

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
 - 10 minutes presentation + 5 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 general registration period starts on **02.04.2025** and continues until **16.04.2025**
- Make sure to **register!**

Course Schedule

You are here

Day	Topic
20.03	Project: Introduction to Student Projects
27.03	Exercise: Introduction to LangGraph 2
03.04	Exercise: Introduction to AutoGen
Sunday, April 6th, 23:59 Submission of Project Outlines	
10.04	Project: Project Coaching
30.04	Project: Project Coaching
08.05	Project: Project Coaching
15.05	Project: Project Coaching
22.05	Project: Project Coaching
Sunday, May 25th, 23:59 Submission of Project Report	
28.05	Project: Presentation of Project Results
Wednesday, May 28th, 23:59 Submission of Presentation Slides	

Team Formation

- Team formation is already done based on your preferences
 - You have been assigned to your team in ILIAS with your team number and assigned topic
 - It was difficult to satisfy everyone's wishes at once, so receiving the top priority topic area was not always possible
- 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: 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:**
 - <https://github.com/THUDM/AgentBench>
 - <https://github.com/microsoft/SmartPlay>
- **Evaluation:**
 - Existing evaluation when using agent benchmark
 - Otherwise depends on game, could be win rate, ELO ranking, ...

Topic 2: 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
 - <https://github.com/THUDM/AgentBench>
 - <https://webarena.dev/>
 - <https://webshop-pnlp.github.io/>
- **Evaluation:**
 - Relevant evaluation metrics in used benchmarks

Topic 3: 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 4: Text-to-BPMN

- **Task:** Build a (Multi-)Agent system that can create a business process model given a natural language description
- **Dataset/APIs:**
 - Use existing benchmarks and relevant metrics
 - Select a representative set of description-process pairs
- **Evaluation:**
 - Relevant evaluation metrics in the field
 - Human Evaluation

Topic 5: Browser Interaction

- **Task:** Build a (Multi-)Agent system that can perform various tasks on websites.
 - Site search/link following/extraction of relevant info
 - Identify problems with each approach
- **Dataset/APIs:**
 - Use existing benchmarks and relevant metrics
 - Or define your own tasks in the real web
 - <https://github.com/THUDM/AgentBench>
 - <https://webarena.dev/>
- **Evaluation:**
 - Relevant evaluation metrics in used benchmarks
 - Measure task fulfillment

Where to find Datasets/Benchmarks?

- Agent Benchmarks
 - [AgentBench](#)
 - [WebShop](#)
 - [WebArena](#)
 - [Meta RAG KDD Cup](#)
 - [SmartPlay](#)
 - [EvalPlus](#)
- General Task Benchmarks
 - [Papers with Code Datasets](#) (filter by task)
 - [Huggingface Datasets](#)

Even More Benchmarks

Method	Evaluation	Model Types	Abilities/Domain	Data Source
Benchmark	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
	APIBench [717]	Base/Fine-tuned	Tool manipulation	Web
	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
	SciBench [728]	Fine-tuned	Complex reasoning	Human exam/practice
Model	AlpacaEval [729]	Fine-tuned	Instruction following	Synthesis
	MT-bench [727]	Fine-tuned	Human alignment	Human annotation
	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

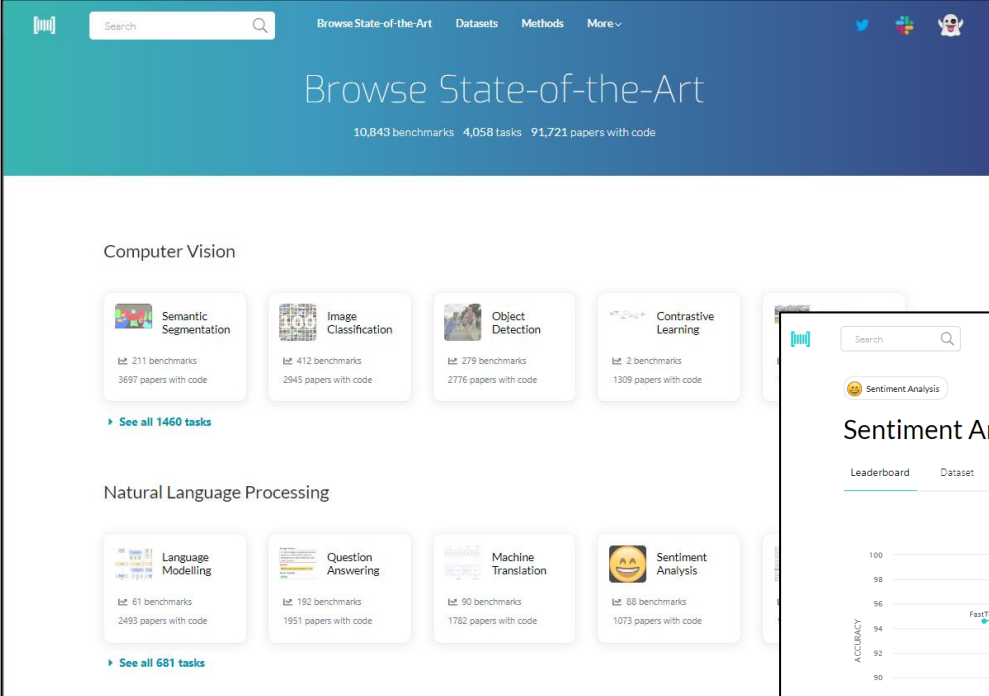
Zhao et al.: A Survey of Large Language Models. 2024. arXiv:2303.18223

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



State-of-the-Art for Specific Tasks



Browse State-of-the-Art
10,843 benchmarks 4,058 tasks 91,721 papers with code

Computer Vision

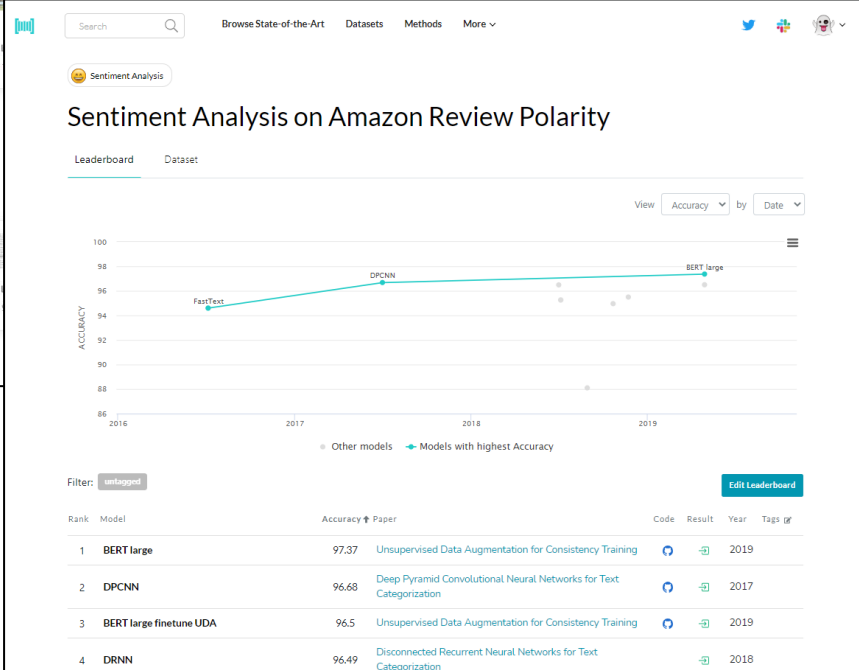
- Semantic Segmentation**: 211 benchmarks, 3697 papers with code
- Image Classification**: 412 benchmarks, 2945 papers with code
- Object Detection**: 279 benchmarks, 2776 papers with code
- Contrastive Learning**: 2 benchmarks, 1309 papers with code

▶ See all 1460 tasks

Natural Language Processing

- Language Modelling**: 61 benchmarks, 2493 papers with code
- Question Answering**: 192 benchmarks, 1951 papers with code
- Machine Translation**: 90 benchmarks, 1782 papers with code
- Sentiment Analysis**: 88 benchmarks, 1073 papers with code

▶ See all 681 tasks



Sentiment Analysis on Amazon Review Polarity

Leaderboard Dataset

View Accuracy by Date

ACCURACY

2016 2017 2018 2019

FastText DPCNN BERT large

Other models Models with highest Accuracy

Filter: sentiment Edit Leaderboard

Rank	Model	Accuracy	Paper	Code	Result	Year	Tags
1	BERT large	97.37	Unsupervised Data Augmentation for Consistency Training	🔗	📄	2019	
2	DPCNN	96.68	Deep Pyramid Convolutional Neural Networks for Text Categorization	🔗	📄	2017	
3	BERT large finetune UDA	96.5	Unsupervised Data Augmentation for Consistency Training	🔗	📄	2019	
4	DRNN	96.49	Disconnected Recurrent Neural Networks for Text Categorization	🔗	📄	2018	

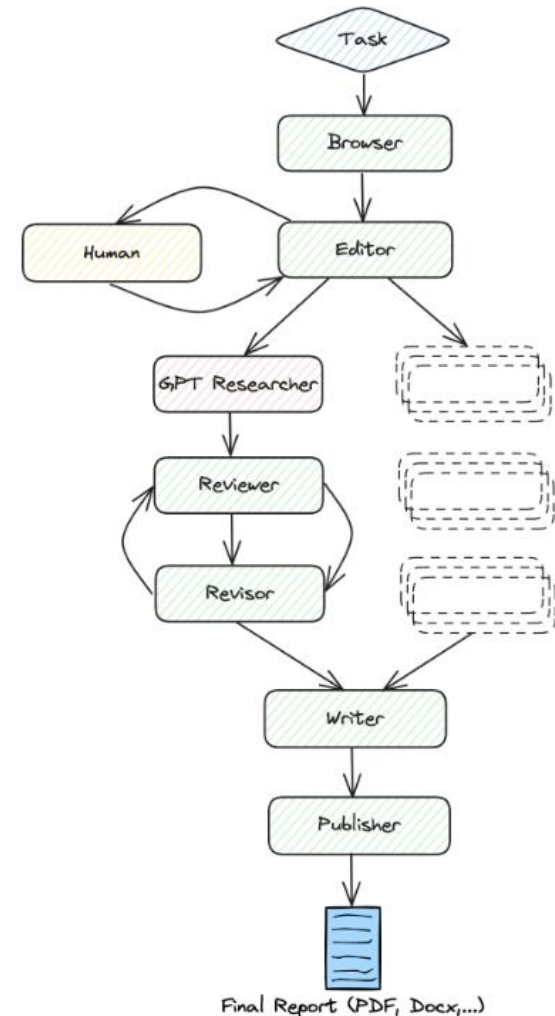
<https://paperswithcode.com/sota>

Suggestions for APIs/Tools

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

GPT-Researcher

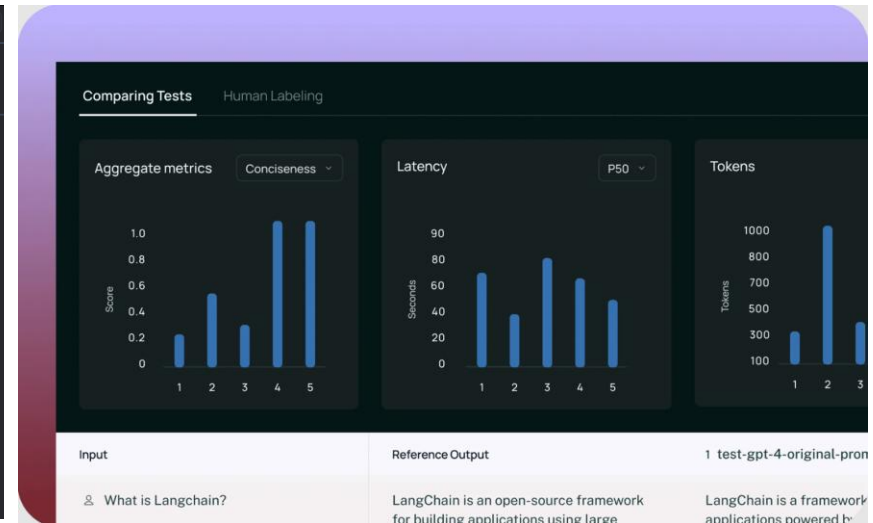
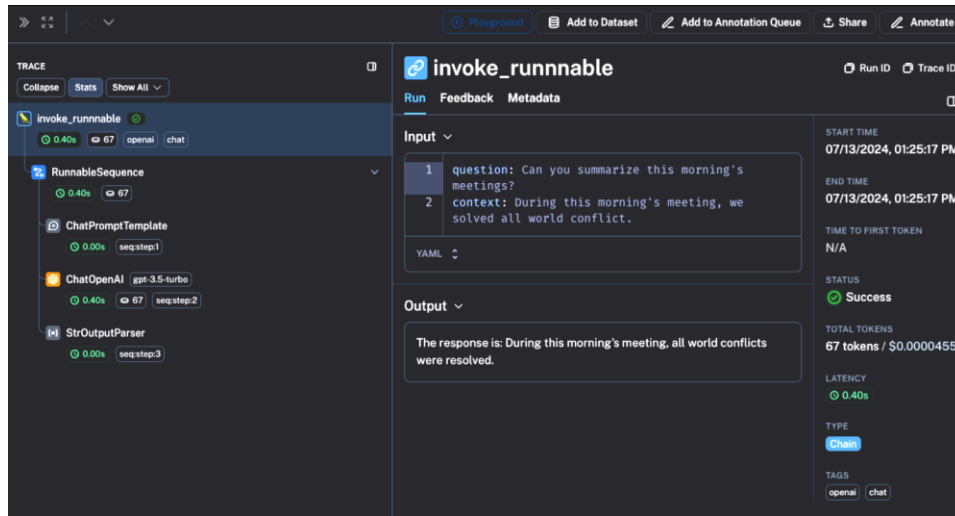
- Multi-agent system for online research
- Uses “Plan-and-Solve” prompting to divide task into subtasks...
- Which are carried out by multiple agents in parallel using web crawling as a tool.
- Written in LangGraph
- Can be used as inspiration



<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 LangSmith
 - 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... 😊)

Project Outlines

- Maximum 4 pages (sharp!) including title page
 - Using DWS Seminar thesis layout (PDF!)
 - Include a project name, your team number and name on the first page!
- Due **Sunday, April 6th, 23:59**
- Upload to ILIAS submission (Every team member can submit, last submission will be evaluated)
- Feedback about your project outlines if required:
Thursday, 10.04.2025, lecture time (15:15-17:00)
 - We will inform you Wednesday, 09.04.2025 with some feedback via mail and let you know if you need to attend the session on Thursday
 - Then the projects start!

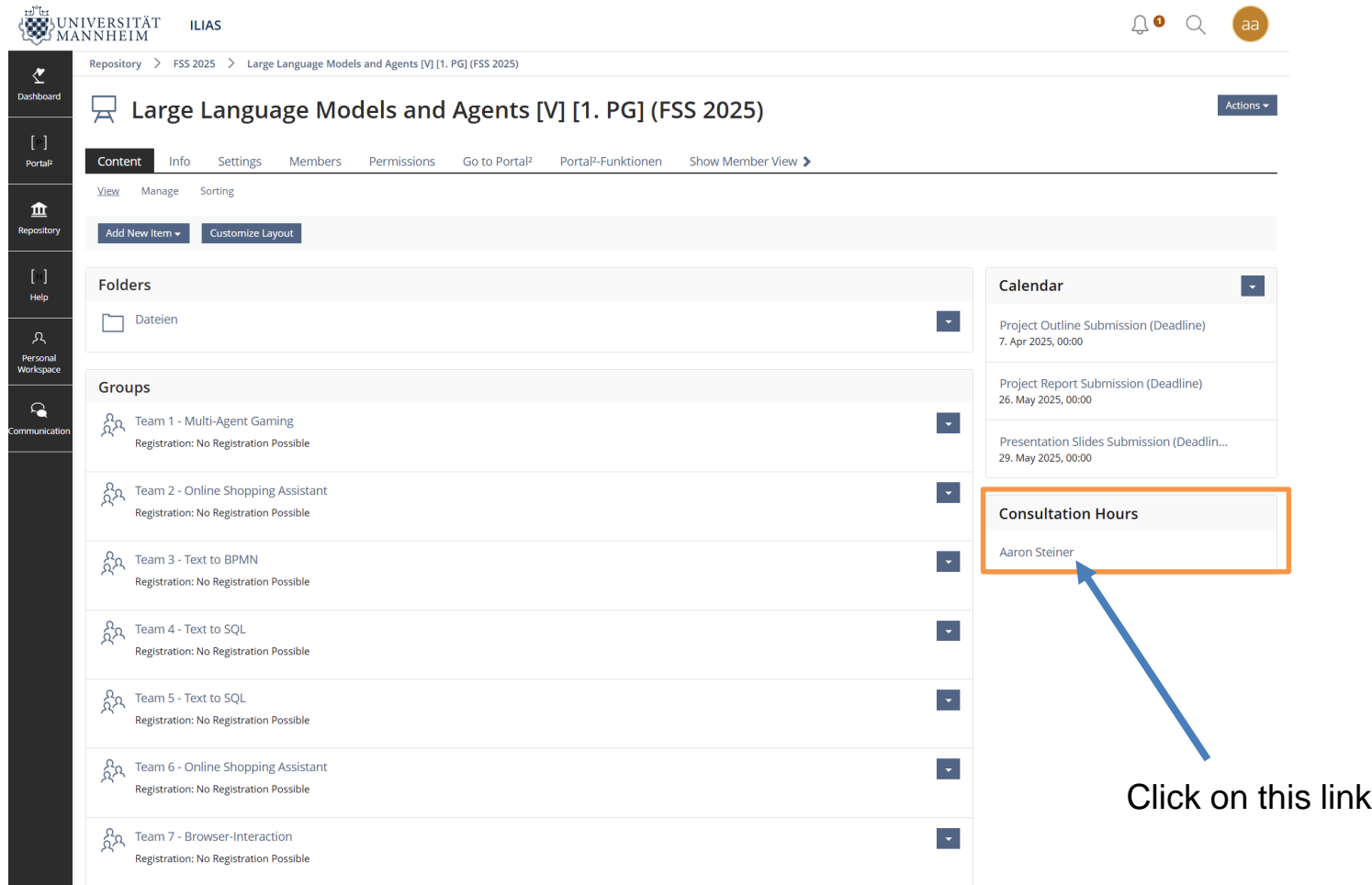
Project Outlines

- 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

- We give you tips and answer questions about your project.
- At the time of the lecture (Thursdays 15:15-17:00)
- **Registration** is mandatory if you want coaching!
- Make sure to register until Monday (23:59) of the week you want to attend the coaching session
- Each coaching session lasts for 15 minutes
 - Include your questions when booking the session, so we can prepare!
 - Most time efficient use of the session
 - We will of course also answer any question you pose directly in the session

Booking Coaching Sessions: How-to



UNIVERSITÄT MANNHEIM ILIAS

Repository > FSS 2025 > Large Language Models and Agents [V] [1. PG] (FSS 2025)

Large Language Models and Agents [V] [1. PG] (FSS 2025)

Content Info Settings Members Permissions Go to Portal² Portal²-Funktionen Show Member View >

View Manage Sorting

Add New Item Customize Layout

Folders

- Dateien

Groups

- Team 1 - Multi-Agent Gaming
Registration: No Registration Possible
- Team 2 - Online Shopping Assistant
Registration: No Registration Possible
- Team 3 - Text to BPMN
Registration: No Registration Possible
- Team 4 - Text to SQL
Registration: No Registration Possible
- Team 5 - Text to SQL
Registration: No Registration Possible
- Team 6 - Online Shopping Assistant
Registration: No Registration Possible
- Team 7 - Browser-Interaction
Registration: No Registration Possible

Calendar

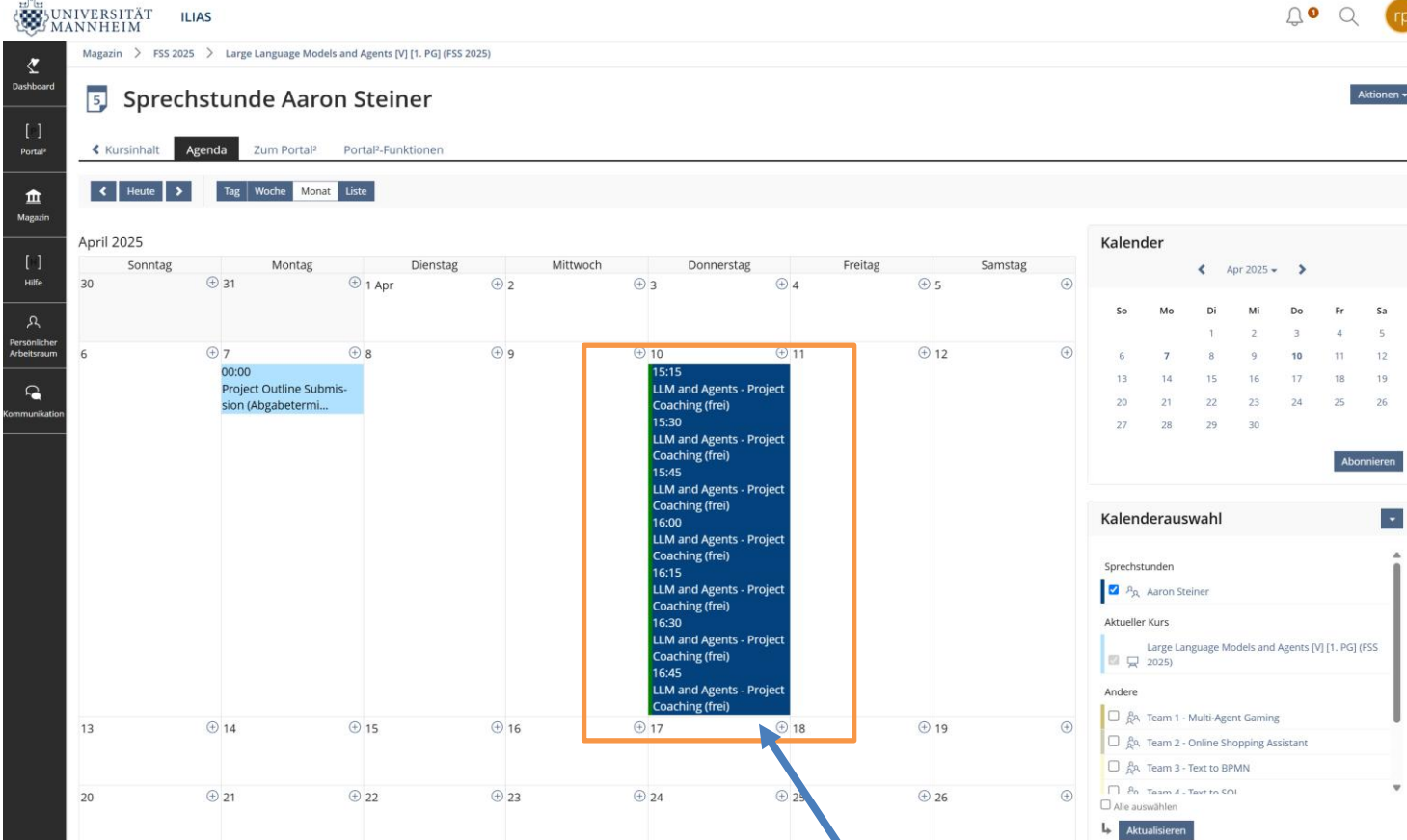
- Project Outline Submission (Deadline)
7. Apr 2025, 00:00
- Project Report Submission (Deadline)
26. May 2025, 00:00
- Presentation Slides Submission (Deadlin...
29. May 2025, 00:00

Consultation Hours

- Aaron Steiner

Click on this link

Booking Coaching Sessions: How-to



Magazin > FSS 2025 > Large Language Models and Agents [V] [1. PG] (FSS 2025)

Sprechstunde Aaron Steiner

Kursinhalt **Agenda** Zum PortalP PortalF-Funktionen

Heute Tag Woche Monat Liste

April 2025

Sonntag	Montag	Dienstag	Mittwoch	Donnerstag	Freitag	Samstag
30	31	1 Apr	2	3	4	5
6	7 00:00 Project Outline Submission (Abgabetermi...	8	9	10 15:15 LLM and Agents - Project Coaching (fre) 15:30 LLM and Agents - Project Coaching (fre) 15:45 LLM and Agents - Project Coaching (fre) 16:00 LLM and Agents - Project Coaching (fre) 16:15 LLM and Agents - Project Coaching (fre) 16:30 LLM and Agents - Project Coaching (fre) 16:45 LLM and Agents - Project Coaching (fre)	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26

Kalender

Apr 2025

So	Mo	Di	Mi	Do	Fr	Sa
	1	2	3	4	5	
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21	22	23	24	25	26
27	28	29	30			

Kalenderauswahl

Sprechstunden

Aaron Steiner

Aktueller Kurs

Large Language Models and Agents [V] [1. PG] (FSS 2025)

Andere

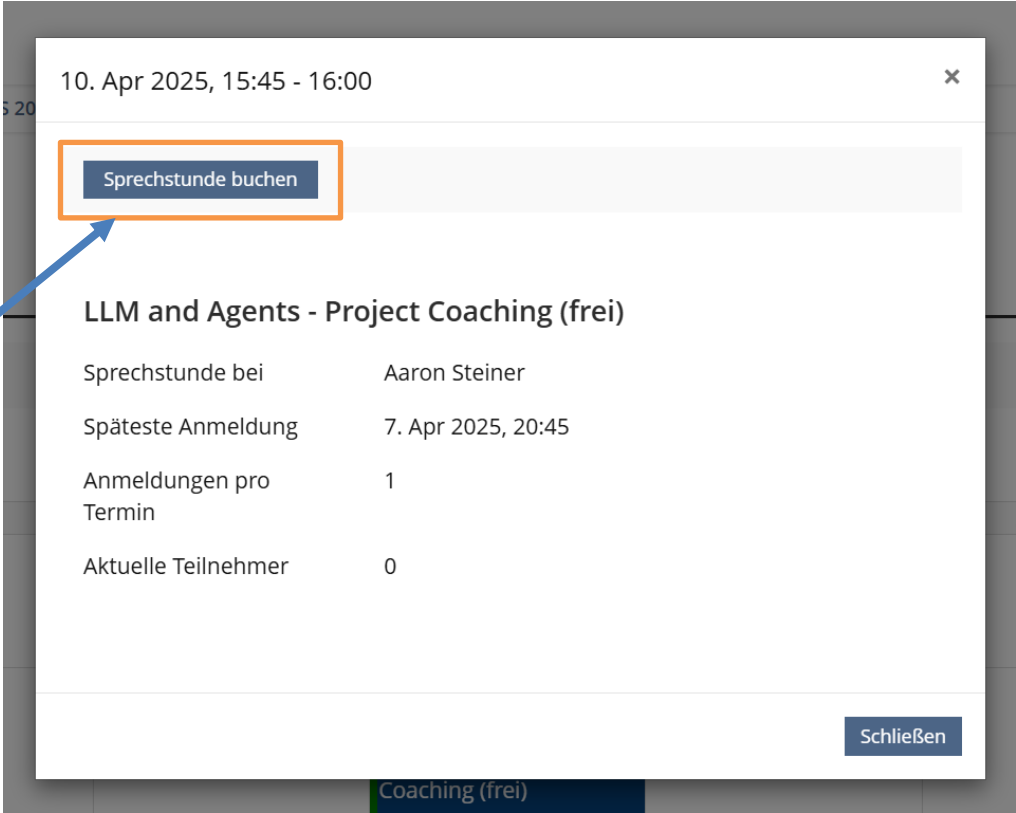
- Team 1 - Multi-Agent Gaming
- Team 2 - Online Shopping Assistant
- Team 3 - Text to BPMN
- Team 4 - Text to SQL

Alle auswählen

Aktualisieren

Choose and click the slot you want to attend (as long as it is still free)

Booking Coaching Sessions: How-to



10. Apr 2025, 15:45 - 16:00

Sprechstunde buchen

LLM and Agents - Project Coaching (frei)

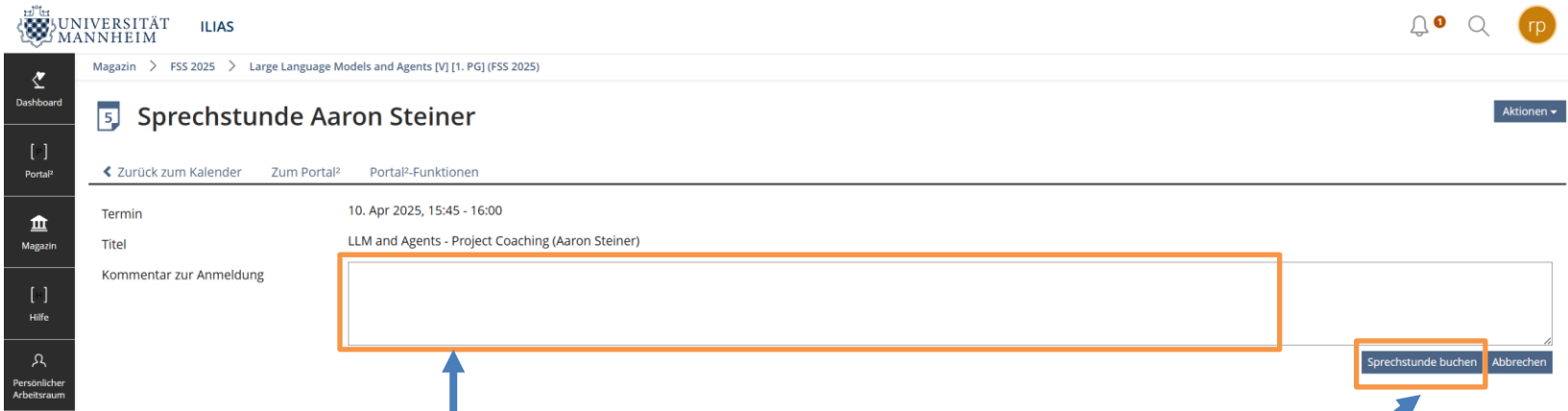
Sprechstunde bei	Aaron Steiner
Späteste Anmeldung	7. Apr 2025, 20:45
Anmeldungen pro Termin	1
Aktuelle Teilnehmer	0

Schließen

Coaching (frei)

Click to book

Booking Coaching Sessions: How-to



The screenshot shows the ILIAS interface for booking a coaching session. The page title is "Sprechstunde Aaron Steiner". The session details are: Termin: 10. Apr 2025, 15:45 - 16:00; Titel: LLM and Agents - Project Coaching (Aaron Steiner). Below the title is a large text input field for "Kommentar zur Anmeldung", which is highlighted with an orange border. To the right of this field are two buttons: "Sprechstunde buchen" (highlighted with an orange border) and "Abbrechen". A blue arrow points from the text "1. Enter your questions and topics to discuss during the session" to the input field. Another blue arrow points from the text "2. Finalize the booking of the coaching session" to the "Sprechstunde buchen" button.

1. Enter your questions and topics to discuss during the session

2. Finalize the booking of the coaching session

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. 20.04.25 Data exploration done and first simple baseline implemented
 - e.g. 26.04.25 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
 - 10 minutes presentation + 5 minutes discussion
 - Everyone must attend as its part of your grade
- Upload your presentation in **PDF format**
 - Via your ILIAS Group
 - Until **Wednesday, May 28th, 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, May 25th, 23:59**
- Upload in **PDF format** via your ILIAS group

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 DWS Seminar Theses
 - Please cite sources properly and use your references page
 - Also submit your Python code and (a subset) of your data
 - Also submit a **who-did-what table** together with your report/code!
 - 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	+
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?
 - <https://www.youtube.com/watch?v=TxTbIROT9IY>
- How to structure your project report?
 - <https://www.youtube.com/watch?v=DZNwO-p5PGY>
- How to present the results of your project?
 - <https://www.youtube.com/watch?v=GGx7klcahZ>



Christopher Potts

Deadlines - Overview

- Project outline until **Sunday, April 6th, 23:59**
- Coaching Sessions
 - Every week
 - Registration for coaching must be done by Monday 23:59 of that week the latest!
- Project report as PDF until **Sunday, May 25th, 23:59**
- Project presentation as PDF until **Wednesday, May 28th, 23:59**

Questions?



Thank you

