# **Introduction to the Student Projects**



#### IE686 Large Language Models and Agents



# **Student Projects**



#### Goals

- Gain practical experience with developing multi-agent LLM applications for a specific problem
- Learn about problem domain
  - What domain-specific problems exist and make the task hard
  - Come up with ways to try to solve them with agents

#### • Expectation

- You explore your assigned topic area and its problems with data/benchmarks/LLMs of your own choosing
- You solve the problem by
  - Building on the general methods we have seen in the course
  - Exploring additional methods that may be specific to the domain
  - Comparing multiple approaches using some form of evaluation

#### Procedure



- Teams of four to five students
  - realize a project using LLMs as agents
  - write a 12-page summary of the project and the methods employed in the project
  - present the project results to the other students
    - 12 minutes presentation + 3 minutes discussion
- Final mark for the course
  - 70 % written final report about the project
  - 30 % project presentation

#### **Exam Registration**



- There is no "real" exam at the end of this course
- You still need to formally register with the Prüfungsamt separate of your course registration
- The registration period started on 23.10.2024 and continues until 06.11.2024
- Make sure to **register soon**, so you don't forget!

### **Course Schedule**



You	u are here				
Day	Торіс				
24.10	Project: Introduction to Student Projects				
Monday, October 28 <sup>th</sup> , 23:59 Submission of Project Outlines					
31.10	Project: Project Coaching				
07.11	Project: Project Coaching				
14.11	Project: Project Coaching				
21.11	Project: Project Coaching				
28.11	Project: Project Coaching				
Sunday, December 1 <sup>st</sup> , 23:59 Submission of Project Report					
05.12	Project: Presentation of Project Results				
Thursday, December 5 <sup>th</sup> , 23:59 Submission of Presentation Slides					

## **Team Formation**



- Team formation is already done based on your preferences
  - You should have received a mail with your team number and topic area
  - It was difficult to satisfy everyone's wishes at once, so receiving the top priority topic area was often not possible
  - I had to get a bit creative but nearly everyone has gotten one of their preferences
- Now it is time to define your own project in your assigned area!
- Meet with your team after the session to organize your work!
  - Decide project specifics
  - Organize writing of project outline

# **Topic 1: Question Answering**



- Task: Build a (Multi-)Agent QA system.
  - Create an agent-based environment for QA challenges
  - Make use of RAG, search engines, etc.
  - Identify problems with each approach
- Dataset/APIs:
  - Use existing benchmarks and relevant metrics
  - <u>https://www.aicrowd.com/challenges/meta-comprehensive-rag-benchmark-kdd-cup-2024</u>

#### • Evaluation:

- Relevant evaluation metrics in used benchmarks
- Compare different setups

# **Topic 2: Multi-Agent Gaming**



- **Task:** Build a Multi-Agent gaming application.
  - Have teams of agents cooperate to solve a game or play against each other
  - Seminar student in last semester tried playing Dungeons and Dragons with Agents with good success
  - Can be explorative or you use an existing agent benchmark
- Dataset/APIs:
  - <u>https://github.com/THUDM/AgentBench</u>
  - <u>https://github.com/microsoft/SmartPlay</u>
- Evaluation:
  - Existing evaluation when using agent benchmark
  - Otherwise depends on game, could be win rate, ELO ranking, ...

# **Topic 3: Online Shopping Assistant**



- **Task:** Build a (Multi-)Agent system that supports a user in making shopping decisions, e.g. for a new TV.
  - Search and present relevant products
  - Based on the users wishes
  - Present in structured format or directly perform the transaction

#### • Dataset/APIs:

- Use existing benchmarks and relevant metrics
- <u>https://github.com/THUDM/AgentBench</u>
- <u>https://webarena.dev/</u>
- <u>https://webshop-pnlp.github.io/</u>
- Evaluation:
  - Relevant evaluation metrics in used benchmarks

# **Topic 4: Job Hunting Agents**



- **Task:** Build a (Multi-)Agent application to search for job postings relevant to a user.
  - start with a user query
  - search for relevant postings
  - extract relevant facts about the jobs
  - present to user in a structured, easy-to-browse format
- Dataset/APIs:
  - select 2 job providers/APIs to work with
  - <u>https://apislist.com/category/28/jobs</u>
- Evaluation:
  - Explorative Topic, likely no relevant benchmarks exist
  - Human-based evaluation by the group members

# **Topic 5: Text to SQL**



- **Task:** Build a (Multi-)Agent system that converts natural language queries to SQL.
  - Convert query
  - Query database
  - Refine based on result/errors
- Dataset/APIs:
  - https://paperswithcode.com/task/text-to-sql
  - https://github.com/THUDM/AgentBench
- Evaluation:
  - Relevant evaluation metrics in used benchmarks

# **Topic 6: Code-driven Developer**



- **Task:** (Multi-)Agent system for software engineering. Start with informal description of functionality.
  - Generate code solutions based on descriptions
  - Generate and run tests on solutions
  - Identify potential discrepancies in behavior of code solutions
  - Recommend final solution

#### Dataset/APIs/Language:

- <u>https://evalplus.github.io/</u>
- Python or Java
- Evaluation:
  - Relevant evaluation metrics in used benchmarks

# Where to find Datasets/Benchmarks?



- Agent Benchmarks
  - AgentBench
  - WebShop
  - <u>WebArena</u>
  - Meta RAG KDD Cup
  - <u>SmartPlay</u>
  - <u>EvalPlus</u>
- General Task Benchmarks
  - Papers with Code Datasets (filter by task)
  - Huggingface Datasets

#### **Even More Benchmarks**



Method	Evaluation	Model Types	Abilities/Domain	Data Source
	MMLU [364]	Base/Fine-tuned/Specialized	General	Human exam/practice
	BIG-bench [70]	Base/Fine-tuned/Specialized	General	Human annotation
	HELM [520]	Base/Fine-tuned/Specialized	General	Benchmark collection
	Open LLM Leaderboard [707]	Base/Fine-tuned/Specialized	General	Benchmark collection
	AGIEval [708]	Base/Fine-tuned/Specialized	General	Human exam/practice
	MMCU [709]	Base/Fine-tuned/Specialized	General	Human exam/practice
	M3KE [710]	Base/Fine-tuned/Specialized	General	Human exam/practice
	C-Eval [711]	Base/Fine-tuned/Specialized	General	Human exam/practice
	Xiezhi [712]	Base/Fine-tuned/Specialized	General	Human exam/practice
	OpenCompass [713]	Base/Fine-tuned/Specialized	General	Benchmark collection
	Chain-of-Thought Hub [714]	Base/Fine-tuned	General	Benchmark collection
	KoLA [715]	Base/Fine-tuned	Knowledge utilization	Web
	ARB [716]	Fine-tuned	Complex reasoning	Human exam/practice
Benchmark	APIBench [717]	Base/Fine-tuned	Tool manipulation	Web
Denciunark	APIBank [718]	Fine-tuned	Tool manipulation	Synthesis
	ToolAlpaca [719]	Base/Fine-tuned	Tool manipulation	Synthesis
	T-Bench [720]	Fine-tuned	Tool manipulation	Synthesis
	ToolBench [721]	Fine-tuned	Tool manipulation	Synthesis
	BOLAA [722]	Base/Fine-tuned	Environment interaction	Benchmark collection
	AgentBench [723]	Base/Fine-tuned	Environment interaction	Human annotation/Synthesis
	HaluEval [602]	Base/Fine-tuned	Human alignment	Human annotation/Synthesis
	PromptBench [724]	Base/Fine-tuned	Robustness	Benchmark collection
	HumanEval [105]	Base/Fine-tuned/Specialized	Code synthesis	Human annotation
	MultiMedQA [356]	Specialized	Healthcare	Benchmark collection
	FLUE [725]	Specialized	Finance	Benchmark collection
	LegalBench [726]	Specialized	Legal	Human annotation
Human	Chatbot Arena [727]	Base/Fine-tuned/Specialized	Human Alignment	Human annotation
Tuman	SciBench [728]	Fine-tuned	Complex reasoning	Human exam/practice
	AlpacaEval [729]	Fine-tuned	Instruction following	Synthesis
	MT-bench [727]	Fine-tuned	Human alignment	Human annotation
Model	TrustGPT [730]	Base/Fine-tuned	Human alignment	Benchmark collection
	LMExamQA [731]	Base/Fine-tuned	Knowledge utilization	Synthesis
	ChatEval [732]	Base/Fine-tuned	Knowledge utilization	Benchmark collection

Zhao et al.: <u>A Survey of Large Language Models</u>. 2024. arXiv:2303.18223

# **Where to Find Additional Information**



- Check out the LLMs/solutions to your problem that other people have tried.
  - by looking at leaderboards of the relevant benchmarks
  - by investigating the state-of-the-art for your your task on Papers with Code
  - or search for relevant scientific papers using Google Scholar
- Use them for inspiration and/or comparison

[IIII] Papers With Code

# **State-of-the-Art for Specific Tasks**



[111]	Search	Q Browse State-of-the	e-Art Datasets Methods	More∽	🛩 🕂 👷
		Brows 10,843 ber	e State-of Ichimarks 4,058 tasks 91,721	-the-Art papers with code	
	Computer Vision	51776 Incore		The Contraction	
	Segmentation Let 211 benchmarks 3697 papers with code	Le 412 benchmarks 2945 papers with code	2776 papers with code	Learning	Imil     Search     Q     Browse State of the Art     Datasets     More ~     Imiliar Analysis
	<ul> <li>See all 1460 tasks</li> <li>Natural Language</li> </ul>	Processing			Sentiment Analysis on Amazon Review Polarity Leaderboard Dataset
	Language Modelling Let 61 benchmarks 2493 papers with code	Let 192 benchmarks 1951 papers with code	Machine Translation Let 90 benchmarks 1782 papers with code	Sentiment Analysis Lez 88 benchmarks 1073 papers with code	Ven Accuracy V of Late V
	> See all 681 tasks				90 92
htt	tps://pap	perswith	ncode.co	om/sota	56 2016 2017 2018 2019 ■ Other models → Models with highest Accuracy Filter: Inntegrad
					Rank     Model     Accuracy + Paper     Code     Result     Year     Tags gr       1     BERT large     97.37     Unsupervised Data Augmentation for Consistency Training     •     •     2019       2     DPCNN     96.68     Categorization     •     •     •     •     2017

4 DRNN

-2018

Disconnected Recurrent Neural Networks for Text

96.49

# **Suggestions for APIs/Tools**



- For finding relevant APIs to query: https://apislist.com/
- Tavily search tool: <u>https://tavily.com/</u>
  - 1000 free requests per month (\*5 team members =  $\bigcirc$ )
  - Also possible to just use search engine
  - then find and extract relevant info from HTML
- Unstructured data transformation tool: <u>https://unstructured.io/</u>
  - Offer many open-source libraries to prepare data for LLMs you can use for free (<u>https://github.com/Unstructured-IO</u>)
- Groq LLM inference API: <u>https://groq.com/</u>
  - Let's you query some powerful open-source models for free (Llama3.1-70B)
  - May run into rate limit or unavailability due to traffic
- OpenAI: gpt4o-mini (or similar Google/Anthropic/Mistral offering)
  - Quite powerful and successfully used in team project for agentic workflow
  - Comparably cheap

University of Mannheim | IE686 LLMs and Agents | Intro to Student Projects | Version 23.10.2024

## **GPT-Researcher**



- Uses "Plan-and-Solve" prompting to divide task into subtasks...
- Which are carried out by multiple agents in parallel using web crawling as a tool.
- Written in LangGraph
- Can be used as inspiration





https://docs.gptr.dev/blog/building-gpt-researcher https://github.com/assafelovic/gpt-researcher/tree/master/multi\_agents

# **Tracking your Experiments**





- Print to console and save to logfiles, or...
- Use <u>LangSmith</u>
  - Commercial offering by LangChain
  - Let's you track and save what's happening during your workflow executions (traces)
  - Easy way to review what's happening and where things go wrong
  - 5000 free traces per month (\*5... 📀)

University of Mannheim | IE686 LLMs and Agents | Intro to Student Projects | Version 23.10.2024

# **Project Outlines**



- Maximum 4 pages (sharp!) including title page
  - Using <u>DWS Seminar thesis layout</u> (PDF!)
  - Include a project name, your team number and name on the first page!
- Due Monday, October 28th, 23:59
- Send by eMail to Ralph
- Feedback about your project outlines if required: Thursday, 31.10.2024, lecture time (15:15-17:00)
  - I will inform you Wednesday, 30.10.2024 with some feedback via mail and let you know if you need to attend the session on Thursday, 31.10.2024
  - Then the projects start!

# **Project Outlines**



- Answer the following questions:
  - 1. What is the problem you are solving?
  - 2. What data will you use?
    - Where will you get it?
    - How will you gather it?
  - 3. How will you solve the problem?
    - What LLMs do you plan to use?
    - Which methods do you plan to apply? Be as specific as you can!
    - What is your idea for a multi-agent workflow for your task?
  - 4. How will you measure success? (Evaluation method)

# **Coaching Sessions**



- I will give you tips and answer questions concerning your project.
- At the time of the lecture (Thursdays 15:15-17:00)
- **Registration via email** is mandatory if you want coaching!
  - Via mail to Ralph
  - Until Monday night (23:59) of the respective week!
  - Including the questions that you would like to discuss (be as clear and specific as possible so I can prepare)
- I will assign you a time slot afterwards and inform you about the slot via email

## **Some Project Management Hints**



- Organize your project in **multiple iterations** 
  - Every artefact will be improved over time!
- Get a **simple process running early** on to have a baseline
- **Parallelize tasks** while keeping centrally track of results
  - e.g. one central document with results plus reference to exact version of the notebooks/datasets that produced these results
  - sub-groups should explore specific ideas for a specified amount of time

# **Some Project Management Hints**



- **Define concrete milestones**: When should what be finished?
  - e.g. 08.11.24 Data exploration done and first simple baseline implemented
  - e.g. 15.11.24 Subgroup using decentralized agent communication adds results to central document

#### • Infrastructure

- use shared folder for result document, versions of data, processes, slideset (e.g. MS Teams, Google Drive, GitHub)
- use LLMs for inspiration about additional methods as well as for coding support

## **Tasks within the Iterations of the Project**



- 1. Data Exploration
- 2. Establish/update baseline (simple LLM agent(s) based method or non LLM baseline)
- 3. Try different methods for applying agents using different LLMs
  - Iteratively improve on methods...
  - solving problems as they come up
  - Track your experiments and all the things you did!
- 4. Perform error analysis in order to understand what is going on!
  - Stream the outputs, see whats going on and going wrong!
  - Analyse mistakes and reason about why system fails

## **Project Presentation**



- Present the project results to the other students
  - 12 minutes presentation + 3 minutes discussion
  - During lecture slot
  - Everyone must attend as its part of your grade
- Send your presentation in **PDF format** 
  - Via mail to Ralph and Prof. Bizer
  - Until Thursday, December 5th, 23:59

# **Project Report**



- Max. 12 pages including title/toc page and reference page
  - max. 10 pages content, no appendix
  - Each extra page and each day of late submission downgrades your mark by 0.3!
- due Sunday, December 1st, 23:59
- send by email to Ralph and Prof. Bizer

# **Project Report**



- Outline for project report:
  - Application area and goals (0.5 pages)
  - Profile (structure and size) of your data (minimum 1 page)
  - Approaches to solving the problem with LLM Agents
    - Describe different approaches that you tried
  - Evaluation
    - Results
    - Including problems that you faced with each approach and what steps you took to fix them
    - Including an analysis of the errors still made, a discussion of the results, and a comparison to state-of-the-art results (together: minimum 2 pages)

# **Project Report**



- Requirements
  - You have to use the latex template for <u>DWS Seminar Theses</u>
  - Please cite sources properly and use your references page
  - Also submit your Python code and (a subset) of your data
  - Include your names and your team number on the first page!
- Usage of additional AI Tools for non-method specific tasks needs to be declared in a separate table

Declaration of Used AI Tools						
Tool	Purpose	Where?	Useful?			
ChatGPT	Rephrasing	Throughout	+			
$\mathrm{DeepL}$	Translation	Throughout	+			
ResearchGPT	Summarization of related work	Sec. 2.2	-			
Dall-E	Image generation	Figs. 2, 3	++			
GPT-4	Code generation	functions.py	+			
ChatGPT	Related work hallucination	Most of bibliography	++			

# Get Additional Advice from a Stanford Professor

- How to evaluate your model?
  - <u>https://www.youtube.com/watch?v=TxTblROT9lY</u>
- How to structure your project report?
  - <u>https://www.youtube.com/watch?v=DZNwO-p5PGY</u>
- How to present the results of your project?
  - <u>https://www.youtube.com/watch?v=GGx7klcahzY</u>







**Christopher Potts** 

## **Deadlines - Overview**



- Project outline until Monday, October 28th, 23:59
- Coaching Sessions
  - Every week
  - Registration for coaching must be done by Monday 23:59 of that week the latest!
- Project report until **Sunday, December 1st, 23:59**
- Project presentation as PDF until **Thursday, December 5th, 23:59**

# **Questions?**





# Thank you



