# Columbia Business School - B9153
# Generative AI: Technical and Social

## Class Schedule

Fall 2024
Thursdays 2:20 pm – 5:20 pm
Location: Geffen 430
Credits: 3 Credits

## Instructors

Professors
Hannah Li: `hannah.li@columbia.edu`
Tianyi Peng: `tianyi.peng@columbia.edu`

Teaching Assistant
Haozhe Chen: `haozhe.chen@columbia.edu`

## Course Description

This course examines generative AI technologies through both a technical and social lens.

We will introduce the foundational architecture of generative AI and large language models (LLMs) and simultaneously investigate the deployment of these models in various sectors.

Designed for students interested in advancing AI technology responsibly, this course encourages critical thinking about AI's broader effects. Participants will gain practical skills and a deeper understanding of how AI tools can be developed and utilized ethically and effectively.

On the technical side, students will develop hands-on experience with the technical workings of LLMs. This will range from fundamental aspects such as understanding transformers and training a small language model from scratch to building a chatbot by leveraging prompt engineering, retrieval-augmented generation, and fine-tuning to enhance usability and reliability.

On the social side, we will examine the implications of these technologies on applications in education, labor markets, online platforms, and data science. Students will identify the various use cases of AI in human-augmenting, human-replacing, and human-understanding roles. Additionally, students will identify the risks and benefits, as well as avenues of improvement.

This course is oriented toward identifying open research questions in the area of generative AI, particularly at the intersection with statistics, economics, and operations research (OR).

# Prerequisites

This class is open to PhD students across Columbia with some prior background in operations research, machine learning, statistics, or economics. An introductory level of coding experience with Python (or a similar programming language, such as C/C++, Matlab, Java) is recommended. A key milestone in this class will be coding up a transformer model from scratch, but we will provide a step by step walk through on how to do this over several lectures. We do not presume existing knowledge for coding up LLMs.

The most important prerequisite is to "Stay hungry and stay foolish," as we explore topics together in the new domain of generative AI.

# Course Schedule (Subject to Change)

Note: Of the readings listed, some readings will be required and others will be recommended. In the first week of class, we will solicit suggestions for additional papers that participants wish to discuss. We will consider these suggestions before finalizing the list of required readings.

**Sep 5, 2024 Topic:** Introduction to generative AI and deep learning

> **Readings:** Donoho. "Data science at the singularity." (2024).

**Sep 12, 2024 Topic:** Text classification; Productivity impacts of generative AI tools (part 1)

> **Readings:** Brynjolfsson, Li, and Raymond. "Generative AI at Work." (2023).
>
> Noy and Zhang. "Experimental evidence on the productivity effects of generative artificial intelligence." (2023).

**Sep 19, 2024 Topic:** Coding up a transformer model; Impacts of generative AI on learning and education (part 1)

> **Required Readings:** Bastani, Bastani, Sungu, Ge, Kabakci, and Mariman. "Generative AI Can Harm Learning."(2024).
>
> Sadasivan, Kuar, Balasubramanian, Wang, Feizi. "Can AI-Generated Text be Reliably Detected?" (2023).
>
> **Recommended Readings:**
>
> Mao, Vondrick, Wang, and Yang. "RAIDAR. GeneRative AI Detection viA Rewriting." (2024).
>
> Vaswani, Shazeer., Parmar., Uszkoreit, Jones, Gomez, Kaiser, and Polosukhin. "Attention is all you need." (2017)
>
> Devlin, Chang, Lee, Toutanova. "Bert: Pre-training of deep bidirectional transformers for language understanding." (2018). (Bert)
>
> Alec, Wu, Child, Luan, Amodei, and Sutskever. "Language models are unsupervised multitask learners." (2019). (GPT2)

**Sep 26, 2024 Topic:** Coding up a transformer model; Impacts of generative AI on learning and education (part 2)

> **Required Readings:** N/A
>
> **Recommended Readings:**
>
> Lee, Hicke, Yu, Brooks, and Kizilcec. "The life cycle of large language models in education: A framework for understanding sources of bias." (2024).
>
> AlphaProof and AlphaGeometry teams. "AI achieves silver-medal standard solving International Mathematical Olympiad problems." (2024).
>
> Brown, Mann, Ryder, Subbiah, et.al. "Language models are few-shot learners." (GPT3)

**Oct 3, 2024 Topic:** Chatbot; LLMs for science and experimentation

    **Required Readings:** N/A

    **Recommended Readings:**

    Lu, Lu, Lange, Foerster, Clune, and Ha. "The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery." (2024). (LLM as AI scientist)

    Huang, Vora, Liang, and Leskovec. "MLAgentBench: Evaluating Language Agents on Machine Learning Experimentation." (2024). (LLM as data scientist)

**Oct 10, 2024 Topic:** Prompt engineering and Retrieval Augmentation Generation; LLM for science and experimentation (part 2)

    **Required Readings:**

    Manning, Zhu, and Horton. "Automated Social Science: Language Models as Scientist and Subjects." (2024). (LLMs as social scientists)

    Romera-Paredes, Barekatain, Novikov, Balog, Kumar, Dupont, Ruiz, Ellenberg, Wang, Fawzi, Kohli, and Fawzi. "Mathematical discoveries from program search with large language models." (2024). (LLMs as mathematicians)

    **Recommended Readings:**

    Wei, Wang, Schuurmans, Bosma, Xia, Chi, Le, and Zhou. "Chain-of-thought prompting elicits reasoning in large language models." (2022).

    Anthropic's Prompt Engineering Interactive Tutorial: `https://github.com/anthropics/courses/tree/master/prompt_engineering_interactive_tutorial`

    Lewis, Perez, Piktus, Petroni, Karpukhin, Goyal, Küttler et al. "Retrieval-augmented generation for knowledge-intensive nlp tasks." (2020).

**Oct 17, 2024 Topic:** No Class

**Oct 24, 2024 Topic:** Retrieval Augmented Generation and Reinforcement Learning with Human Feedback; Chatbot applications and customer interactions

    **Required Readings:**

    Costello, Pennycook, and Rand. "Durably reducing conspiracy beliefs through dialogues with AI." (2024).

    Haque and Rubya. "An Overview of Chatbot-Based Mobile Mental Health Apps: Insights From App Description and User Reviews." (2023).

    **Recommended Readings:**

    Yao, Zhao, Yu, Du, Shafran, Narasimhan, Cao. "React: Synergizing reasoning and acting in language models." (2022).

    Ouyang, Wu, Jiang, Almeida, Wainwright, Mishkin, Zhang et al. "Training language models to follow instructions with human feedback." (2022). (RLHF)

**Oct 31, 2024 Topic:** Reinforcement Learning with Human Feedback; Chatbot applications and customer interactions

    **Required Readings:**

    Sanner, Balog, Radlinski, Wedin, and Dixon. "LLMs are Competitive Near Cold-start Recommenders for Language- and Item-based Preferences." (2023).

    Rio-Chanona, Laurentsyeva, Wachs. "Are Large Language Models a Threat to Digital Public Goods? Evidence from Activity on Stack Overflow." (2023).

    **Recommended Readings:**

"How Amazon continues to improve the customer reviews experience with generative AI." (https://www.aboutamazon.com/news/amazon-ai/amazon-improves-customer-reviews-with-generative-ai)

Zhai, Liao, Liu, Wang, Li, Cao, Gao et al. "Actions speak louder than words: Trillion-parameter sequential transducers for generative recommendations." (2024). (LLM for recommendation)

**Nov 07, 2024 Topic:** Evaluation of LLM models; LLMs for Operations

> **Required Readings:** AhmadiTeshnizi, Gao, Brunborg, Talaei, and Udell. "OptiMUS-0.3: Using Large Language Models to Model and Solve Optimization Problems at Scale." (2024).
>
> Park, O'Brien, Cai, Morris, Liang, and Bernstein. "Generative agents: Interactive simulacra of human behavior." (2023).
>
> **Recommended Readings:**
>
> Tang, Huang, Zheng, Hu, Wang, Ge, and Wang. "ORLM: Training Large Language Models for Optimization Modeling." (2024).
>
> Li, Wang, Zhang, Li, Lai, Kang, Ma, and Liu. "Agent hospital: A simulacrum of hospital with evolvable medical agents." (2024).
>
> Huggingface LLM leaderboard `https://huggingface.co/open-llm-leaderboard`
>
> Chiang, Zheng, Sheng, Angelopoulos, Li, Li, Zhang et al. "Chatbot arena: An open platform for evaluating llms by human preference." (2024).

**Nov 14, 2024 Topic:** Operations for LLMs: resource allocation and serving optimization

> **Required Readings:**
>
> Chen, Zaharia, and Zou. "FrugalGPT: How to Use Large Language Models While Reducing Cost and Improving Performance." (2023).
>
> Agrawal, Kedia, Panwar, Mohan, Kwatra, Gulavani, Tumanov, and Ramjee. "Taming throughput-latency tradeoff in llm inference with sarathi-serve." (2024).
>
> **Recommended Readings:**
>
> Zhu, Sheng, Zheng, Barrett, Jordan, and Jiao. "Towards Optimal Caching and Model Selection for Large Model Inference." (2024).
>
> Li, Zheng, Zhong, Liu, Sheng, Jin, Huang et al. "AlpaServe: Statistical multiplexing with model parallelism for deep learning serving." (2023).
>
> Griggs, Liu, Yu, Kim, Chiang, Cheung, and Stoica. "Mélange: Cost Efficient Large Language Model Serving by Exploiting GPU Heterogeneity." (2024).

**Nov 21, 2024 Topic:** Productivity impacts of LLMs; human and AI interaction

> **Required Readings:** Otis, Clarke, Delecourt, Holtz, and Koning. "The Uneven Impact of Generative AI on Entrepreneurial Performance." (2024).

**Nov 28, 2024 Topic:** No Class (Thanksgiving Holiday)

**Dec 5, 2024 Topic:** Student project presentations

# Requirements

## Class participation - 20%

Active class participation is essential to this class, as we explore open questions and potential research problems in this domain together. Additionally, each student will be responsible for presenting a paper in class and guiding the discussion for this paper. Attendance is required.

## Homework assignments - 30%

There will be two homework assignments that involve implementing the core technical concepts from class.

## Final project - 50%

Students will complete a class project that delves into a research question in the field of Generative AI. Projects can include coding, empirical research, theoretical analysis, or modeling. We encourage exploring with enthusiasm, and reports of failures are welcomed as valuable learning experiences.

Deliverables for this project include a project proposal, an in-class presentation, and a final writeup. The final project writeup is expected to be 5-10 pages, excluding the references and appendix (if any). Students can work individually or in pairs.

The project proposal is due on November 1st. The In-class presentation is due on December 5th. The final project write-up is due on December 7th.