of 1-bromo-4-isobutylbenzere Image: Constructi

## Generalization

### - Question -

How to test the generalizability of models after SFT?

## Training generalization: Ask questions on the test set

LLMs	Dataset	Improved performance		
		training set question	test set question	
GPT4o-mini	PDF_synthetic	5.86>7.29 +14.3%	5.25>5.96 +7.1%	
	Knowledge_ distillation	5.86>7.18 +13.2%	5.25>7.29 +20.4%	
Qwen2.5-7B	Knowledge_ distillation	4.40>6.18 +17.8%	4.38>5.58 +12%	

#### Response before SFT

User

how to improve the final yield and purity of ibuprofen?

#### Assistant

To improve the final yield and purity of ibuprofen, consider the following strategies:

#### 1. Optimize Reaction Conditions: Control temperature, pH, and reaction time carefully.

- . Optimize neaction conditions, control temperature, pri, and reaction time carefully.
- 2. Purification Techniques: Utilize recrystallization or chromatography for better separatio
- 4. Workup Efficiency: Improve extraction methods to reduce losses during purification.
- 5. Monitor Reaction Progress: Use analytical techniques (e.g., HPLC) to track purity and yi
- 6. Prevent Side Reactions: Avoid conditions that lead to byproducts.

3. Use High-Quality Reagents: Ensure raw materials are pure and compatible.

Implementing these strategies can enhance both yield and purity.

#### Response after SFT



how to improve the final yield and purity of ibuprofen?

Assistant To enhance the yield and purity a bifen, ensure precise control of reaction temperatures (30-40°C for promination and 20-40°C for hydrelysis) and maintain uniform stirring. Adjust pH (1-2 for ibuprofen precipitation) to optimally form the product. Thoroughly wash with solvents (water, NaHCO3, time) to remove impurities and ensure complete methanol solvent removal to prevent side reactions or solvent interference. Employ careful distillation to obtain pure ibuprofen. These improvements can enhance product quality and yield by minimizing side reactions and impurities.

## Verification: Improved training generalization of fine-tuned models

## **Test & evaluation method**

Question

How to test the **generalizability** of models after SFT?

## Language generalization:ask different questions on same topic

Variants generation



## Test result by rephrased question

## Response before SFT -

#### User

In the preparation of 1-bromo-4-isobutylbenzene, what solvent is utilized?

LLMs	Dataset	Improved perform		Assistant ≎ In the preparation of 1-bromo-4-isobutylbenzene, a common solvent used is dichloromethane (DCM chloroform. These solvents effectively dissolve both the starting materials and reagents, facilitating b bromination reaction. Alternatively, other polar aprotic solvents like acetone may also be used, deper	
		training set question	test set question	the specific reaction conditions required. rephrased question	
GPT4o-mini	PDF_synthetic	5.86>7.29 +14.3%	5.25>5.96 +7.1%	5.88>7.16 +12.8% > Response after SFT	
	Knowledge_ distillation	5.86>7.18 +13.2%	5.25>7.29 +20.4%	User In the preparation of 1-bromo-4-isobutylbenzene, what solvent is utilized?	
Qwen2.5-7B	Knowledge_ distillation	4.40>6.18 +17.8%	4.38>5.58 +12%	In the preparation of 1-bromo-4-isobutylbenzene, carbon tetrachloride is used as the solvent.	

## Verification: Improved language generalization of fine-tuned models

## **Test & evaluation method**

- Question

How to test the **generalizability** of models after SFT?

Knowledge generalization:ask for task solving

## Quiz example

Question: In the preparation of 1-bromo-4-isobutylbenzene, if you start with 26.8g (0.2 mol) of isobutylbenzene, how much bromine (Br2) in moles is required to ensure the correct molar ratio is maintained?

A. 0.100 mol

B. 0.200 mol

C. 0.210 mol

D. 0.105 mol



Reasoning process

Step 1: From the provided knowledge, the molar ratio of isobutylbenzene to bromine is 1:1.05. This means that for every mole of isobutylbenzene, 1.05 moles of bromine are required.

Step 2: You start with 0.2 moles of isobutylbenzene. To find the required moles of bromine, multiply the moles of isobutylbenzene by the ratio:

 $[ \text{Required moles of Br}_2 = 0.2 \, \text{mol} \times 1.05 = 0.21 \, \text{mol} ] Step 3: Therefore, 0.210 moles of bromine are required to maintain the correct molar ratio in the reaction. Thus, the correct choice is C.$ 



success rate

## Test result by task solving

LLMs	Dataset	Improved performance		
		Task	<b>Task+</b> prompting	
GPT4o-	PDF_synthetic	Success Rate:	Success	
mini		58/87>50/87	Rate:65/87	
	Knowledge_	Success Rate:	Success Rate:	
	distillation	58/87>57/87	70/87	
Qwen2.5-	Knowledge_	Success Rate:	Success Rate:	
7B	distillation	54/87>58/87	70/87	

#### Verification: The fine-tuned model + **incontextual learning(prompting)** enhances knowledge generalization

## Task solving without prompting

#### User ≎

please choose the option you think is correct from the following multiple-choice questions and tell me the reasoning process.

Question: In the preparation of 1-bromo-4-isobutylbenzene, if you start with 26.8g (0.2 mol) of isobutylbenzene, how much bromine (Br2) in moles is required to ensure the correct molar ratio is maintained?

4. 0.100 mol
3. 0.200 mol
C. 0.210 mol
D. 0.105 mol

#### Assistant

To maintain the correct molar ratio in the preparation of 1-bromo-4-isobutylbenzene, isobutylbenzene (0.2 mol) reacts with bromine. The reaction requires a 1:1 molar ratio. Therefore, to match the 0.2 mol of bromine.

#### Option B (0.200 mol) is correct.

#### task solving with prompting

#### User

please choose the option you think is correct from the following multiple-choice questions and tell me the reasoning process.Please use the knowledge to solve the problem.The knowledge applied in this question is:In the preparation reaction of 1-bromo-4-isobutylbenzene, the molar ratio of isobutylbenzene: bromine = 1:1.05.

Question: In the preparation of 1-bromo-4-isobutylbenzene, if you start with 26.8g (0.2 mol) of isobutylbenzene, how much bromine (Br2) in moles is required to ensure the correct molar ratio is maintained?

A. 0.100 mol B. 0.200 mol C. 0.210 mol D. 0.105 mol

#### Assistant

to ensure the molar ratio of isobutylbenzene to bromine is 1:1.05, we start with 0.2 mol of isobutylbenzene. Defoulate the required moles of bromine:

\text{Moles of Br}\_2 = 0.2 \text{ mol} \times 1.05 = 0.21 \text{ mol}

Thus, the correct choice is C. 0.210 mol, ensuring the appropriate molar ratio is maintained for the reaction.

00

# **Overall Results**

## **Orthogonal experiment results**

#### Catastrophic forgetting? \_\_\_\_

Models may forget old knowledge or tasks while learning a new task.

LLMs	Dataset	Improved performance					
		train set QA	test set QA	test with rephrased QA	Task	Task+ prompting	Task(MMLU)
GPT4o-	PDF_synthetic	5.86>7.29	5.25>5.96	5.88>7.16	Success Rate:	Success	56.0>51.8
mini		+14.3%	+7.1%	+12.8%	58/87>50/87	Rate:65/87	-4.2%
	Knowledge_	5.86>7.18	5.25>7.29	5.88>7.31	Success Rate:	Success Rate:	56.0>55.5
	distillation	+13.2%	+20.4%	+14.3%	58/87>57/87	70/87	-0.5%
Qwen2.5-	Knowledge_	4.40>6.18	4.38>5.58	3.62>5.89	Success Rate:	Success Rate:	20.3>23.7
7B	distillation	+17.8%	+12%	+22.7%	54/87>58/87	70/87	+3.7%
BioMedLM[8]	Clinical cases generation with GPT-4	3.13	2.71	3.01	Success Rate: 44/87	N.A.	N.A.

Conclusion: The phenomenon of catastrophic forgetting is NOT significant

# Findings and Outlook

## **Findings and Outlook**

- Achievement:
  - > SFT successfully experimented:
    - ➤ 3 fine-tuned Models:
      - language generalization 1,
      - training generalization f,
      - knowledge generalization-,
    - > fine-tuning+incotext learning(Prompting):
      - knowledge generalization↑
    - ➤ Cost:182€

## **Findings and Outlook**

- Findings:
  - Comparative experiment:knowledge\_distillation dataset performs better than PDF synthesis dataset
  - > Vertical domain-specific LLM (BioMedLM) does not have sufficient performance in our use case.
  - > In task solving: incontext learning(prompting) **outperforms** fine-tuning approach
- Outlooks:
  - Better finetune LLM with all data
  - Powerful base LLM+incontext learning(prompting)



**University of Stuttgart** Institut of Industrial Automation and Software Engineering

## Thank you!



#### Juntao Lin

e-mail st176526@stud.uni-stuttgart.de phone +49 (0) 711 685fax +49 (0) 711 685-

University of Stuttgart Institute of Automation and Software Systems Pfaffenwaldring 47, 70550 Stuttgart



## References

- [1] O. Ovadia, M. Brief, M. Mishaeli, O. Elisha, et al., "Fine-Tuning or Retrieval? Comparing Knowledge Injection in LLMs," arXiv e-prints, 2023, arXiv:2312.05934.
- [2] Z. Wang, B. Bi, S. Kumar Pentyala, K. Ramnath, S. Chaudhuri, S. Mehrotra, Z. (James)Zhu, X. Mao, S. Asur, N. (Claire)Cheng, et al., "A Comprehensive Survey of LLM Alignment Techniques: RLHF, RLAIF, PPO, DPO and More," arXiv e-prints, 2024, arXiv:2407.16216.
- [3] C. Peng, X. Yang, K. E Smith, Z. Yu, A. Chen, J. Bian, Y. Wu,et al., "Model Tuning or Prompt Tuning? A Study of Large Language Models for Clinical Concept and Relation Extraction," arXiv e-prints, 2023, arXiv:2310.06239.
- [4] D. Yin, X. Han, B. Li, H. Feng, J. Bai,et al., "Parameter-efficient is not sufficient: Exploring Parameter, Memory, and Time Efficient Adapter Tuning for Dense Predictions," arXiv e-prints, 2023, arXiv:2306.09729.
- [5] H. Zhong, Z. Shan, G. Feng, W. Xiong, X. Cheng, L.Zhao, D. He, J. Bian, L. Wang, "DPO Meets PPO: Reinforced Token Optimization for RLHF," arXiv e-prints, 2024, arXiv:2404.18922.
- [6] Q. Yin, X. He, L. Deng, C. Tou Leong, F. Wang, Y. Yan, X. Shen, Q. Zhang, et al., "Deeper Insights Without Updates: The Power of In-Context Learning Over Fine-Tuning," arXiv e-prints, 2024, arXiv:2410.04691.
- [7] J. Gu, X. Jiang, Z. Shi, H. Tan, X. Zhai, C. Xu, W. Li, Y. Shen, S. Ma, H. Liu, S. Wang, K. Zhang, Y. Wang, W. Gao, L. Ni, J. Guo, et al., "A Survey on LLM-as-a-Judge," arXiv e-prints, 2025, arXiv:2411.15594.
- [8] J. Khlaut, C. Dancette, E. Ferreres, A. Bennani, P. Hérent, P. Manceron, "Efficient Medical Question Answering with Knowledge-Augmented Question Generation," arXiv e-prints, 2024, arXiv:2405.14654.