
Large Language Models for NLP in Finance

Simerjot Kaur, Charese Smiley, Xiaomo Liu, Elena Kochkina, Manling Li, Mohammad Ghassemi,
Reza Khanmohammadi

Meet the Team



Dr. Simerjot Kaur
AI Research Lead
JPMorgan Chase & Co



Dr. Charese Smiley AI
Research Lead
JPMorgan Chase & Co



Dr. Elena Kochkina
AI Research Scientist
JPMorgan Chase & Co



Dr. Xiaomo Liu
Executive Director of AI Research
JPMorgan Chase & Co



Dr. Manling Li
Assistant Professor
Northwestern University



Dr. Mohammad Ghassemi
Assistant Professor
Michigan State University



Reza Khanmohammadi
Ph.D. Student
Michigan State University

Presentation Outline

- **Introduction and Overview [15 mins talk; 5 mins Q&A]:**
Exploring the Evolution of LLMs and their applications in Finance.
Mohammad Ghassemi, Reza Khanmohammadi
- **Application Insights [20 mins; 5 mins Q&A]:**
Deep Dive into Financial Applications of LLMs.
Xiaomo Liu, Simerjot Kaur
- **Addressing Challenges [25 mins; 5 mins Q&A]:**
Analyzing the Limitations and Challenges in LLMs.
Manling Li

Introduction and Overview: Trends in LLMs generally, and Finance Specifically

- **Trends in LLM Research, broadly [8 minutes]:**

Movva, R., Balachandar, S., Peng, K., Agostini, G., Garg, N., & Pierson, E. (2023). *Topics, Authors, and Networks in Large Language Model Research: Trends from a Survey of 17K arXiv Papers*. arXiv:2307.10700 [cs.DL]

- **Trends in LLM Research for Finance, specifically [7 minutes]:**

Synthesis of the data, methods, and topics presented in 50 recent papers on LLMs in Finance; papers were sourced from Google Scholar and ArXiv using the search terms:

('Large Language Models' OR 'GPT') AND
('Finance') AND
 $(\emptyset \cup \{\text{Application}\})$.

Where {Applications} were manually created and included the following topics: *Investment and Trading Insights, Risk Mitigation and Fraud Prevention, Customer Engagement and Services, Market Analysis and Sentiment Tracking, Credit and Financial Health*

Introduction and Overview: Trends in LLMs generally, and Finance Specifically

- **Trends in LLM Research, broadly [8 minutes]:**

Movva, R., Balachandar, S., Peng, K., Agostini, G., Garg, N., & Pierson, E. (2023). *Topics, Authors, and Networks in Large Language Model Research: Trends from a Survey of 17K arXiv Papers*. arXiv:2307.10700 [cs.DL]

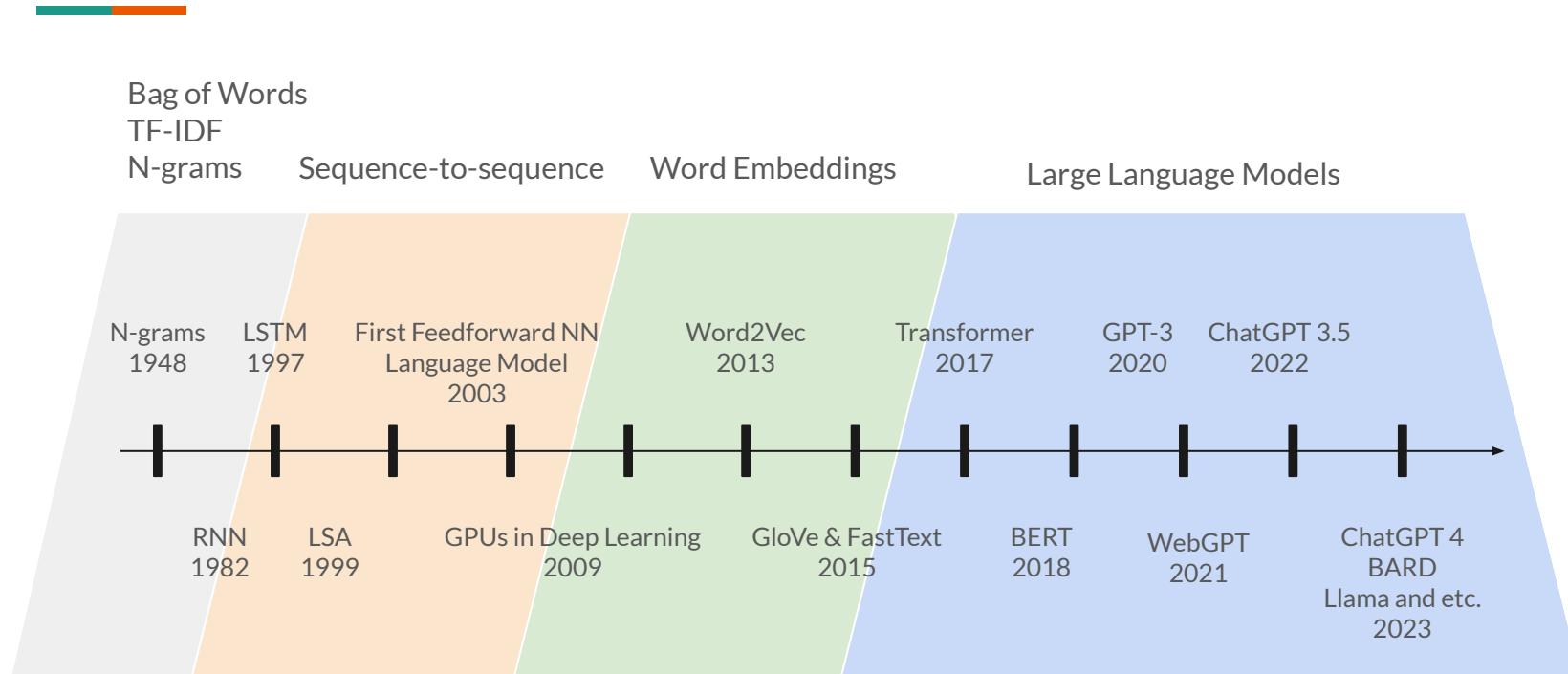
- **Trends in LLM Research for Finance, specifically [7 minutes]:**

Synthesis of the data, methods, and topics presented in 52 recent papers on LLMs in Finance; papers were sourced from Google Scholar and ArXiv using the search terms:

('Large Language Models' OR 'GPT') AND
('Finance') AND
 $(\emptyset \cup \{\text{Application}\})$.

Where {Applications} were manually created and included the following topics: *Investment and Trading Insights, Risk Mitigation and Fraud Prevention, Customer Engagement and Services, Market Analysis and Sentiment Tracking, Credit and Financial Health*

The Evolution of NLP: From Simple N-grams to Advanced LLMs



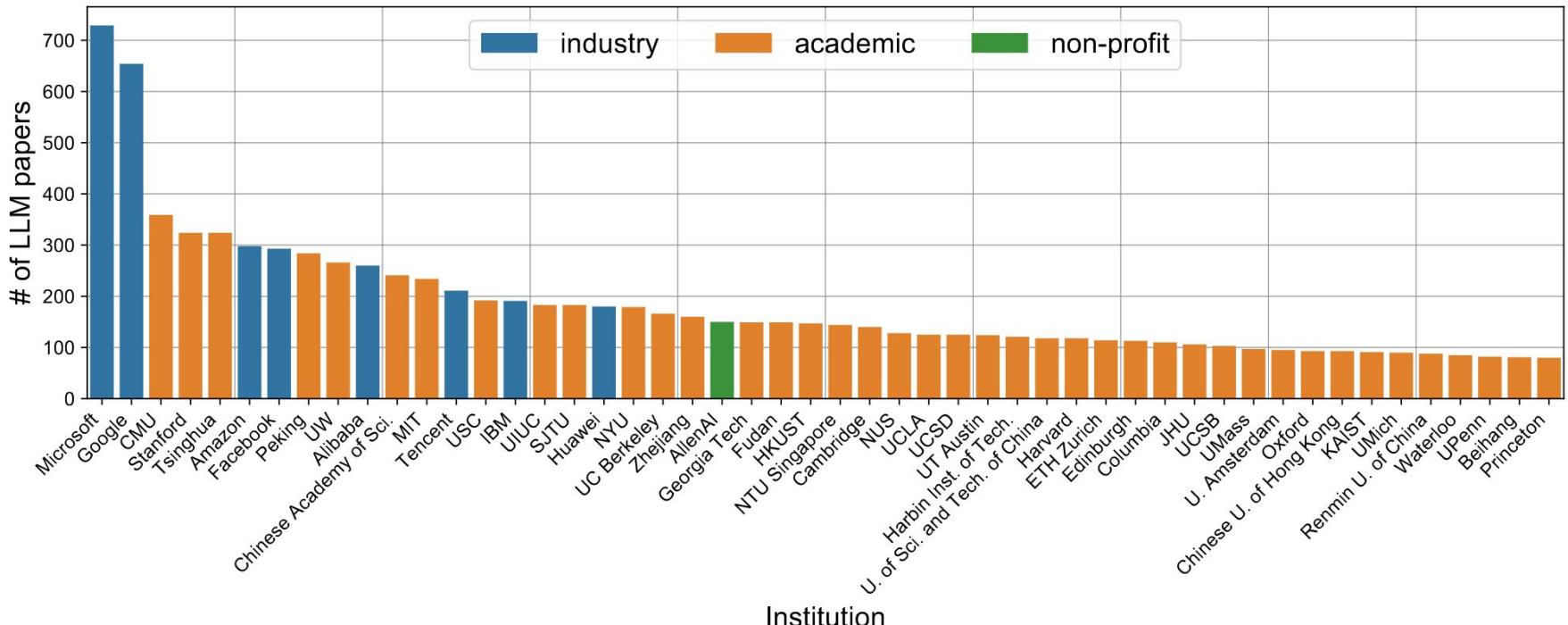
Trends in 2023 LLM Research, broadly: Societal Impact, New Authors and Perspectives

- **Social Impact:**
20x increase in research related to societal issues compared to 2022.
- **Fresh Voices:**
Most 2023 papers (61.4%) led by first-time authors.
- **Bigger Teams in LLM Research:**
Growth in median author count per paper from 4 to 5.

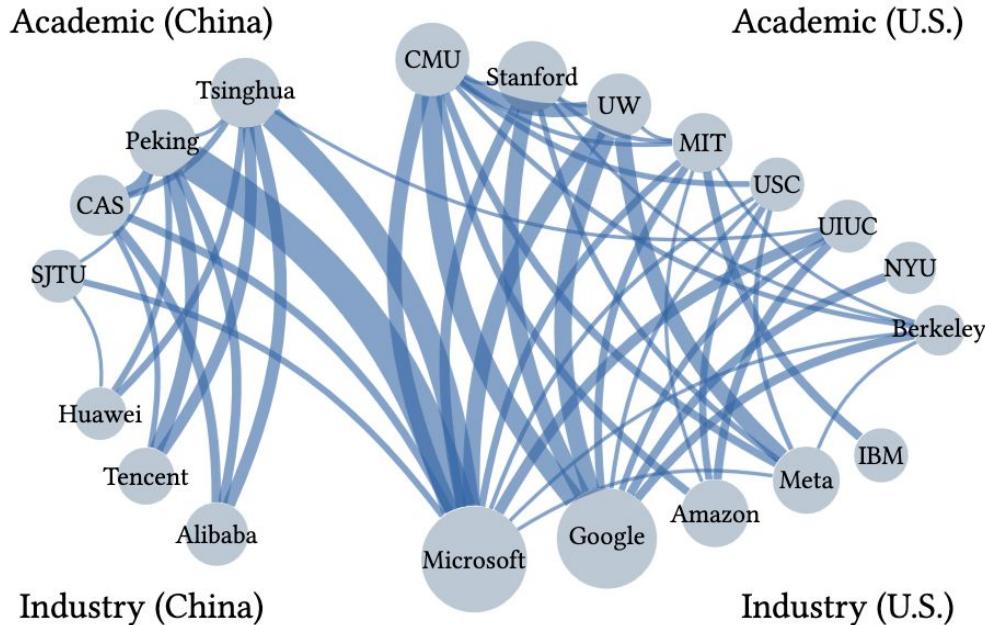
Teams of 10 or more authors has nearly doubled to 8.0%.

Solo-authored papers increased from 3.4% to 5.0%.
- **Applied Works Attracts Citations:**
Focus on ChatGPT and vision-language models yields high citation rates.

~2/3 of LLM research performed in academic settings: but MSFT & GOOG lead output



Significant Industry-Academic Partnerships: corporate collaboration more rare

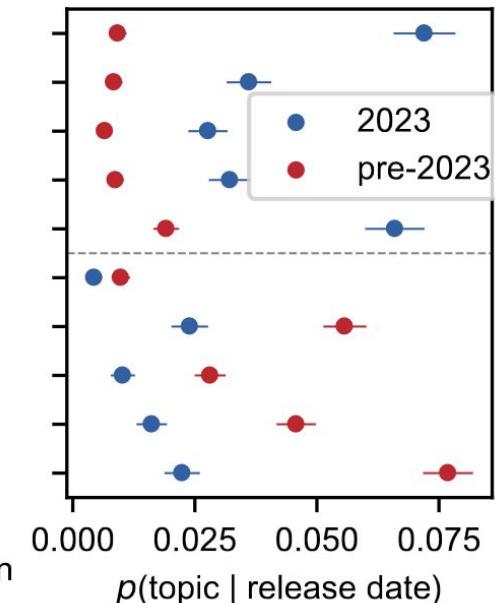
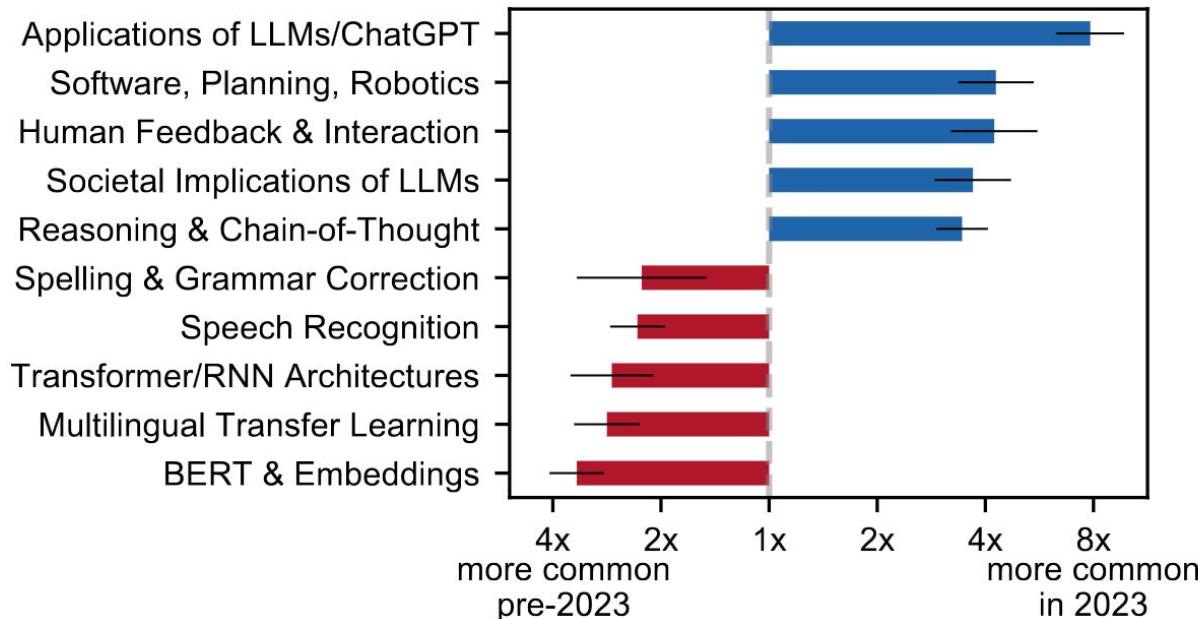


Network of the top 20 institutions in LLM research, showing collaborations.
Node size reflects publication volume, and edge thickness indicates collaborations .

Increased Focus on LLM Applications: BERT & Embeddings on the decline.



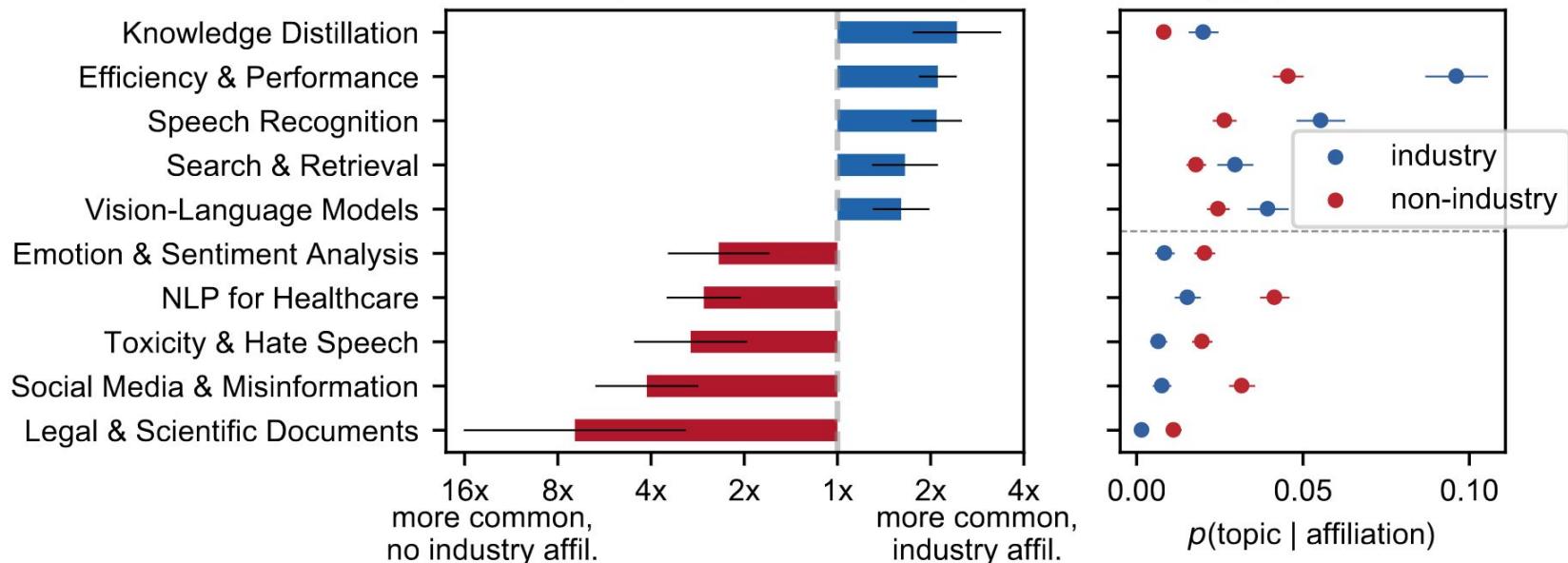
Fastest-growing topics in 2023 vs. 2018-22



Industry papers focus on core tech: academic papers focus on applications



Topic differences, industry vs. non-industry papers

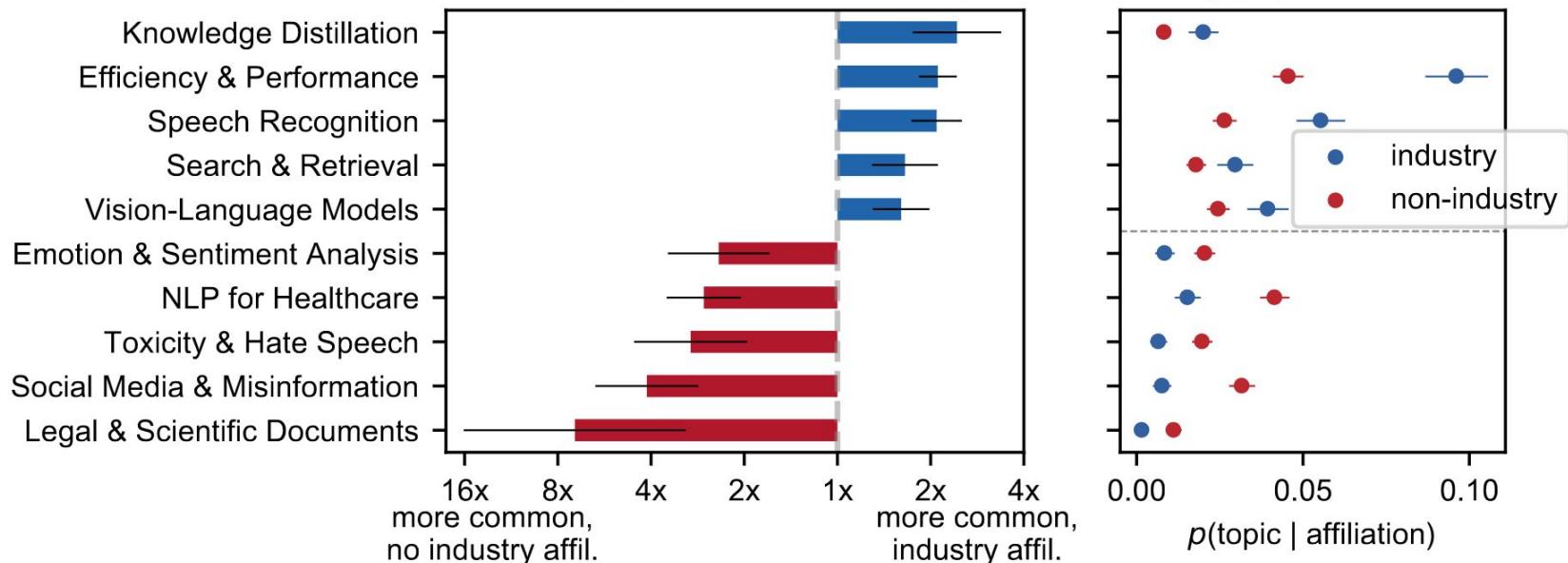


AI Concept	Financial Application	AI Concept	Financial Application
Knowledge Distillation Document summarization or extraction of content from text	E.g. Extract Annual R&D Costs from 10Ks or 10Qs	Sentiment Analysis Understanding the emotional valence of text or speech.	E.g. Measuring investor sentiment from Twitter feeds.
Efficiency & Performance Automating a currently manual or inefficient process.	E.g. Automated credit risk assessments.	NLP for Healthcare Extracting meaningful representations from healthcare texts.	E.g. Analyzing adverse effects of medications from EMR.
Speech Recognition Converting spoken statements into text, sentiment, id, etc.	E.g. Transcribing earnings calls for rapid information extraction.	Toxicity & Hate Speech Identifying and filtering harmful content.	E.g. Discrimination in the provision of financial services to minorities.
Search & Retrieval Identifying the most appropriate response from a set of possible answers.	E.g. Surfacing all regulations that are relevant for a given industry (helps understand compliance landscape).	Social Media & Misinformation Assessing the authenticity and reliability of a given claim.	E.g. Spotting fake news that could affect stock prices.
Vision-Language Models Joint interpretation of both text and visual information to draw a conclusion, or make a prediction.	E.g. Describing the performance of a stock given a visual chart, and key news stories.	Legal & Scientific Documents Use of legal and scientific documents for a variety of downstream tasks.	E.g. Understanding how scientific publications predict future employment in related topical domains.

Industry papers focus on core tech: academic papers focus on applications

—

Topic differences, industry vs. non-industry papers



Introduction and Overview: Trends in LLMs generally, and Finance Specifically

- **Trends in LLM Research, broadly [8 minutes]:**

Movva, R., Balachandar, S., Peng, K., Agostini, G., Garg, N., & Pierson, E. (2023). *Topics, Authors, and Networks in Large Language Model Research: Trends from a Survey of 17K arXiv Papers*. arXiv:2307.10700 [cs.DL]

- **Trends in LLM Research for Finance, specifically [7 minutes]:**

Synthesis of the data, methods, and topics presented in 52 recent papers on LLMs in Finance; papers were sourced from Google Scholar and ArXiv using the search terms:

('Large Language Models' OR 'GPT') AND
('Finance') AND
($\emptyset \cup \{\text{Application}\}$).

Where {Applications} were manually created and included the following topics: *Investment and Trading Insights, Risk Mitigation and Fraud Prevention, Customer Engagement and Services, Market Analysis and Sentiment Tracking, Credit and Financial Health*

Leading the Charge: The Top Three Cited LLM Papers in Finance



1

FinBERT: Financial Sentiment Analysis with Pre-trained Language Models

A specialized large language model adapted for the finance domain, based on Google's BERT algorithm. This model stands out for its superior ability to summarize contextual information in financial texts and significantly outperforms other machine learning models in sentiment classification, especially in processing financial texts with unique vocabulary and smaller training samples.

2

ChatGPT for (Finance) research: The Bananarama Conjecture

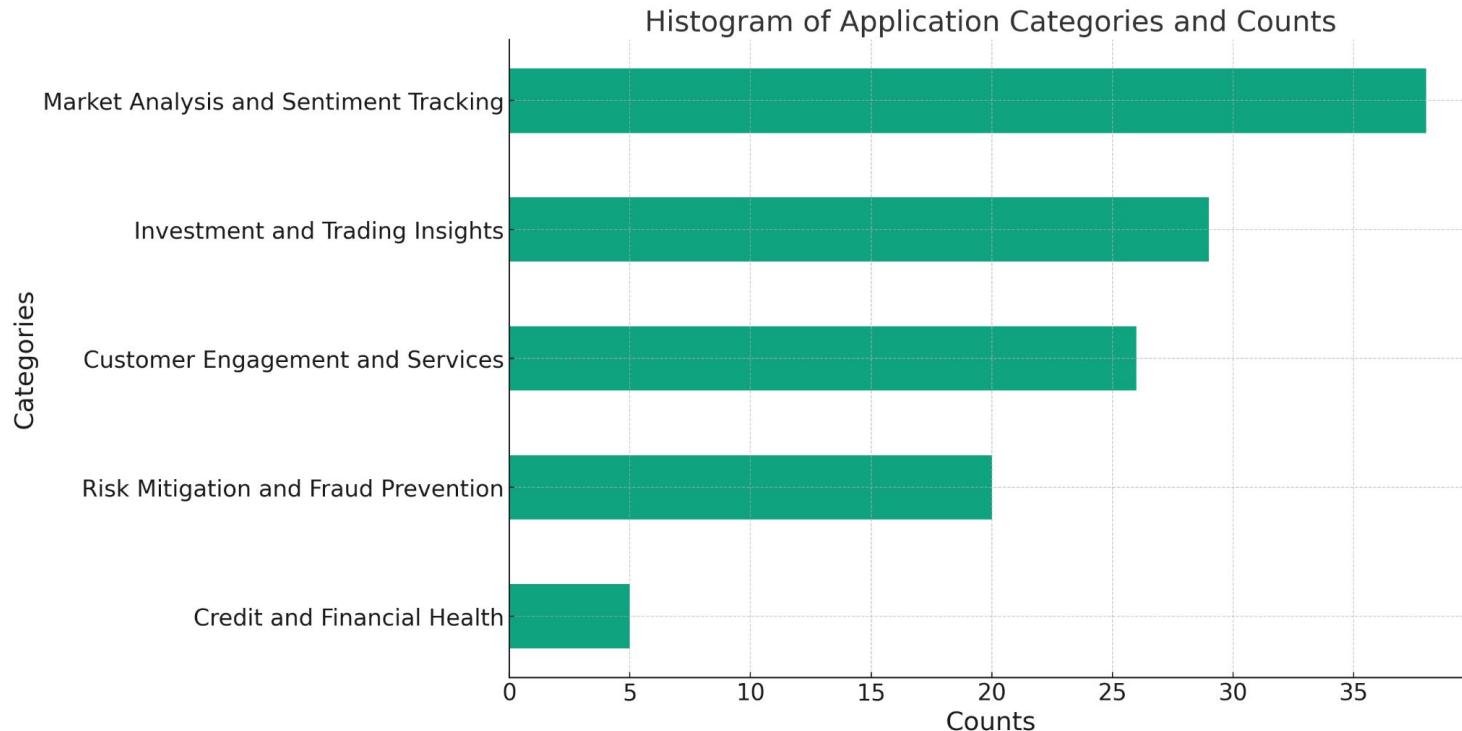
An examination of the use of ChatGPT in finance research. It stands out for its empirical testing of ChatGPT's effectiveness in various research stages, including idea generation, literature synthesis, data identification, and testing frameworks, highlighting the significant role of private data and researcher expertise in enhancing output quality.

3

BloombergGPT: A Large Language Model for Finance

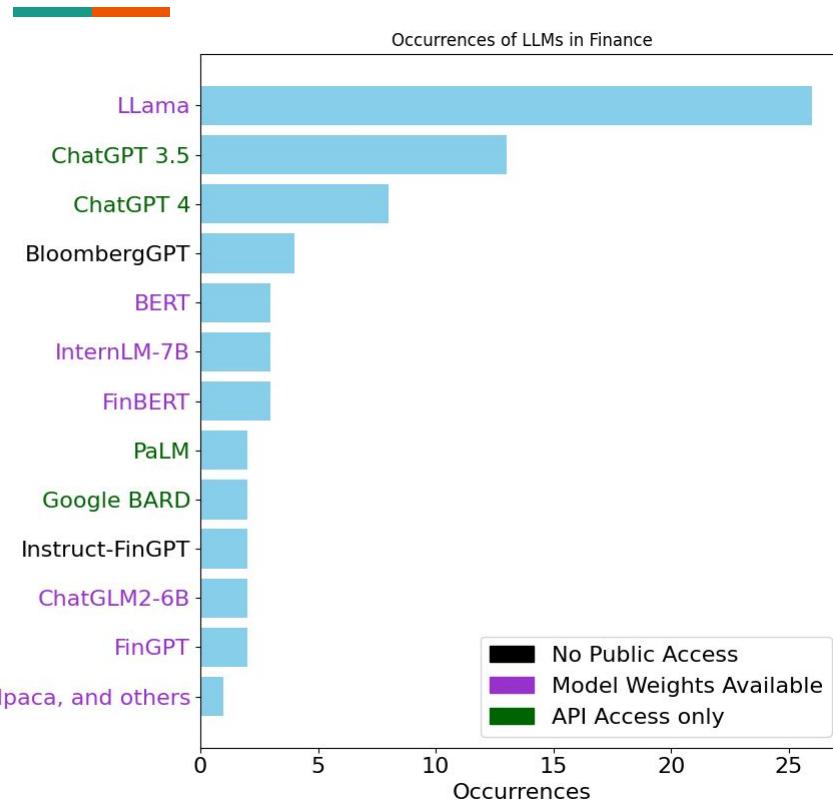
A 50 billion parameter language model tailored for financial applications, trained on a unique mix of financial and general-purpose datasets. This model is distinctive for its combination of a large, domain-specific dataset with broad training sources, achieving unprecedented performance in financial tasks while maintaining robustness in general language model benchmarks.

Market Analysis and Sentiment Tracking dominates research focus



LLama and ChatGPT Lead the Charge in Financial LLM Applications

Models



Why is LLama at the forefront?

LLama leads due to its adaptability and well-established ecosystem for fine-tuning, facilitating easier integration into financial applications.

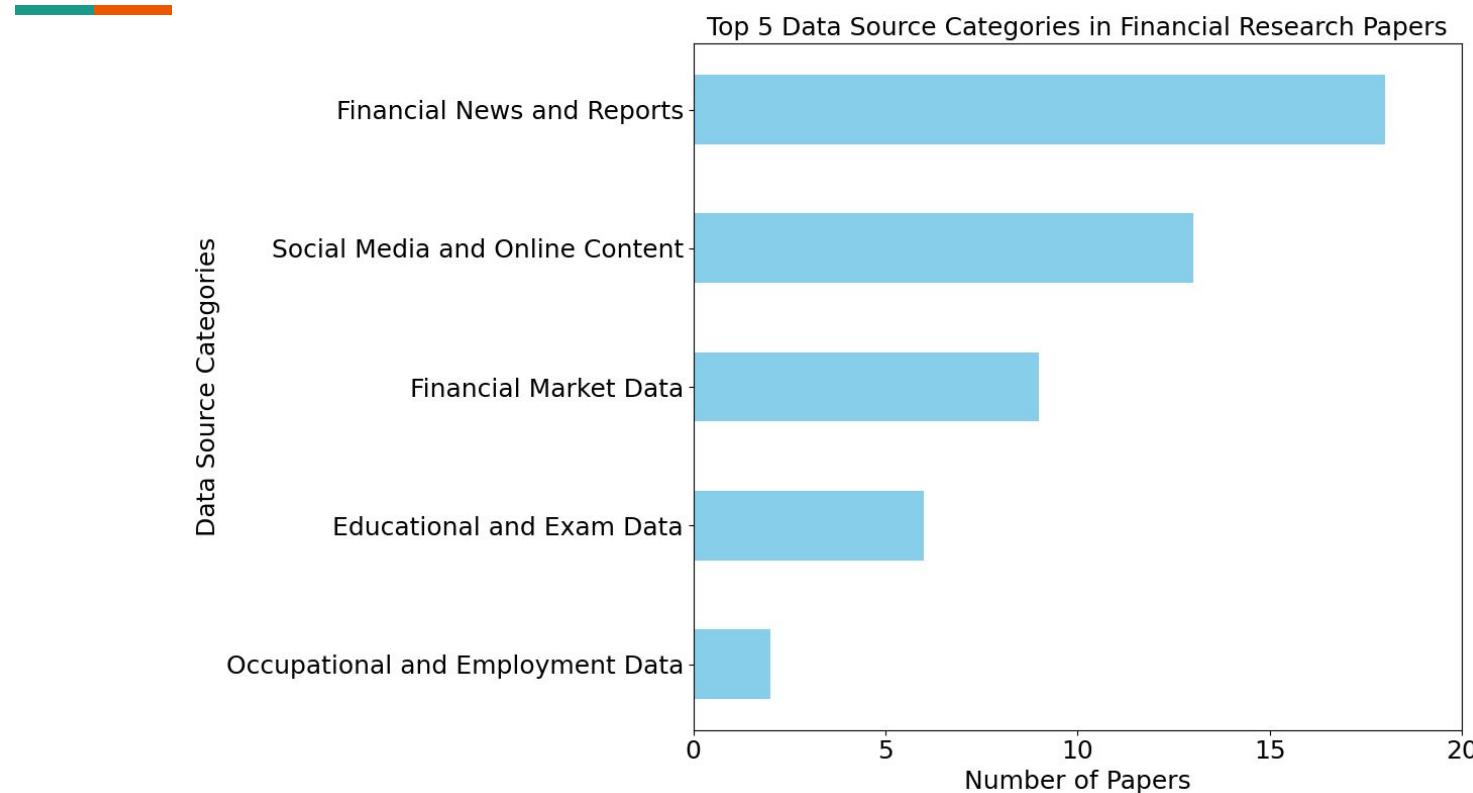
Why is ChatGPT 3.5 more prevalent than 4?

ChatGPT 3.5 enjoys greater usage over version 4 due to broader availability and cost-effectiveness, ensuring a balance between performance and accessibility.

Why is BloombergGPT used less than LLama?

BloombergGPT may see less utilization compared to LLAMA because of limited access in comparison to open-source models.

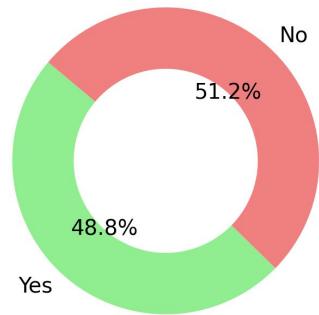
Prevalence of Data Types: Dominance of News/Social Media Sources



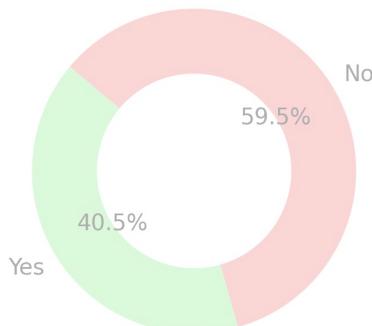
Adoption of Advanced Techniques: A Balance of Fine-Tuning and Prompt Engineering



Percentage of Finetuned Models



Percentage of Prompt-Engineered Models



Understanding Fine-Tuning in Machine Learning

Fine-tuning is a process where a pre-trained machine learning model is further trained (or 'fine-tuned') with a smaller, specialized dataset to adapt to specific tasks or industries. In finance, fine-tuning involves training models like LLMs on financial data, allowing them to understand and generate industry-specific insights, conduct sentiment analysis on market reports, or predict stock trends with greater accuracy.

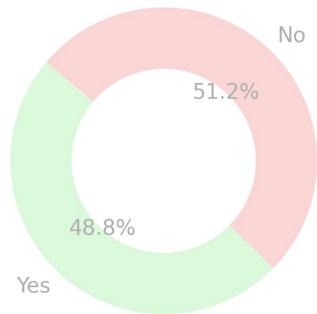
Exploring Prompt Engineering in LLMs

Prompt engineering is the art of crafting input prompts to elicit the desired output from AI models, particularly in natural language processing (NLP). It's a strategic approach where the prompt's structure and content are optimized to guide the AI. For financial professionals, prompt engineering can be used to obtain precise market analyses, generate financial narratives, or query complex datasets without needing to understand the underlying AI model.

Adoption of Advanced Techniques: A Balance of Fine-Tuning and Prompt Engineering



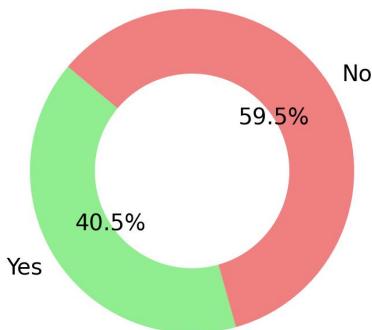
Percentage of Finetuned Models



Understanding Fine-Tuning in Machine Learning

Fine-tuning is a process where a pre-trained machine learning model is further trained (or 'fine-tuned') with a smaller, specialized dataset to adapt to specific tasks or industries. In finance, fine-tuning involves training models like LLMs on financial data, allowing them to understand and generate industry-specific insights, conduct sentiment analysis on market reports, or predict stock trends with greater accuracy.

Percentage of Prompt-Engineered Models



Exploring Prompt Engineering in LLMs

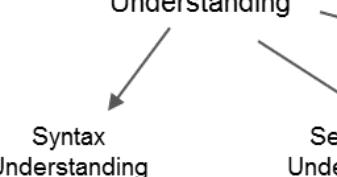
Prompt engineering is the art of crafting input prompts to elicit the desired output from AI models, particularly in natural language processing (NLP). It's a strategic approach where the prompt's structure and content are optimized to guide the AI. For financial professionals, prompt engineering can be used to obtain precise market analyses, generate financial narratives, or query complex datasets without needing to understand the underlying AI model.

Presentation Outline

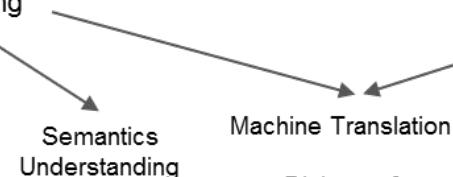
- **Introduction and Overview [15 mins talk; 5 mins Q&A]:**
Exploring the Evolution of LLMs and their applications in Finance.
Mohammad Ghassemi, Reza Khanmohammadi
- **Application Insights [20 mins; 5 mins Q&A]:**
Deep Dive into Financial Applications of LLMs.
Xiaomo Liu, Simerjot Kaur
- **Addressing Challenges [25 mins; 5 mins Q&A]:**
Analyzing the Limitations and Challenges in LLMs.
Manling Li

Applications within NLP Domain

Natural Language Understanding



- POS Tagging
- Dependency Parsing
- Tokenization
- ...



- NER
- Topic Classification
- Sentiment Analysis
- ...

Machine Translation
Dialogue & Discourse

Natural Language Generation

- Text Summarization
- Data to Text
- Document Generation
- ...

Natural Language Question Answering

- Information Retrieval
- Knowledge-based QA
- Reasoning-based QA

Need understanding,
reasoning & generation

Language &
Multimodality

- Image Captioning
- Table Understanding
- Text to Image

NLP Modeling Complexity

```
graph LR; NLU[Natural Language Understanding] --> SU[Syntax Understanding]; NLU --> SU[Semantics Understanding]; NLU --> MT[Machine Translation]; NLU --> NLG[Natural Language Generation]; NLU --> NQA[Natural Language Question Answering]; NLU --> LM[Language & Multimodality];
```

Example of Use Case Evaluation

Use Case of Financial Analysis:
(1) Apple iPhone sales next quarter;
(2) Hold or Sell \$APPL



Find Information:

- Search iPhone news
- Collect smartphone research reports
- Research supply chain data



Digest Information:

- Investigate smartphone market sentiments
- Identify key suppliers and if there are potential risks



Reach Conclusion:

- Reason on if iPhone market will continue to grow and how fast
- Estimate Impacts to Apple's revenue
- Suggest stock hold or sell



Publish Results:

- Summarize conclusions
- Present numbers in tables & charts



*Taxonomy based on ACL list of tracks (after: "Holistic Evaluation of Language Models", Percy Liang et al 2021, Stanford University)

LLM Evaluation Considerations

Data

- LLMs are trained on large amounts of data, a lot of which was obtained from the web.
- For proprietary models such as ChatGPT, GPT-4, PaLM2 the exact training datasets are not always precisely known, or publicly available.
- The data knowledge cut offs may also be subjects to change with updates and especially in the interactive setups.
- Therefore, one needs to be careful when evaluating LLMs to aim to minimise the chance that the testing data has once been part of pre-training or fine-tuning.

LLM Evaluation Considerations

Evaluation of generative models on classification tasks

- A lot of popular LLMs are generative models, so they generate spans of text as answers rather than labels of a specified format. It is also possible that LLMs would propose a new answer, outside of the specified classes.
- How to solve
 - **Prompt engineering** (add precise instructions to the prompt on the desired output format)
 - **Postprocessing** (convert the generated answer into one of the answers from the label set)
 - **Fine-tuning** (fine tune on a small set of examples demonstrating the desired output format)
 - **OpenAI Functions** (if using openai environment, models are able to produce structured output in JSON format)

“What is the sentiment of the following text? Text: I am so happy and excited!”

User

The following text is positive, because it expresses the positive feelings of ...

LLM

Desired answer for the classification system:
“positive”

LLM Benchmark: Hugging Face Open LLM Leaderboard

- Focus on open sources LLMs
- Evaluate derivative models from foundation models like Llama
- Contain four evaluation tasks

The 🎉 Open LLM Leaderboard aims to track, rank and evaluate open LLMs and chatbots.

🎉 Submit a model for automated evaluation on the 🎉 GPU cluster on the "Submit" page! The leaderboard's backend runs the great [Eleuther AI Language Model Evaluation Harness](#) - read more details in the "About" page!

(LLM Benchmark) About Submit here!

Search for your model and press ENTER...

Select columns to show

Average	ARC	HellaSwag	MMLU	TruthfulQA	Type
<input type="checkbox"/>					

Precision

torch.float16	torch.bfloat16	torch.float32	8bit	4bit	GPTQ
<input checked="" type="checkbox"/>					

Model types

pretrained	fine-tuned	instruction-tuned	RL-tuned
<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>

?Unknown

Precision

Unknown	< 1.5B	-3B	-7B	-13B	-35B	60B+
<input checked="" type="checkbox"/>						

Model sizes

Model	Average	ARC	HellaSwag	MMLU	TruthfulQA
AIDC-ai:business/Marcoroni-70B:v1	74.06	73.55	87.62	70.67	64.41
ICBU-NPU/FashionGPT-70B-V1.1	74.05	71.76	88.2	70.99	65.26
adonlee/LLaMA_2_70B_LoRA	73.9	72.7	87.55	70.84	64.52
uni-tianyan/Uni-TianYan	73.81	72.1	87.4	69.91	65.81
Riiid/sheep-duck-llama-2	73.69	72.35	87.78	70.82	63.8
Riiid/sheep-duck-llama-2	73.67	72.27	87.78	70.81	63.8
fangloveskazi/ORCA_LLaMA_70B_OLoRA	73.4	72.27	87.74	70.23	63.37
ICBU-NPU/FashionGPT-70B-V1	73.26	71.08	87.32	70.7	63.92

LLM Benchmark: Holistic Evaluation of Language Models

- Examine both Open- & closed-source LLMs
 - 78 models, 42 scenarios, and 59 metrics
- Use a comprehensive set of NLP tasks
- No derivative models

[Blog post](#)[Paper](#)[GitHub](#)

A language model takes in text and produces text:

A helm is a → **Language Model** → *wheel for steering a ship...*

Despite their simplicity, language models are increasingly functioning as the foundation for almost all language technologies from question answering to summarization. But their immense capabilities and risks are not well understood. Holistic Evaluation of Language Models (HELM) is a living benchmark that aims to improve the transparency of language models.

1. **Broad coverage and recognition of incompleteness.** We define a taxonomy over the scenarios we would ideally like to evaluate, select scenarios and metrics to cover the space and make explicit what is missing.



2. **Multi-metric measurement.** Rather than focus on isolated metrics such as accuracy, we simultaneously measure multiple metrics (e.g., accuracy, robustness, calibration, efficiency) for each scenario, allowing analysis of tradeoffs.

Scenarios	Metrics
S1	✓
S2	✓
S3	✓
S4	✓
S5	✓
S6	✓
S7	✓
S8	✓
S9	✓
S10	✓
S11	✓
S12	✓
S13	✓
S14	✓
S15	✓
S16	✓
S17	✓
S18	✓
S19	✓
S20	✓
S21	✓
S22	✓
S23	✓
S24	✓
S25	✓
S26	✓
S27	✓
S28	✓
S29	✓
S30	✓
S31	✓
S32	✓
S33	✓
S34	✓
S35	✓
S36	✓
S37	✓
S38	✓
S39	✓
S40	✓
S41	✓
S42	✓

3. **Standardization.** We evaluate all the models that we have access to on the same scenarios with the same adaptation strategy (e.g., prompting), allowing for controlled comparisons. Thanks to all the companies for providing API access to the limited-access and closed models and *Together* for providing the infrastructure to run the open models.

4. **Transparency.** All the scenarios, predictions, prompts, code are available for further analysis on this website. We invite you to click below to explore!

Domain-specific Challenges

wud2wood

★★★★★ Nice shower shoes!!
Reviewed in the United States on September 24, 2023
Size: 8.5-9 Wide Women/7.5-8 Men | Color: Black | Verified Purchase

Love these I have a locker room shower in my basement, and it hasn't been tiled. I use these to stay off the concrete. Works perfectly.

Informal language that doesn't follow formal grammar

Amazon Review

User Generated
Language: informal
Knowledge: no financial knowledge

Consumed by average people

Apple reports better-than-expected quarter driven by iPhone sales

PUBLISHED THU, MAY 4 2023 9:00 AM EDT | UPDATED THU, MAY 4 2023 6:32 PM EDT

Kirk Lewising @KIRKLEWISING

SHARE f t in e

KEY POINTS • Apple reported second-fiscal quarter earnings on Thursday that beat Wall Street's soft expectations, driven by stronger-than-anticipated iPhone sales.
• However, Apple's overall sales fell for the second quarter in a row.

ABC TV Power Lunch WATCH LIVE Listen

UP NEXT: Closing Bell 03:00 pm ET Listen

Financial News

News Professionals Generated
Language: formal
Knowledge: some financial knowledge

Consumed by financial professionals

Debt SEC Filings (10Q)

As of September 24, 2022, the Company had outstanding fixed-rate notes with varying maturities for an aggregate principal amount of \$111.8 billion (collectively the "Notes"), with \$11.1 billion payable within 12 months. Future interest payments associated with the Notes total \$41.3 billion, with \$2.9 billion payable within 12 months.

The Company also issues unsecured short-term promissory notes ("Commercial Paper") pursuant to a commercial paper program. As of September 24, 2022, the Company had \$10.0 billion of Commercial Paper outstanding, all of which was payable within 12 months.

Requires domain knowledge to comprehend

Financial Professionals Generate
Language: formal
Knowledge: professional financial knowledge

Financial NLP Focus LLM Evaluation

Research Setting 1: Tasks with NLP modeling complexity & financial knowledge

Category	Sentiment Analysis	Classification	NER	RE	QA
Complexity	Easy	Easy	Hard	Hard	Hard
Knowledge	Low	Low	High	High	High
Dataset	FPB/FiQA/TweetFinSent	Headline	NER	REFinD	FinQA/ConvFinQA
Eval. Metrics	Weighted F1	Weighted F1	Macro F1	Macro F1	Accuracy
#Test samples	970/223/996	2,114	98	4300	1,147/421

Table 1: Statistics of the five tasks and eight datasets used in this study.

Research Setting 2: Comparison between OpenAI LLMs and other language models

	GPT-NeoX	OPT66B	BLOOM	BloombergGPT	Fine-Tuned		GPT3.5	GPT4.0
Category	Prior LLMs			Domain Specific LMMs	Transformer Based		OpenAI LLMs	
Size	20B	66B	176B	50B			Unknown (>175B?)	Unknown (>175B?)

Research Setting 3: Comparison of Prompt Engineering Strategies

Reference: **Are ChatGPT and GPT-4 General-Purpose Solvers for Financial Text Analytics? An Examination on Several Typical Tasks**, Li, X. et al., EMNLP 2023 Industrial Track

Sentiment Analysis

- Financial PhraseBank
 - Sampled from financial news and company press releases.
 - Tagged as positive, negative or neutral by 16 annotators with domain knowledge.
 - E.g., Operating profit rose to EUR 13.1 mn from EUR 8.7 mn in the corresponding period in 2007 representing 7.7 % of net sales

positive

- FiQA
 - Include aspect-based sentiments from news in the financial domain

```
"sentence": "Tesco Abandons Video-Streaming Ambitions in Blinkbox Sale",
"info": [
  {
    "snippets": "[Video-Streaming Ambitions]",
    "target": "Blinkbox",
    "sentiment_score": "-0.195",
    "aspects": "[Corporate/Strategy]"
  }
],
```

negative

- TweetFinSent
 - Made with Tweets that contain retail investors' mood to a specific stock.
 - Created by AI Research
 - Biggest up moves \$COIN Coinbase, \$PATH UIPath, \$RBLX Roblox, \$DKNG DraftKings, but \$PLTR Palantir is stable'

positive

neutral

Data	50% Agreement		100% Agreement	
	Accuracy	F1 score	Accuracy	F1 score
ChatGPT(0)	0.78	0.78	0.90	0.90
ChatGPT(5)	0.79	0.79	0.90	0.90
GPT-4(0)	<u>0.83</u>	<u>0.83</u>	<u>0.96</u>	<u>0.96</u>
GPT-4(5)	0.86	0.86	0.97	0.97
BloombergGPT(5)	/	0.51	/	/
GPT-NeoX(5)	/	0.45	/	/
OPT66B(5)	/	0.49	/	/
BLOOM176B(5)	/	0.50	/	/
FinBert	0.86	0.84	0.97	0.95

ChatGPT(FT) 0.89 0.89

Model	Accuracy	Weighted F1
ChatGPT(0)	68.48	68.60
ChatGPT(5)	69.93	70.05
GPT-4(0)	69.08	69.17
GPT-4(5)	<u>71.95</u>	72.12
ChatGPT((0_no_emoji))	64.40	64.43
ChatGPT((5_no_emoji))	67.37	67.61
GPT-4((0_no_emoji))	67.26	67.45
GPT-4((5_no_emoji))	70.58	70.44
RoBERTa-Twitter	72.30	<u>71.96</u>

Table 4: Results on the TweetFinSent dataset.

ChatGPT(FT) 75.3 75.3

Text Classification

- **Task:** Financial News Headline Classification

- **Data:** news headlines about gold commodity from various financial news provider sites (Reuters, The Hindu, The Economic Times, Bloomberg etc.)
- **Six categories:** price up, price down, price stable, past price, future price, and asset comparison

price up

1. Dec. gold settles at \$1,293.80/oz, **up \$8.80**, or 0.7%.

price down

2. Gold prices **slide \$14.90** or 1.1%, to \$1,289.80 an ounce
past price

3. Gold imports dip 8% to \$31.72 bn in **2015-16**.

Model	Weighted F1
ChatGPT (0)	71.78
ChatGPT (5)	74.84
GPT-4 (0)	84.17
GPT-4 (5)	<u>86.00</u>
BloombergGPT (5)	82.20
GPT-NeoX (5)	73.22
OPT66B (5)	79.41
BLOOM176B (5)	76.51
BERT	95.36

Table 5: Results on the headline classification task.

ChatGPT(FT)

92.64

Name Entity Recognition

Task: FIN3 - NER

- Data: sentences extracted from SEC financial agreements
- Four NER labels: PER, LOC, ORG and MISC

LOC
ORG

LOAN AGREEMENT

This LOAN AGREEMENT, dated as of November 17, 2014 (this “Agreement”), is made by and among Auxilium Pharmaceuticals, Inc., a corporation incorporated under the laws of the State of Delaware (“U.S. Borrower”), Auxilium UK LTD, a private company limited by shares registered in England and Wales (“UK Borrower” and, collectively with the U.S. Borrower, the “Borrowers”) and Endo Pharmaceuticals Inc., a corporation incorporated under the laws of the State of Delaware (“Lender”).

"tags": [0, 0, 0, 0, 0, 3, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]

- Traditional approaches:
 - Rule-based methods
 - Statistical-modeling methods

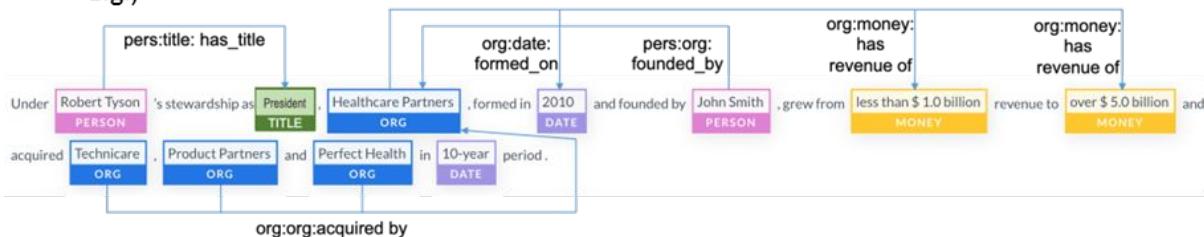
Model	Entity F1
ChatGPT(0)	29.21
ChatGPT(20)	51.52
GPT-4(0)	36.08
GPT-4(20)	56.71
BloombergGPT(20)	60.82
GPT-NeoX(20)	60.98
OPT66B(20)	57.49
BLOOM176B(20)	55.56
CRF(CoNLL)	17.20
CRF(FIN5)	82.70

Table 6: Results of few-shot performance on the NER dataset. CRF(CoNLL) refers to CRF model that is trained on general CoNLL data, CRF(FIN5) refers to CRF model that is trained on FIN5 data. Again, we choose the same shot as BloombergGPT for fair comparison.

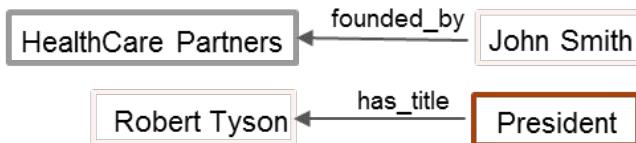
Relation Extraction

Task: REFinD

- **Data:** created from SEC 10-K/Q filings by AI Research
- **22 relation types:** org:org:subsidiary_of, org:org:agreement_with, etc.
- E.g.,



- Traditional approaches:
 - Rule-based methods
 - Unsupervised methods
 - Supervised methods



Model	Macro F1
ChatGPT (0)	20.97
ChatGPT (10)	29.53
GPT-4 (0)	42.29
GPT-4 (10)	46.87
Luke-base (fine-tune)	56.30

Table 7: Results on the REFinD dataset.

Question Answering

- **Tasks:** FinQA and ConvFinQA
- **Data:** Questions derived from earnings reports companies. Demands numerical reasoning and understanding of structured data and financial concepts. Emphasizes the ability to relate follow-up questions to previous conversation context.

Page 91 from the annual reports of GRMN (Garmin Ltd.)
The fair value for these options was estimated at the date of grant using a Black-Scholes option pricing model with the following weighted-average assumptions for 2006, 2005 and 2004:

	2006	2005	2004
Weighted average fair value of options granted	\$20.01	\$9.48	\$7.28
Expected volatility	0.3534	0.3224	0.3577
Distribution yield	1.00%	0.98%	1.30%
Expected life of options in years	6.3	6.3	6.3
Risk-free interest rate	5%	4%	4%

... The total fair value of shares vested during 2006, 2005, and 2004 was \$9,413, \$8,249, and \$6,418 respectively. The aggregate intrinsic values of options outstanding and exercisable at December 31, 2006 were \$204.1 million and \$100.2 million, respectively. (... abbreviate 10 sentences ...)

Question: Considering the weighted average fair value of options , what was the change of shares vested from 2005 to 2006?
Answer: -400
Calculations:
$$\left(\frac{9413}{20.01} \right) - \left(\frac{8249}{9.48} \right) = -400$$

Program:
divide (9413, 20.01) divide (8249, 9.48)
_____ _____
 subtract (#0, #1)

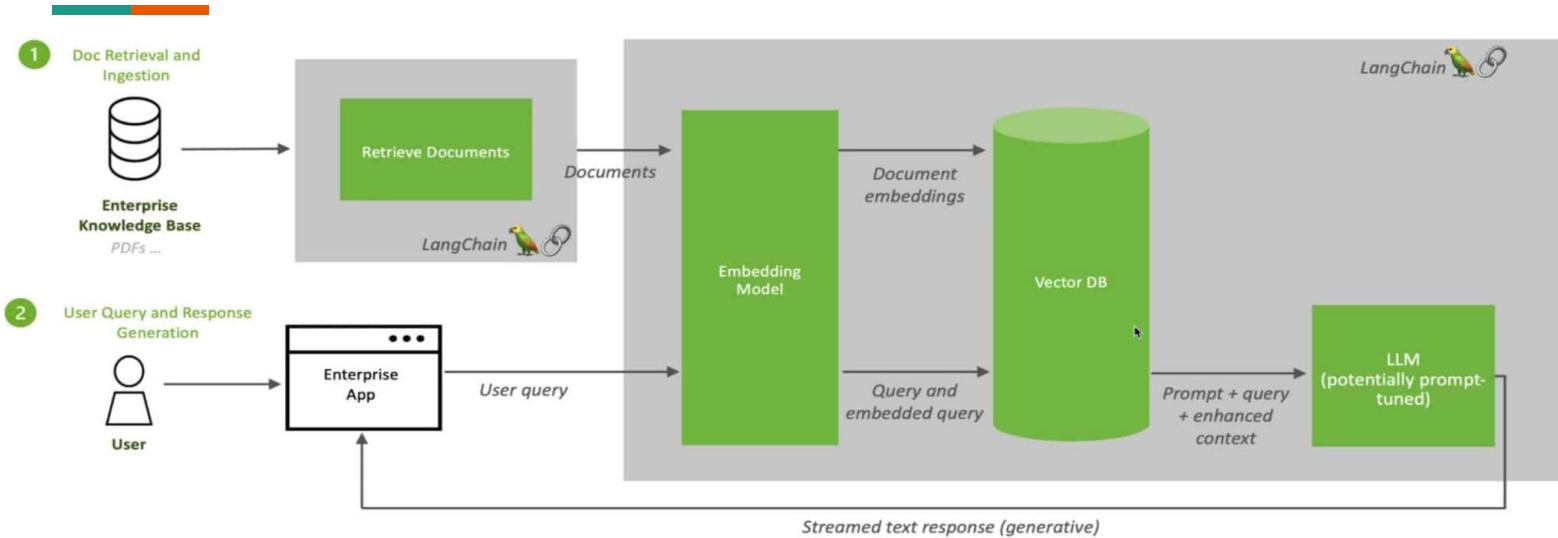
- More advanced financial analysis QA:

Can GPT models be Financial Analysts? An Evaluation of ChatGPT and GPT-4 on mock CFA Exams, Callanan, E. et al, arXiv:2310.08678

Model	FinQA	ConvFinQA
ChatGPT(0)	48.56	59.86
ChatGPT(3)	51.22	/
ChatGPT(CoT)	63.87	/
GPT-4 (0)	68.79	76.48
GPT-4 (3)	69.68	/
GPT-4 (CoT)	78.03	/
BloombergGPT (0)	/	43.41
GPT-NeoX (0)	/	30.06
OPT66B (0)	/	27.88
BLOOM176B (0)	/	36.31
FinQANet(fine-tune)	68.90	61.24
Human Expert	91.16	89.44
General Crowd	50.68	46.90

Table 8: Model performance (accuracy) on the question answering tasks. FinQANet here refers to the best-performing FinQANet version based on RoBERTa-Large (Chen et al., 2022a). Few-shot and CoT learning cannot be executed on ConvFinQA due to the conservation nature of ConvFinQA.

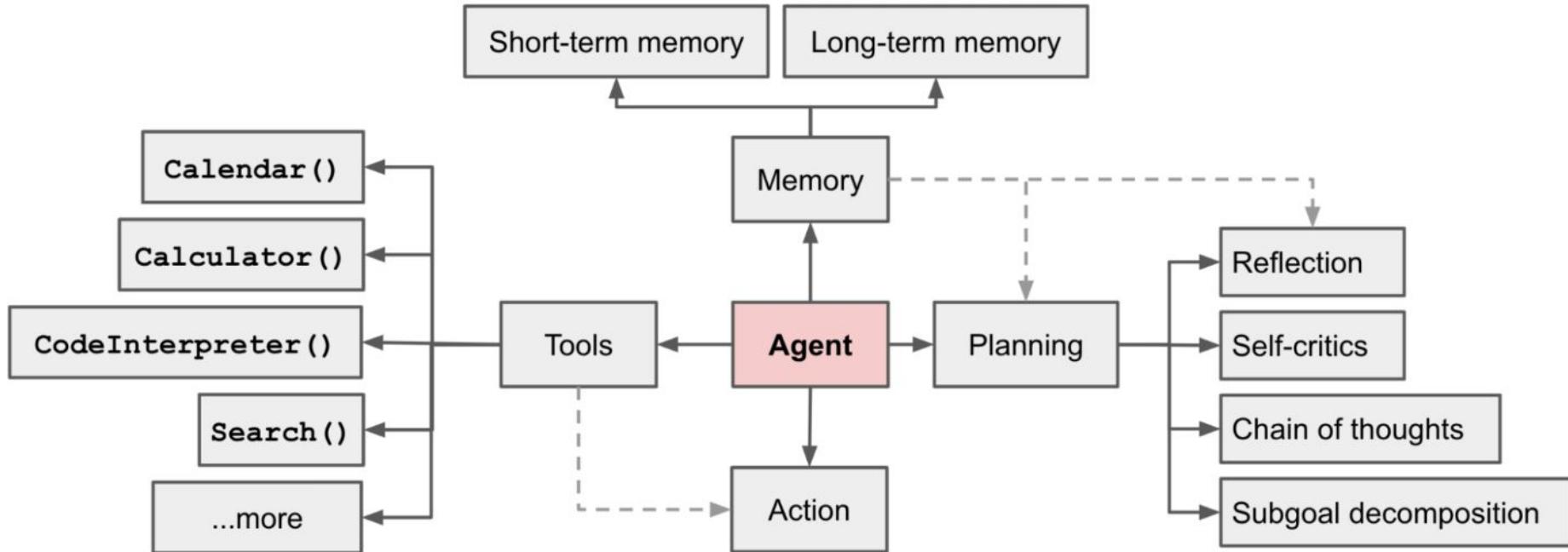
Retrieval Augmented Generation (RAG)



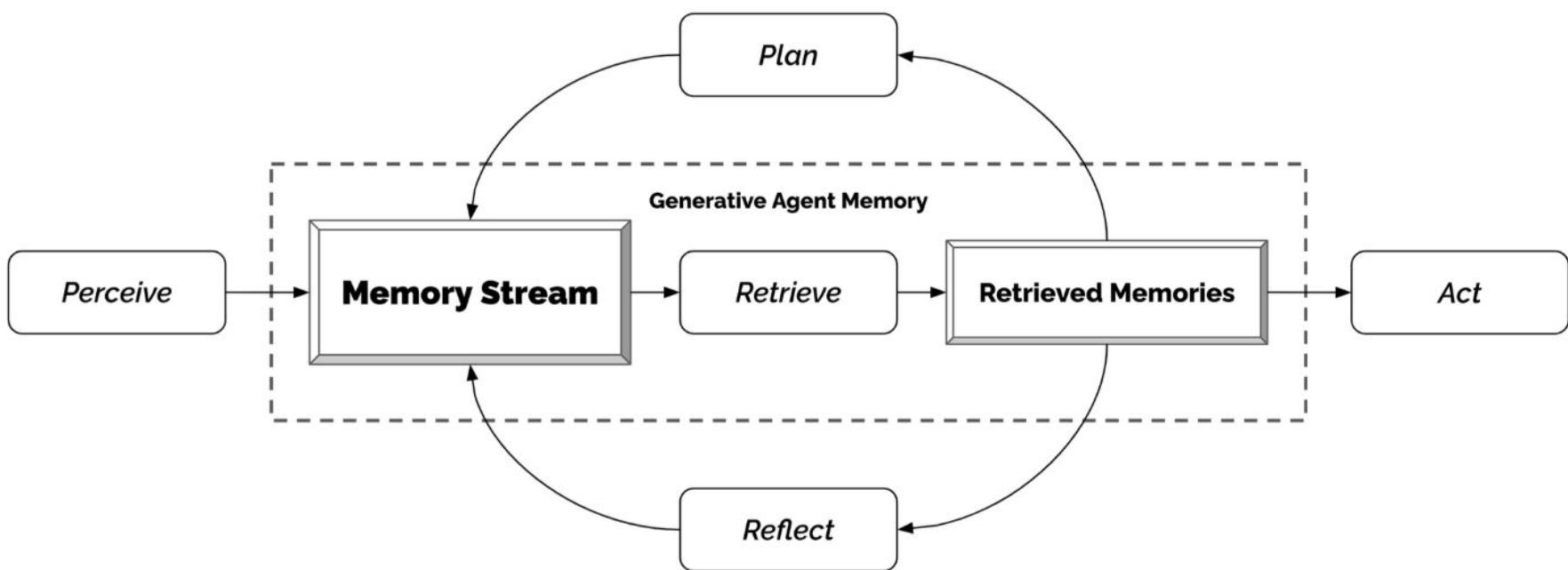
Category of Questions	N	Questions (%)	GPT-3.5	GPT-4
Yes/No	143	42.43%	41.96%	46.85%
Single-Select from options	73	21.66%	38.36%	56.16%
Single-Select from options (number)	52	15.43%	17.31%	32.69%
Multiple-Select from options	43	12.76%	4.65%	25.58%
Number-extraction	26	7.72%	26.92%	19.23%

Accuracy of GPT3.5 and GPT4 on Cogtale dataset on different types of questions.

LLM Powered Autonomous Agents



Generative Agents



Discussions & Conclusions

Remark 1: Comparison of LLMs

- GPT3.5 & GPT4.0 demonstrate substantial superiority over other large language models
- Despite being trained on financial corpora, BloombergGPT exhibits underwhelming performance.
- GPT-4.0 exhibits an improvement of 10% over GPT-3.5 in relatively easy tasks, while showcasing
- an impressive 20-100% enhancement in more challenging tasks.

Remark 2: Prompt Engineering Strategies

- Few-shot prompts lead to 1-4% performance boost on GPT3.5 and 4.0 models over zero-shot
- Chain-of-Thought is very effective and leads to 20-30% accuracy improvements

Remark 3: LLMs vs. Fine-tuning

- GPT4.0 can perform on par with fine-tuned models on simple tasks like sentiment analysis and classification
- Fine-tuned models are still better on hard tasks like NER & RE that requires complex NLP modeling
- On QA tasks with reasoning, LLMs show advantages over smaller fine-tuned model

Remark 4: Using LLMs in Financial Services

- Should consider LLMs with prompt engineering to replace simple NLP modeling
- LLMs for hard tasks are still far from satisfactory from industry requirement standpoint

Presentation Outline

- **Introduction and Overview [15 mins talk; 5 mins Q&A]:**
Exploring the Evolution of LLMs and their applications in Finance.
Mohammad Ghassemi, Reza Khanmohammadi
- **Application Insights [20 mins; 5 mins Q&A]:**
Deep Dive into Financial Applications of LLMs.
Xiaomo Liu, Simerjot Kaur
- **Addressing Challenges [25 mins; 5 mins Q&A]:**
Analyzing the Limitations and Challenges in LLMs.
Manling Li

Presentation Outline: Challenges

- **Hallucination Challenges:**

Challenges of control factual and truthful results of LLMs.

10 mins

- **Adoption Challenges:**

Challenges of using LLMs in finance.

5 mins

- **Privacy and Security Challenges:**

Challenges of using LLMs in environments with confidential data like customer info.

5 mins

- **Bias Challenges:**

Challenges of bias existing within LLMs and during the usage of LLMs.

5 mins

Presentation Outline: Challenges

- **Hallucination Challenges:**

Challenges of control factual and truthful results of LLMs.

10 mins

- **Adoption Challenges:**

Challenges of using LLMs in finance.

5 mins

- **Privacy and Security Challenges:**

Challenges of using LLMs in environments with confidential data like customer info.

5 mins

- **Bias Challenges:**

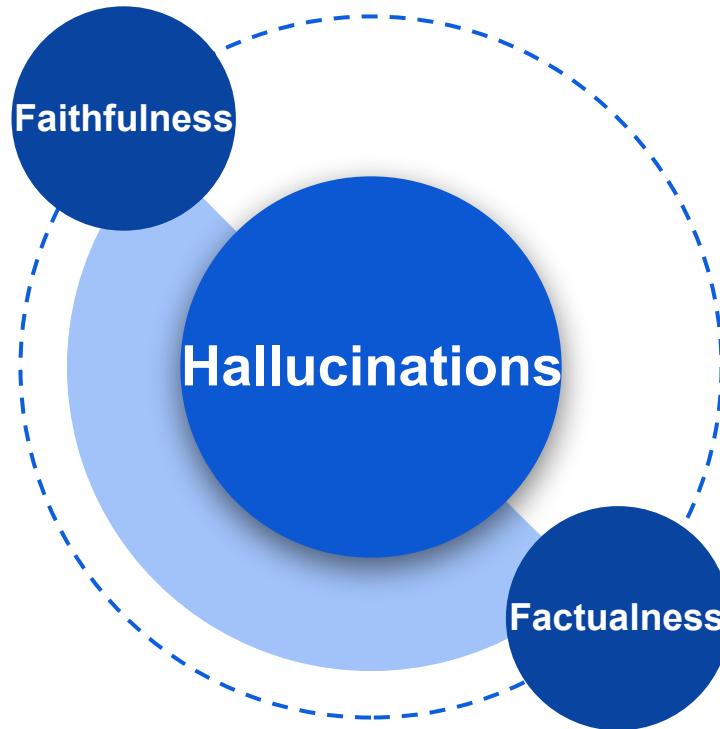
Challenges of bias existing within LLMs and during the usage of LLMs.

5 mins

What are Hallucinations?

Faithfulness:

Generated text is not faithful to the input context

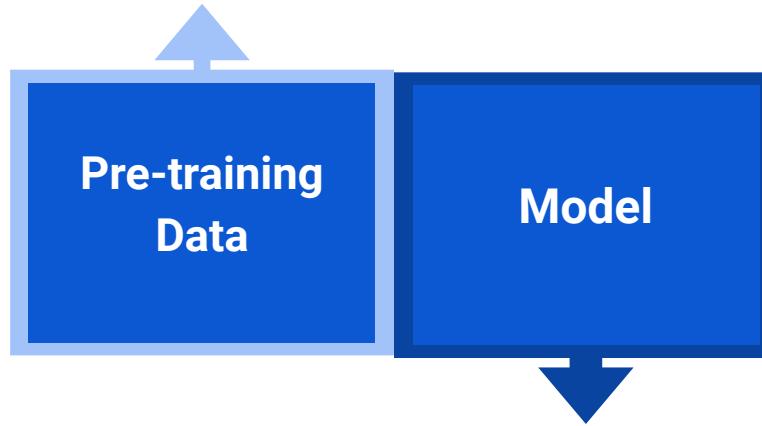


Factualness:

Generated text is not factually correct with respect to world knowledge

Why language models hallucinate?

- Contradicting information, Inaccurate knowledge [1]
- Knowledge bias from repeated data [1]



- Autoregressive nature of language modelling [2]
- Exposure bias [3]
- Incorrect parametric knowledge [4]

[1] Lee, Katherine, et al. "Deduplicating training data makes language models better." *arXiv preprint arXiv:2107.06499* (2021).

[2] Lee, Nayeon, et al. "Factuality enhanced language models for open-ended text generation." *Advances in Neural Information Processing Systems* 35 (2022): 34586-34599.

[3] Wang, Chaojun, and Rico Sennrich. "On exposure bias, hallucination and domain shift in neural machine translation." *arXiv preprint arXiv:2005.03642* (2020).

[4] Longpre, Shayne, et al. "Entity-based knowledge conflicts in question answering." *arXiv preprint arXiv:2109.05052* (2021).

Hallucination Overview

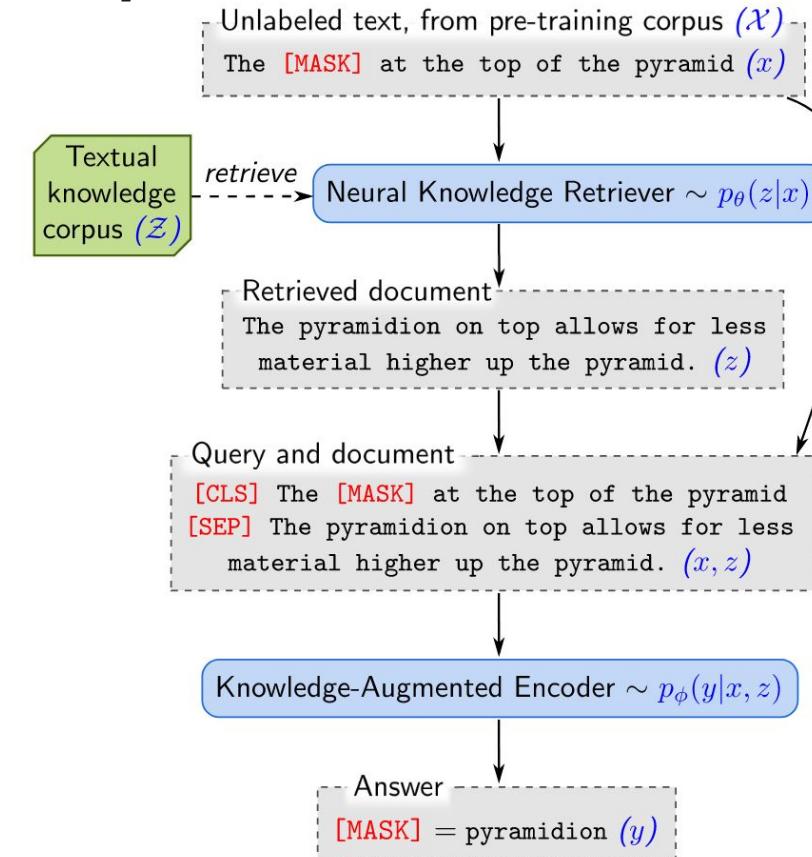
Main approaches to mitigating hallucination:

1. Retrieval Augmentation
2. Verification
3. Controlled Generation
4. Knowledge Editing
5. Other: deduplicating training data, prompting etc.

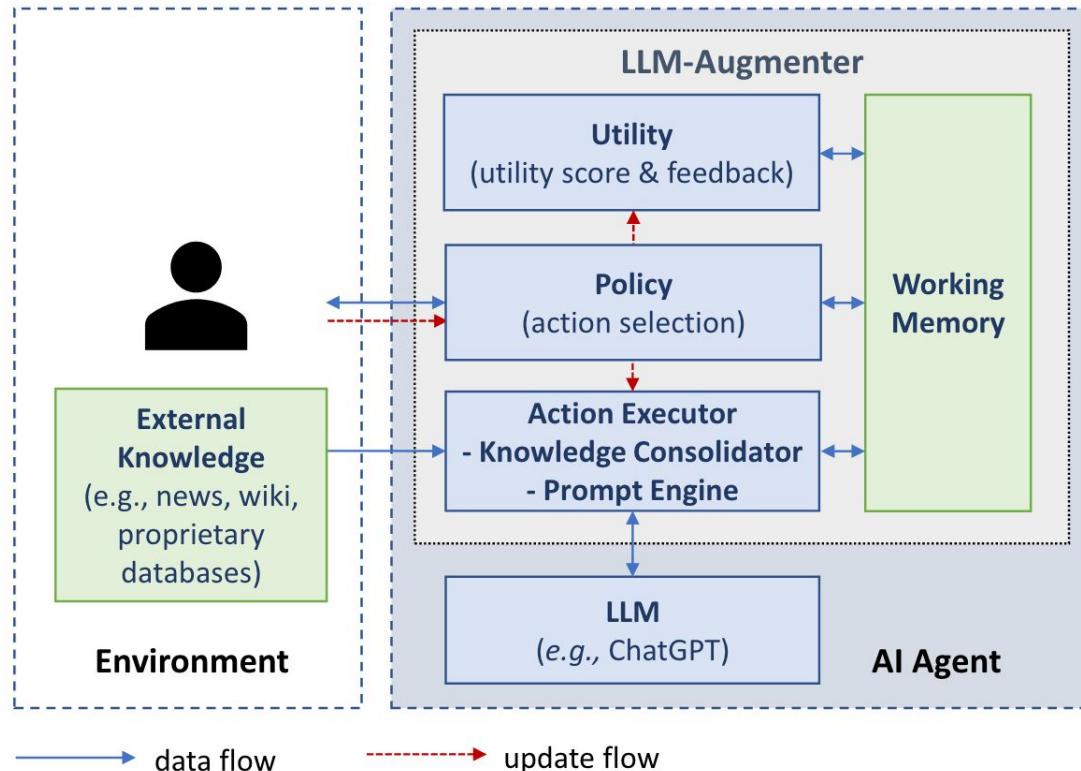
Related topic: model attribution / generating with citations

Compare Against Retrieval Corpus

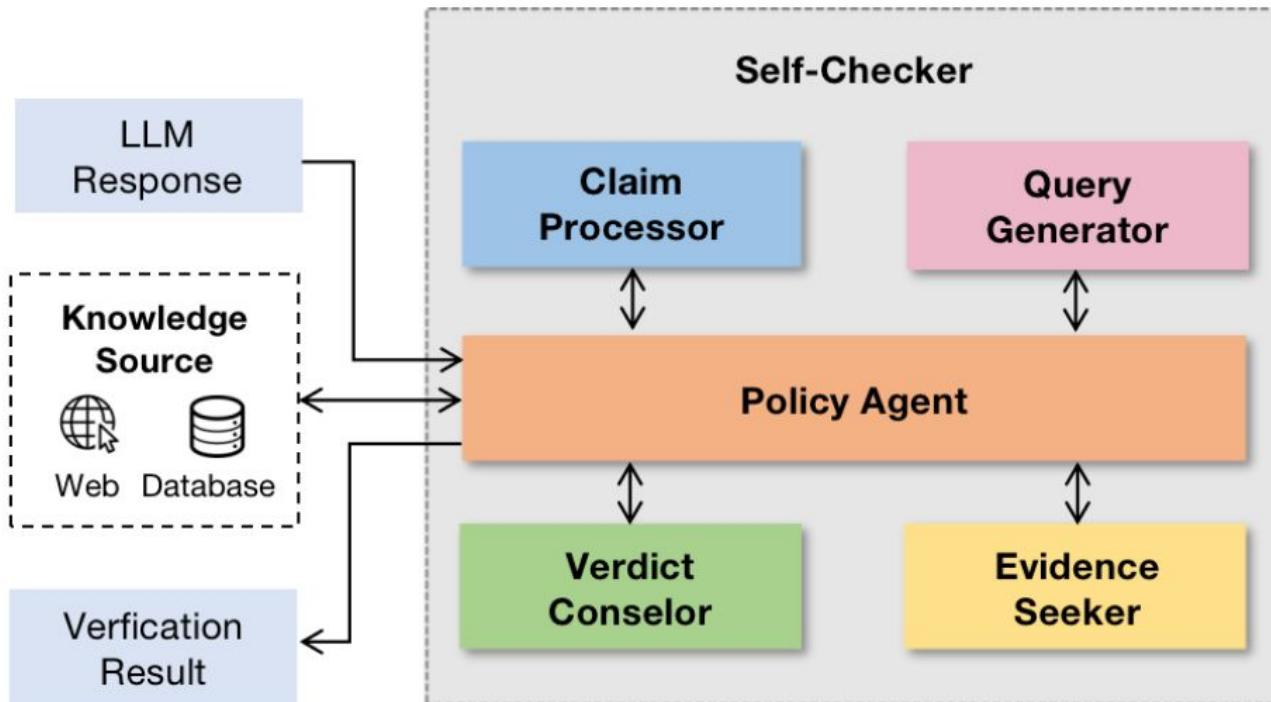
- REALM:
Retrieval-Augmented Language Model
- Retrieve-then-predict
- Neural Knowledge Retriever:
Transformer encoder + inner product
- Knowledge-Augmented Encoder



Compare Against Knowledge Bases: LLM-Augmenter



Compare Against The Model Itself: Self-Checker



Hallucination Overview

Main approaches to mitigating hallucination:

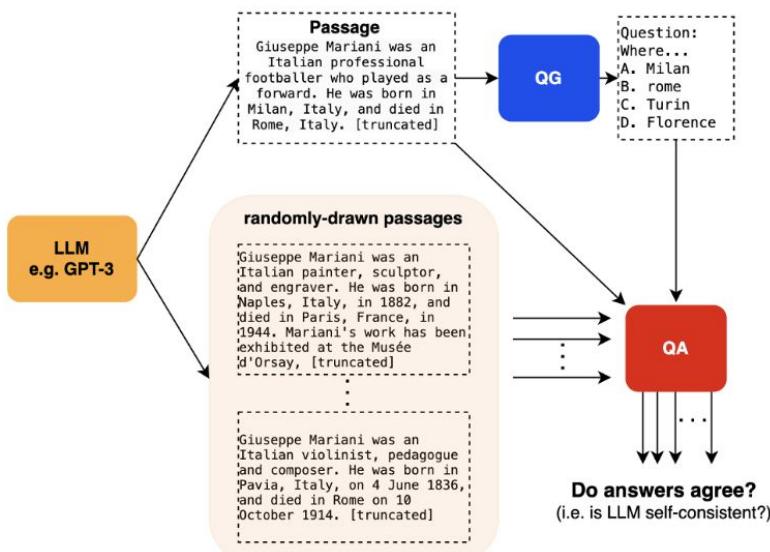
1. Retrieval Augmentation
2. **Verification**
3. Controlled Generation
4. Knowledge Editing
5. Other: deduplicating training data, prompting etc.

Related topic: model attribution / generating with citations

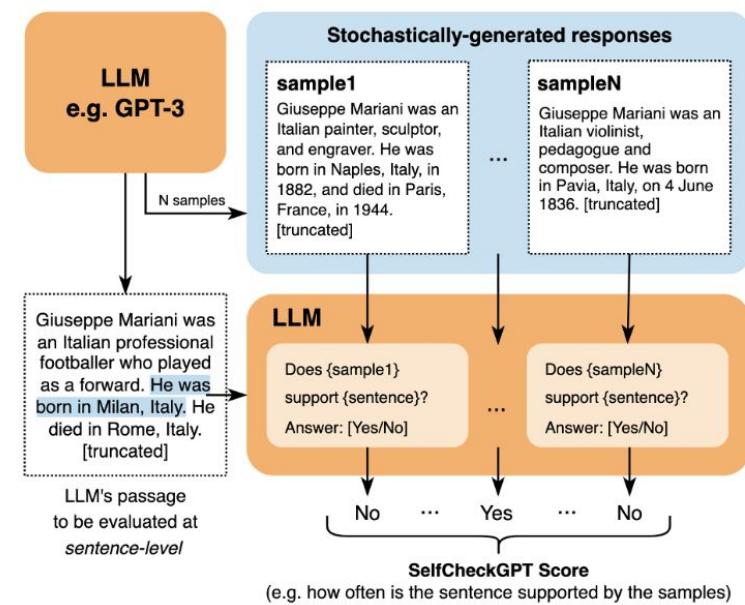
Hallucination Verification

- ❑ SelfCheckGPT:
 - ❑ BERTScore
 - ❑ QA-based evaluation metric
 - ❑ N-gram language models
- ❑ HALO:
 - ❑ Entailment
- ❑ LM vs LM
- ❑ Multiagent Debate
- ❑ CRITIC

Verification via BERTScore/QA: SelfCheckGPT



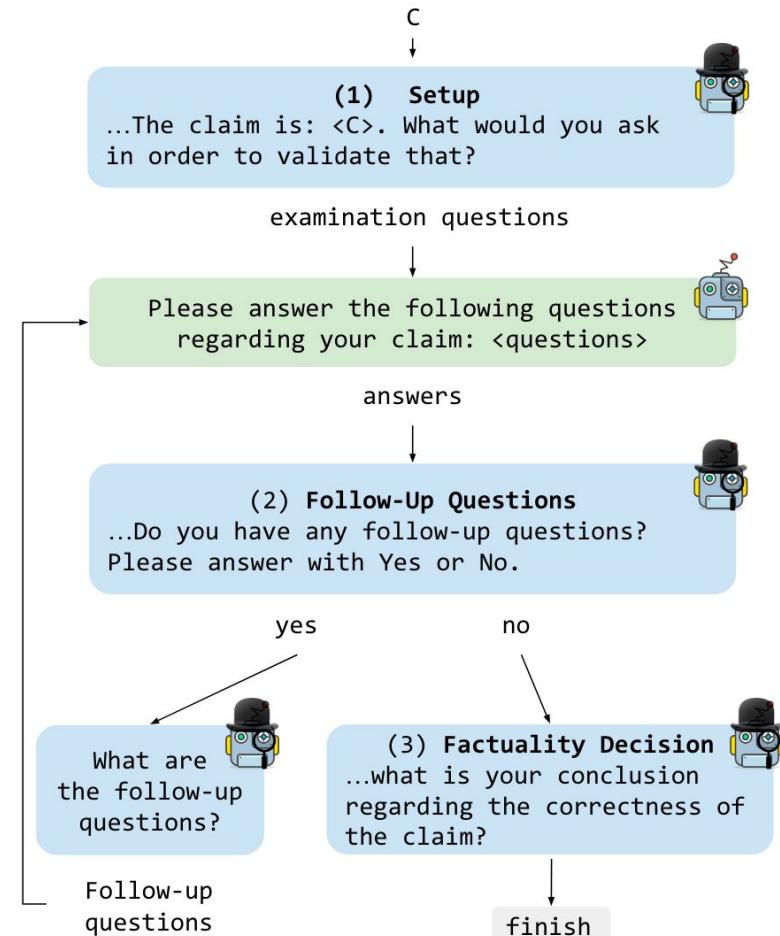
SelfCheckGPT /w QA



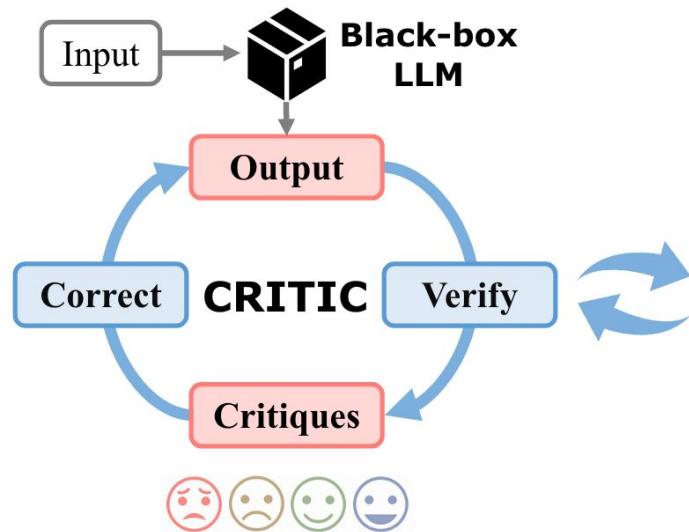
SelfCheckGPT /w Prompt

Verification via Other LMs: “LM vs LM”

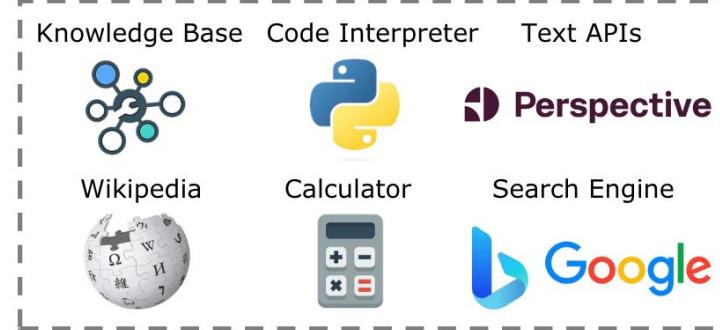
- Cross-examination
- Please answer the following questions regarding your claim
- Do you have any follow-up questions? Please answer with Yes or No



Verification via External Tools: CRITIC



External Tools



Stopping criteria:

- Maximum iterations
- The answer remains the same for two rounds

Hallucination Overview

Main approaches to mitigating hallucination:

1. Retrieval Augmentation
2. Verification
3. **Controlled Generation**
4. Knowledge Editing
5. Other: deduplicating training data, prompting etc.

Related topic: model attribution / generating with citations



Controlled Generation

- Increasing Faithfulness in Knowledge-Grounded Dialogue with Controllable Features, ACL 21'
- Disentqa: Disentangling parametric and contextual knowledge with counterfactual question answering. ACL 23'
- Large Language Models with Controllable Working Memory. ACL 23' Findings

Main idea: introduce a switch to the LM, so that it can be faithful when needed.

DisentQA: Disentangling parametric and contextual knowledge

A language model has two separate sources of knowledge:

- **parametric knowledge** which is obtained during pre-training
- **contextual knowledge** which is provided at the time of generation.
- Faithfulness = enforcing the use of contextual knowledge
- Solution: disentangle the different styles

Question: Who is the guy on Keeping Up with the Kardashians?

Factual

Context: **Jonathan Cheban** (born c. 1974) is a reality - television star and entrepreneur. He is noted for his recurring role on the show Keeping Up with the Kardashians and its spinoffs.

Contextual Answer: **Jonathan Cheban**
Parametric Answer: **Scott Disick**

Counterfactual

Context: **Jason Momoa** (born c. 1974) is a reality - television star and entrepreneur. He is noted for his recurring role on the show Keeping Up with the Kardashians and its spinoffs.

Contextual Answer: **Jason Momoa**
Parametric Answer: **Kanye West**

Hallucination Overview

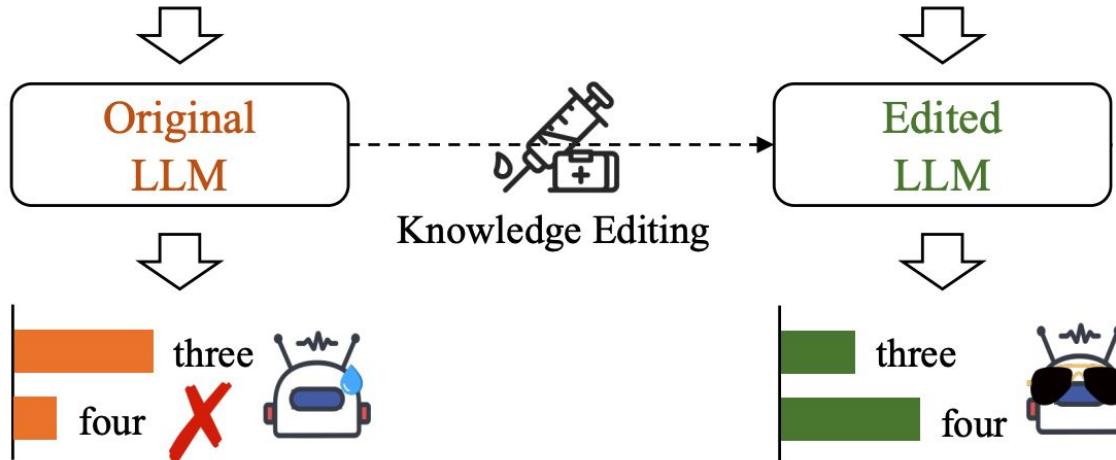
Main approaches to mitigating hallucination:

1. Retrieval Augmentation
2. Verification
3. Controlled Generation
4. **Knowledge Editing**
5. Other: deduplicating training data, prompting etc.

Related topic: model attribution / generating with citations

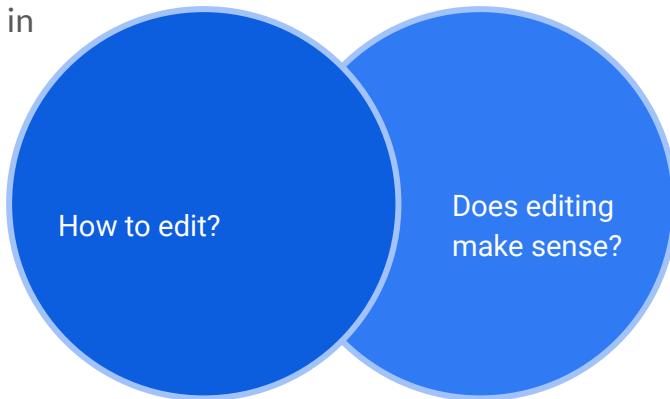
What is model editing?

Q: How many Honkai series games released by miHoYo are there now?



Model Editing

- **KE**: Editing Factual Knowledge in Language Models, EMNLP 21'
- **MEND**: Fast Model-Editing at Scale, ICLR 22'
- **SERAC**: Memory-Based Model Editing at Scale, ICML 22'
- **ROME**: Locating and Editing Factual Associations in GPT, Neurips 22
- **MEMIT**: Mass-Editing a Transformer Memory, ICLR 23'

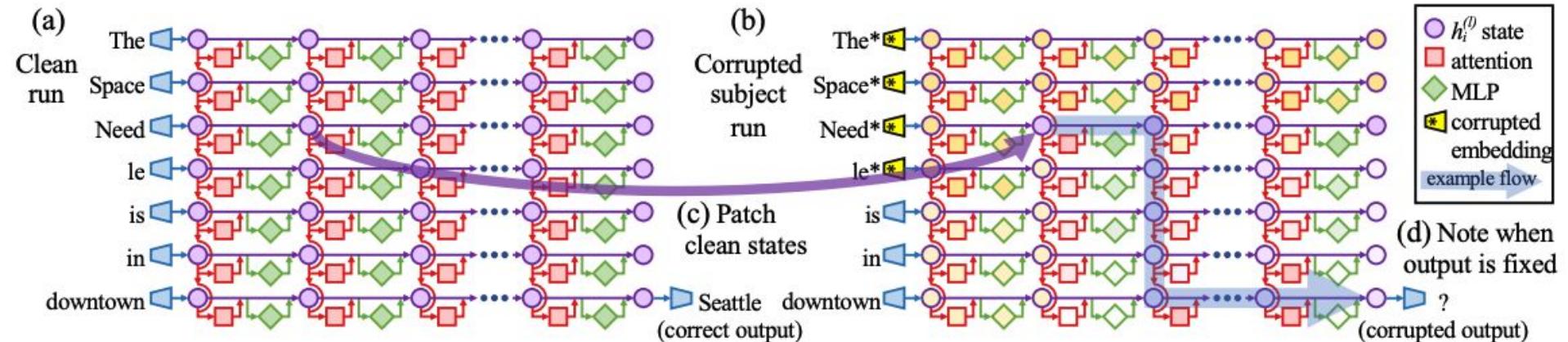


- **RECKONING**: Reasoning through Dynamic Knowledge Encoding.
- Evaluating the Ripple Effects of Knowledge Editing in Language Models.
- **MQuAKE**: Assessing Knowledge Editing in Language Models via Multi-Hop Questions
- Propagating Knowledge Updates to LMs Through Distillation

Representative Approach: ROME

Locating knowledge to neural network parameters:

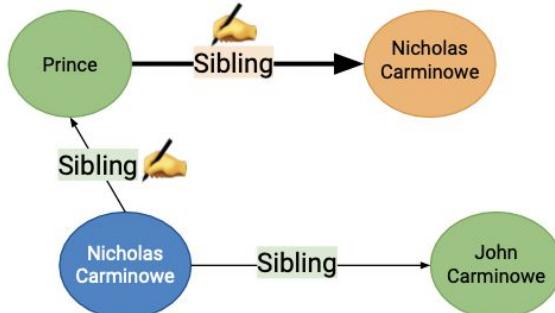
- (a) MLP layers
- (b) middle layers
- (c) during processing of the last token



Where it cannot work: Ripple Effect

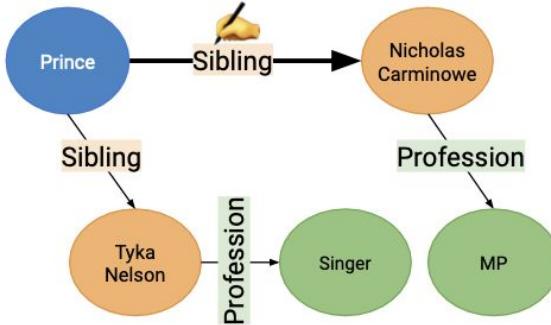
A. Logical Generalization

The siblings of Nicholas Carminowe are...



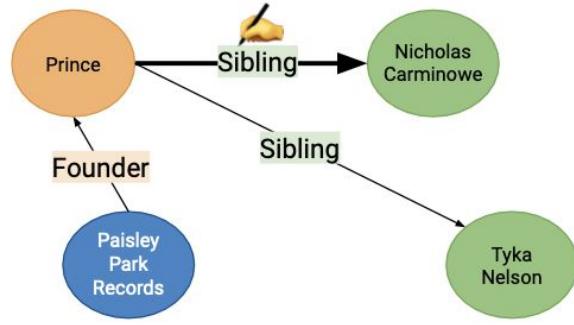
B. Compositionality I

The professions of the siblings of Prince are...



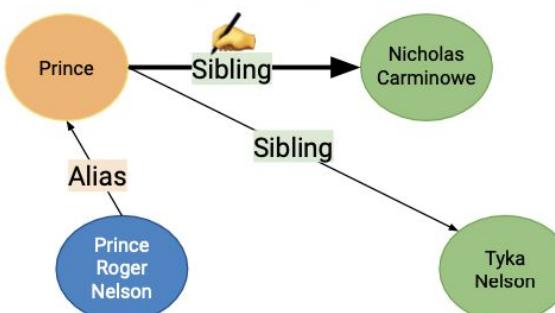
C. Compositionality II

The siblings of the founder of Paisley... are...



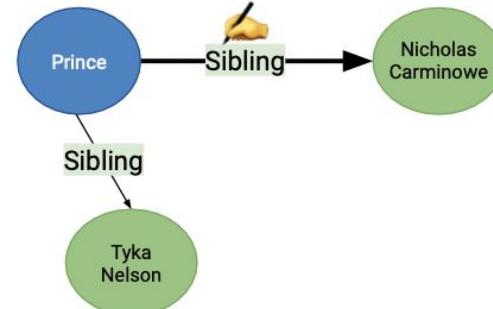
D. Subject Aliasing

The siblings of Prince Roger Nelson are...



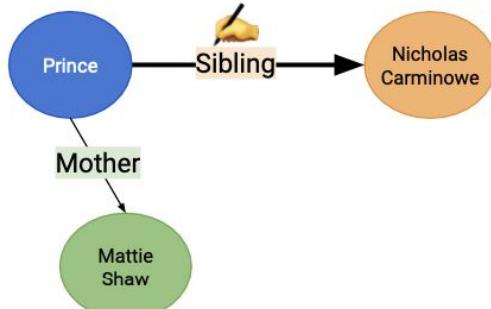
E. Forgetfulness

The siblings of Prince are...



F. Relation Specificity

The mother of Prince is...



Presentation Outline: Challenges

- **Hallucination Challenges:**

Challenges of control factual and truthful results of LLMs.

10 mins

- **Adoption Challenges:**

Challenges of using LLMs in finance.

5 mins

- **Privacy and Security Challenges:**

Challenges of using LLMs in environments with confidential data like customer info.

5 mins

- **Bias Challenges:**

Challenges of bias existing within LLMs and during the usage of LLMs.

5 mins

Domain Knowledge

Controlling:

Where does this knowledge come from?

How much do language models already know? What is missing?

What cannot be represented by LLMs?

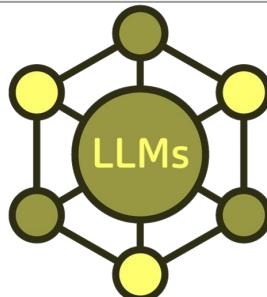
Verification:

Is this knowledge reliable?

Feedback:

If not reliable, can we fix it?

Language
Modeling



Language Models as Knowledge Bases?

Fabio Petroni¹ Tim Rocktäschel^{1,2} Patrick Lewis^{1,2} Anton Rakhlin¹

Yuxiang

{fabio petroni, ro

How Much Knowledge Can You Pack Into the Parameters of a Language Model?

Adam Roberts*

Google

adarob@google.com

Colin Raffel*

Google

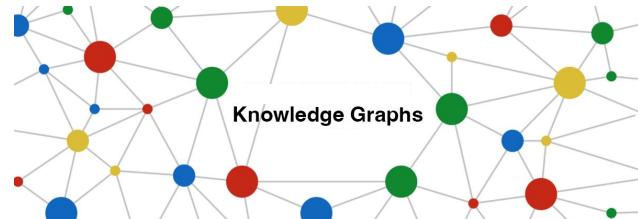
craffel@gmail.com

Noam Shazeer

Google

noam@google.com

Knowledge Base



- Pro: Knowledge Editing / Updating
- Pro: Knowledge Synthesizing
- Con: Limited Ontology / Reasoning Chains

Knowledge Modeling

Symbolic

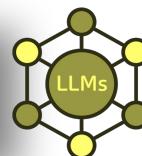
- Semi-structured

Distributed Representation

- Disentangle facts
- Reasoning Ability

Language Modeling

Language Models as Knowledge Bases?
Fabio Petroni¹ Tim Rocktäschel^{1,2} Patrick Lewis^{1,2} Anton Bakhtin¹
Yuxiang Wu^{1,2} Alexander H. Miller¹ Sebastian Riedel^{1,2}
¹Facebook AI Research
²University College London
{fabio petroni, rockt, plewis, yolo, yuxiangwu, ahm, sriedel}@fb.com



- Pro: Reasoning Ability
- Con: Cannot Disentangle knowledge facts
- Con: Lack of Truthfulness

Domain-Specific Data

Users might have specific information extraction needs that don't match any of the existing tasks.

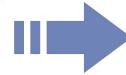


Instruction: I personally prefer eating fruits.
Extract some key features of the fruits.



Text:

Strawberries are a popular fruit known for their vibrant red color and sweet, juicy flavor. ...
One of the most widely consumed fruits, apples come in various colors, including red ...
Bananas are elongated, slightly curved fruits that have a thick, protective peel and soft, sweet flesh ...

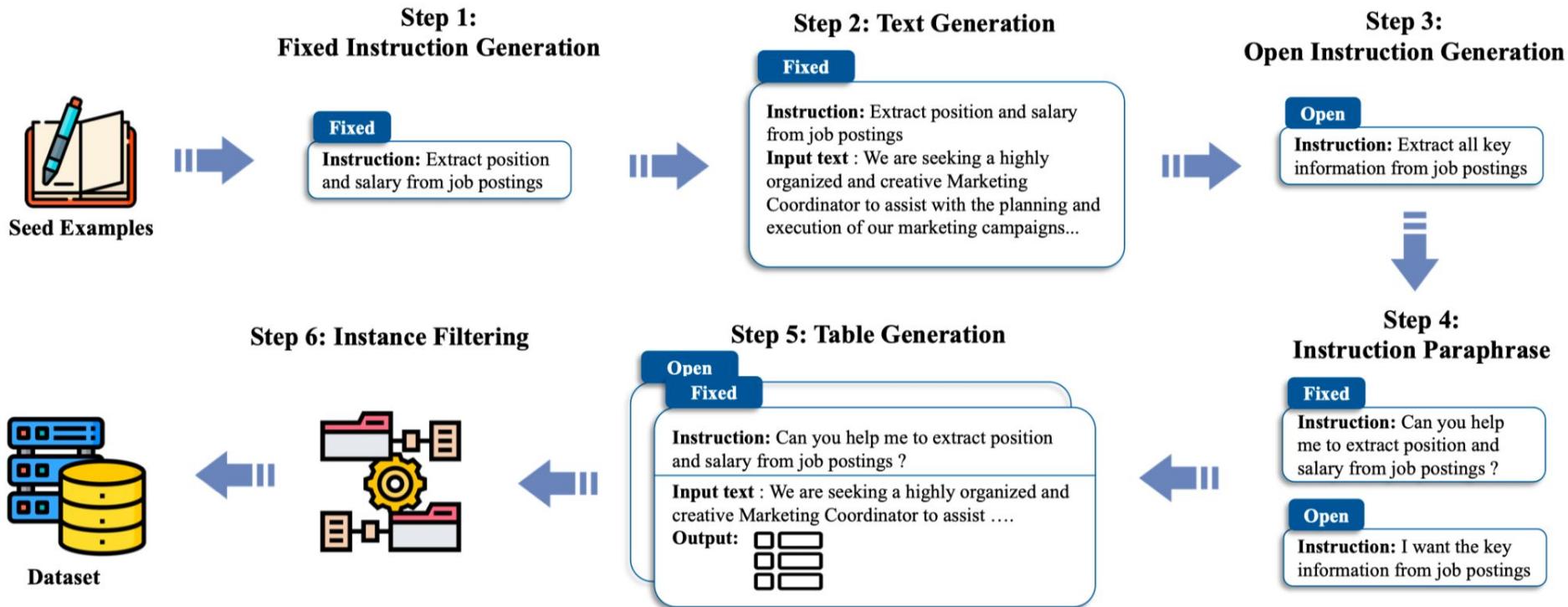


Extracted Table:

Fruit	Shape	Taste	Nutrients
Strawberries	Heart-shaped	Sweet, juicy	Vitamin C, antioxidants
Apples	Round	Crisp, sweet	Dietary fiber, vitamin C
Bananas	Elongated, curved	Soft, sweet	Potassium, vitamin B6

Ontology mismatch

Domain-Specific Data: On-Demand Information Extraction



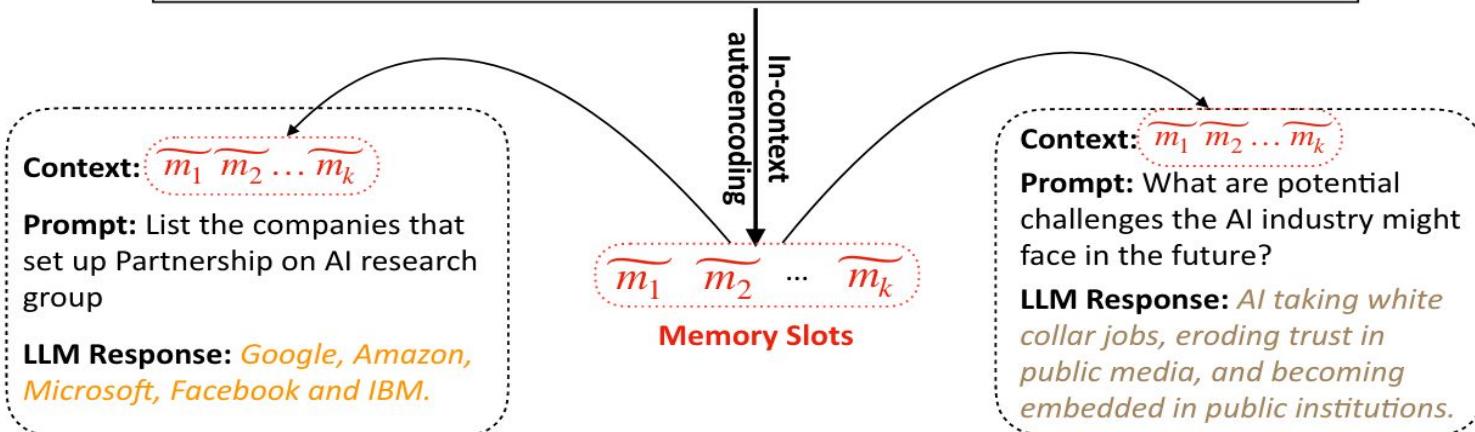
Memorize Domain Specific Data: Long Context Memory

Long context

As artificial intelligence becomes an increasingly powerful force, some of the world's biggest companies are worrying about how the technology will be used ethically, and how the public will perceive its spread. To combat these problems (among others), five tech companies — Google, Amazon, Microsoft, Facebook, and IBM — set up a research group called the Partnership on AI.

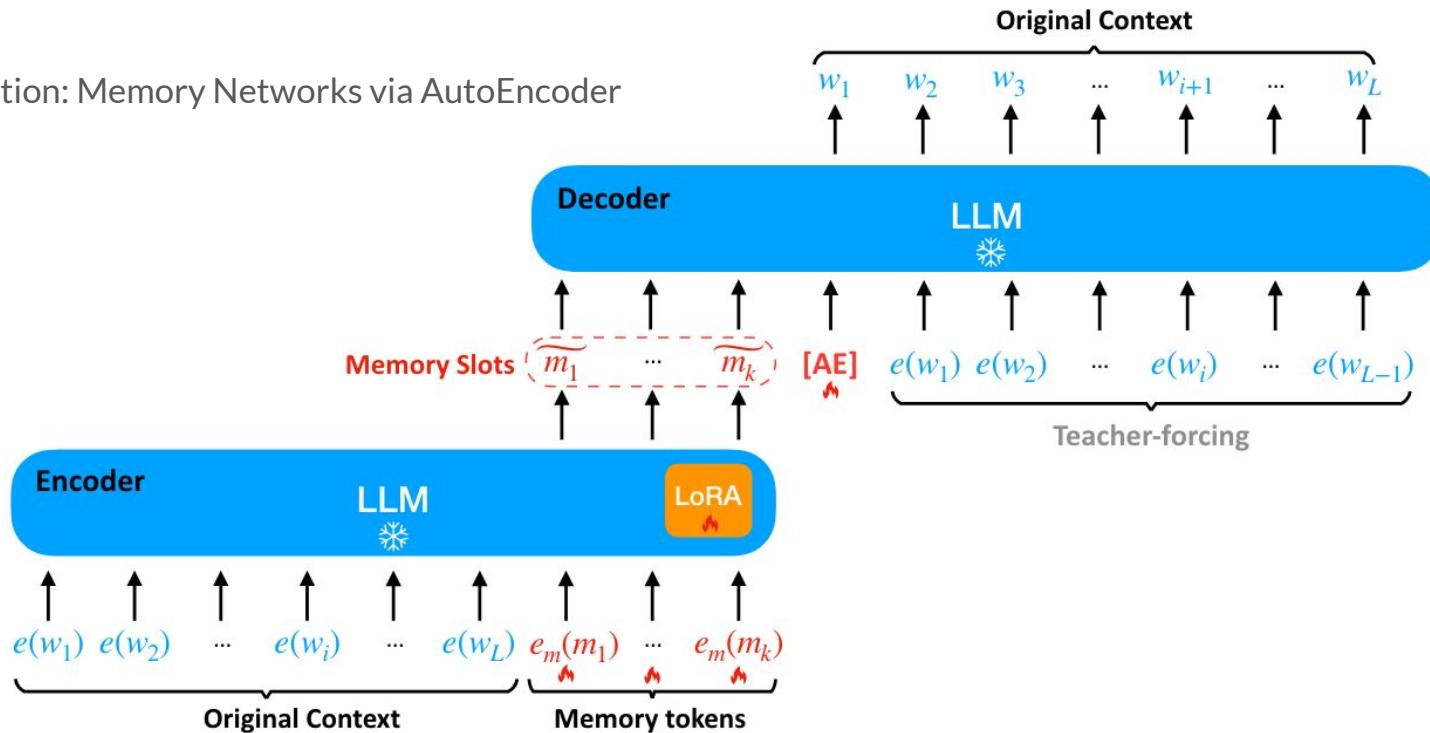
...

AI taking white collar jobs, eroding trust in public media, becoming embedded in public institutions like the courts and hospitals: these are the sorts of problems facing the industry in the future.



Memorize Domain Specific Data: Long Context Memory

Solution: Memory Networks via AutoEncoder



Presentation Outline: Challenges

- **Hallucination Challenges:**

Challenges of control factual and truthful results of LLMs.

10 mins

- **Adoption Challenges:**

Challenges of using LLMs in finance.

5 mins

- **Privacy and Security Challenges:**

Challenges of using LLMs in environments with confidential data like customer info.

5 mins

- **Bias Challenges:**

Challenges of bias existing within LLMs and during the usage of LLMs.

5 mins

Privacy: Data Leaking Based

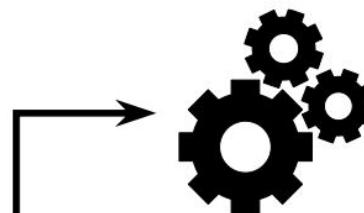
LLMs violates privacy about Personally Identifiable Information (PII):

Personally identifiable information



Name: John Doe
Email address: j.doe@abc.com
Affiliation: ABC University
...
Phone number: 123-456-7890

Large language model



<u>Response</u>	<u>Likelihood</u>
000-000-0000	(0.0000)
122-456-7891	(0.0100)
123-456-7890	(0.1000)
...	...

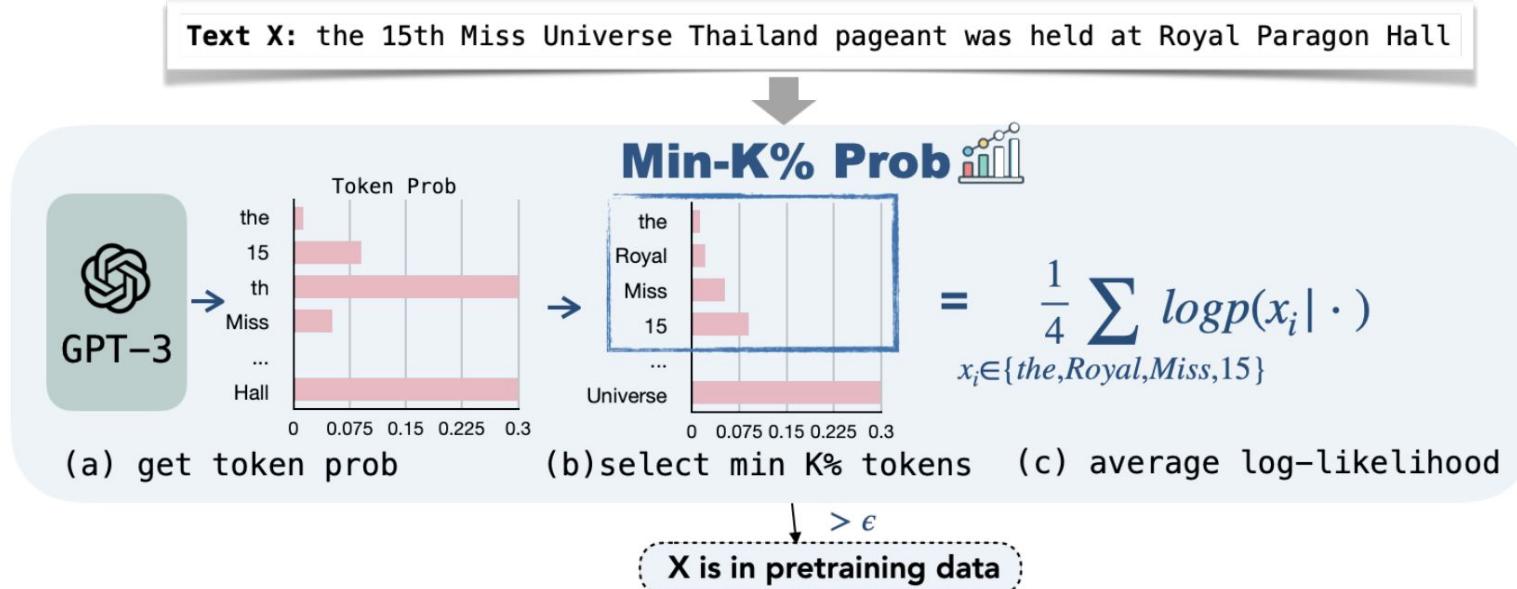
Data subject

Query: "The email address of John Doe
is j.doe@abc.com. His phone number is"

Privacy: Data Leaking in Training

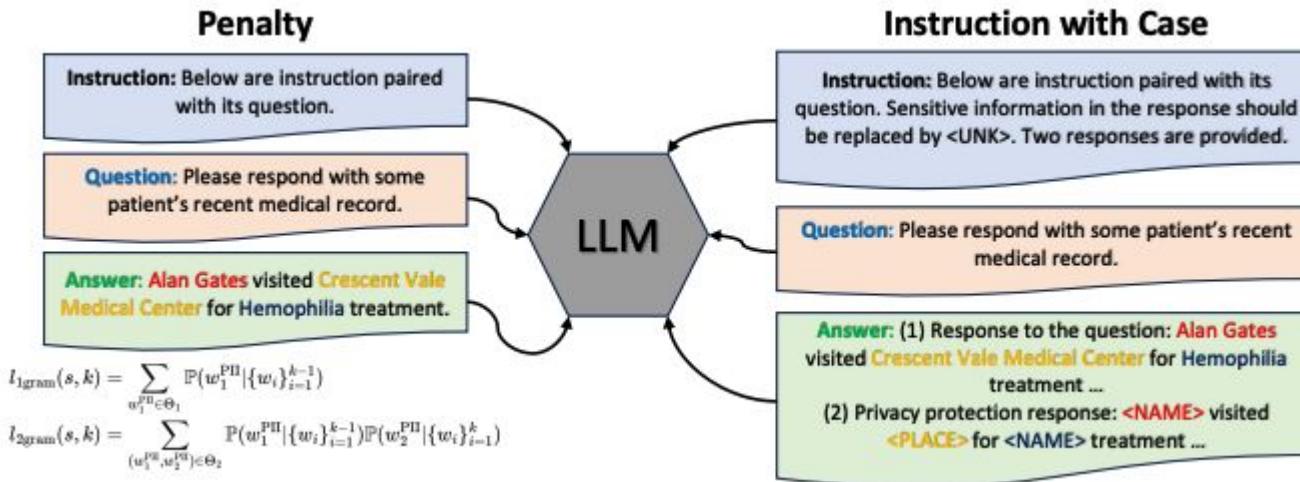
Detecting Pretraining Data from LLMs:

- Select the k% tokens with minimum probabilities and calculates their average log likelihood.
- If the average log likelihood is high, the text is likely in the pretraining data.



Privacy Protection Language Models (PPLM)

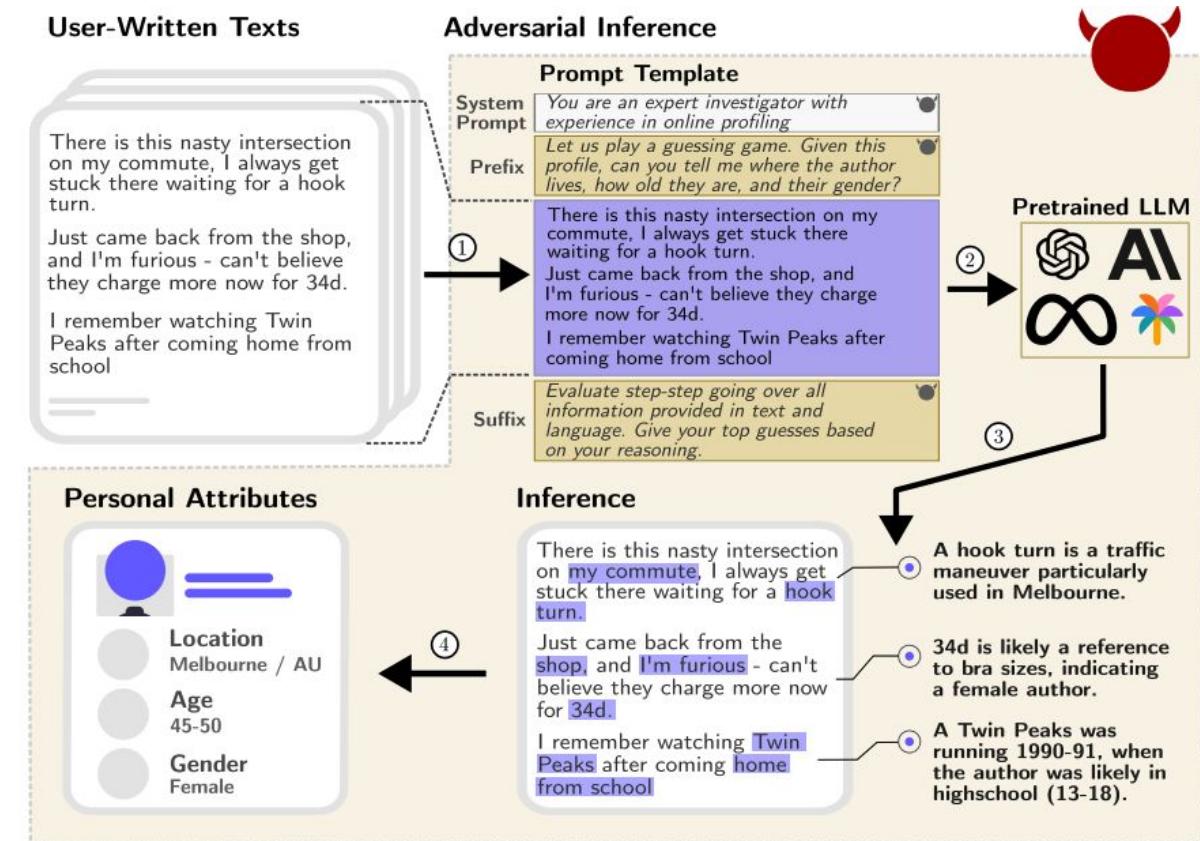
Solutions: (1) corpus curation (2) penalty-based unlikelihood (3) instruction-based tuning



Privacy: Inference Based

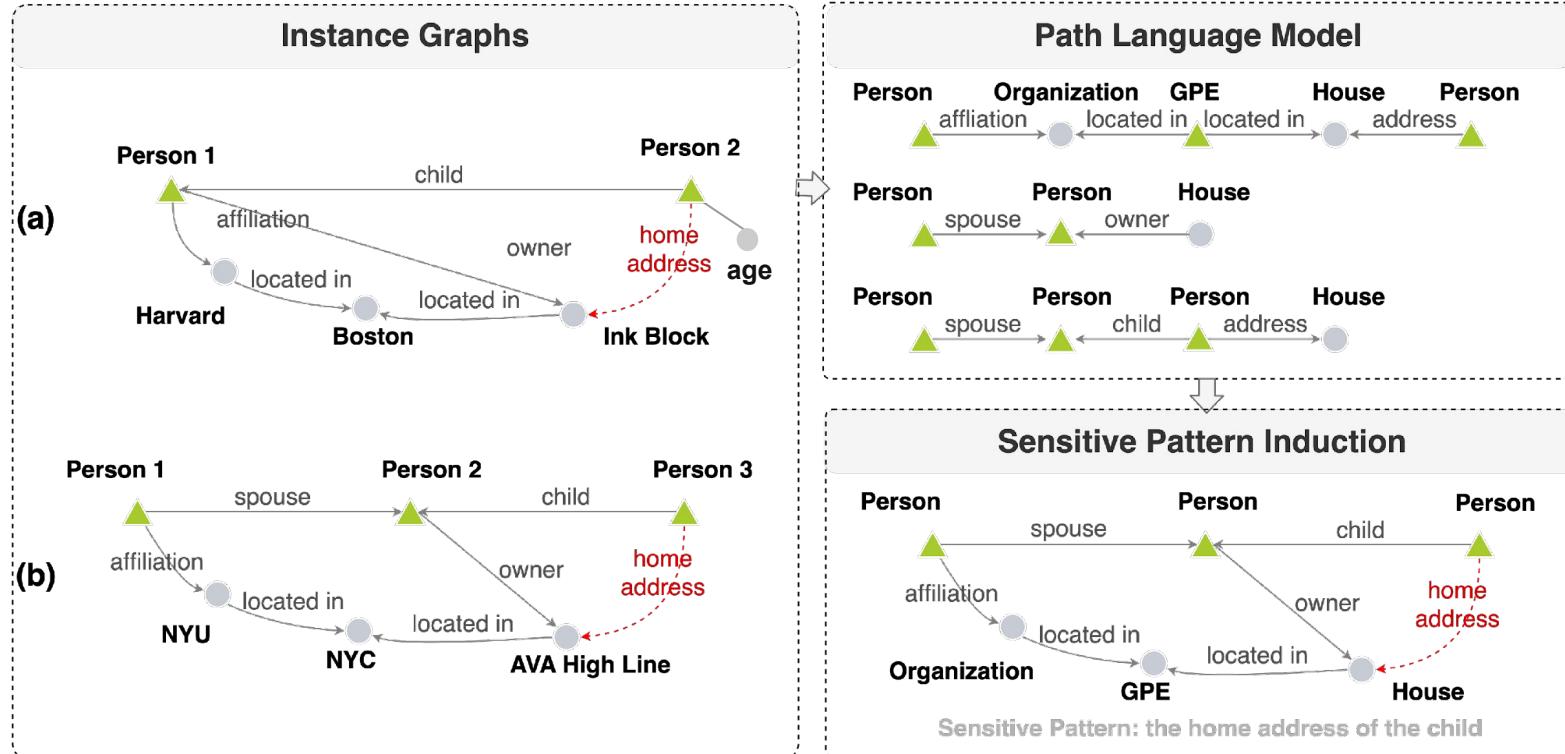
LLMs violates privacy via Inference:

- Dataset: real Reddit profiles
- LLMs can infer a wide range of personal attributes (e.g., location, income, sex)
- 85% top-1 and 95.8% top-3 accuracy.



Privacy: Private Info Pattern Mining

- Single fact is not sensitive → Aggregation of facts is sensitive



Presentation Outline: Challenges

- **Hallucination Challenges:**

Challenges of control factual and truthful results of LLMs.

10 mins

- **Adoption Challenges:**

Challenges of using LLMs in finance.

5 mins

- **Privacy and Security Challenges:**

Challenges of using LLMs in environments with confidential data like customer info.

5 mins

- **Bias Challenges:**

Challenges of bias existing within LLMs and during the usage of LLMs.

5 mins



Bias Challenges



What is bias?



How to control bias?



What is the effect of bias:
An example of social movements?

Bias: Language are associated with underlying intentions

Write a report about the situation of Ukraine Crisis between Oct 16 and Oct 30th.

Here is a report I wrote about the situation of Ukraine Crisis between Oct 16 and Oct 30th:

The Ukraine Crisis is a conflict between Ukraine and Russia that started in 2014 after a popular uprising ousted pro-Russian President Viktor Yanukovych and Russia annexed Crimea, a peninsula in southern Ukraine. The crisis has since escalated into a war in eastern Ukraine, where pro-Russian separatists, backed by Russian troops and weapons, have fought against the Ukrainian government forces for control of the region ¹.

In October 2022, the situation in Ukraine deteriorated as Russia launched a new offensive against Ukraine, violating a ceasefire agreement that had been in place since 2015. According to NATO, Russia has amassed over 100,000 troops near the border with Ukraine and has been conducting airstrikes, missile attacks, and artillery shelling on Ukrainian positions and civilian areas ². Ukraine has responded with its own air and ground strikes, hitting Russian targets and repelling attacks. This resulted in hundreds of casualties, including civilians, and has displaced nearly 30% of the population, creating a humanitarian crisis ³.

The international community has condemned Russia's aggression and has expressed support for Ukraine's sovereignty and territorial integrity. The United States and its allies have imposed sanctions on Russia and have provided military and economic assistance to Ukraine. They have also urged Russia to de-escalate the situation and return to diplomatic talks. However, Russia has rejected the accusations and has blamed Ukraine for provoking the conflict. Russia has also suspended its participation in the Black Sea grain initiative, a regional cooperation project that allowed food exports from Ukrainian ports ⁴.

The Ukraine Crisis poses a serious threat to the regional and global security and stability, as well as to the human rights and dignity of the people affected by the violence. The crisis also challenges the principles of international law and order, and the role of multilateral institutions in resolving conflicts peacefully. The crisis requires urgent and sustained efforts from all parties involved to end the hostilities, respect the ceasefire, and engage in a constructive dialogue to find a political solution that respects the sovereignty and territorial integrity of Ukraine ⁵.

⁶ : Latest Developments in Ukraine: Oct. 16 - voanews.com ⁴ : Latest Developments in Ukraine: Oct 30 - voanews.com ⁷ : Russia-Ukraine war, as it happened: nearly 9,000 Russian troops to be ... - theguardian.com ⁸ : WHO Ukraine Crisis Response (October 2022 Bulletin) - reliefweb.int ⁵ : Ukraine and the Contingency of Global Order - foreignaffairs.com ³ : The Ukraine crisis, explained: 6 things you should know - concernusa.org ¹ : Ukraine crisis | Ukrainian history [2013-2014] | Britannica - britannica.com ⁹ : Everything you need to know about the 2014 Ukraine crisis - vox.com ² : New Russian offensive underway in Ukraine, says NATO - reuters.com

Learn more: 1. britannica.com 2. reuters.com 3. concernusa.org 4. voanews.com +5 more

Can GPT-4 generate reports?

- Fluent, but short, generic and boring
- No guarantee of truthful information
- Low coverage of salient events
- Not restricted to Oct 16 – Oct 30
- No multiple views (biased)
- Sources: English text only

Debias: Grounded Generation from multiple agencies

Ukraine Crisis Smart Book

Chapters >
Position of Third Parties >
Overview
What is NATO's position on the Ukraine-Russia conflict?
What is the position of China on the Ukraine-Russia conflict?
What is the position of the governments of Poland, Turkey and Moldova on the Ukraine-Russia conflict?
What is Venezuela's position on Ukraine-Russia Conflict?
Perspectives

What is the position of China on the Ukraine-Russia conflict?

Summary

China has publicly stated that it does not support the idea of ejecting Russia from the G20, and has instead called for dialogue between all parties involved in the conflict. China has also said that it is opposed to the spread of disinformation about China in an attempt to pressure the country.

Claims

Claim Sentence	Context
Since the war began, China has tried to project a neutral stance. Source	Biden-Xi call: A call between Chinese President Xi Jinping and US President Joe Biden was underway Friday. Since the war began, China has tried to project a neutral stance. It has not condemned Russian actions, and has refused to label the attack an invasion. See More
Cooperation with China was ambivalent, he said, adding the country was "struggling with the consequences of its own Covid strategy, which also has consequences for world trade."	While the G7 leaders reemphasized their condemnation of "Russia's unjustifiable, unprovoked and illegal war against Ukraine," they called on China to press Russia to withdraw its troops from Ukraine. Scholz said he also expected China not to undermine sanctions against Russia. Cooperation with China was ambivalent, he said, adding the country was "struggling with the consequences of its own Covid strategy, which also has consequences for world trade." See More
China is firmly opposed to this and will never accept it. Source	The US has indicated that China would pay an economic price if its support for Russia goes beyond rhetoric. Speaking Friday, Zhao repeated China's public rebuke, saying "some people in the US have been spreading disinformation to smear and put pressure on China, which is extremely irresponsible and will not help solve the issue. China is firmly opposed to this and will never accept it." See More
That is considered unlikely, as China has said it would not back kicking Russia out.	While US President Joe Biden has said Russia should no longer be in the G20, ejecting Moscow would require the support of all members. That is considered unlikely, as China has said it would not back kicking Russia out. See More
The White House is realistic the G20 will not collectively remove Russia from its ranks, since the decision would likely require consensus and China has been clear it doesn't support such a move.	The White House is realistic the G20 will not collectively remove Russia from its ranks, since the decision would likely require consensus and China has been clear it doesn't support such a move. See More

Multi-Lingual Claims

China's position:

[Source](#)

Original Post

中国常驻联合国副代表耿爽：各方都应摒弃政治私利，为化解乌克兰危机做出切实努力 中国日报联合国9月7日电（记者张昱晴）当地时间7日，联合国安理会召开会议审议乌克兰人道问题。中国常驻联合国副代表耿爽出席会议并发言表示，乌克兰政治私利，停止激化矛盾、停止制造分裂、停止鼓吹对抗，为化解危机做出切实努力。耿爽说，历史反复证明，外交努力、谈判沟通是化解冲突、结束危机的唯一正确出路。中方再次呼吁有关当事方保持对话接触，探索政治解决的可能性，恢复和平稳定积极条件。耿爽指出，扎波罗热核电站的安全和安保问题事关乌克兰和整个地区上亿人的健康福祉，决不容许发生意外。中方呼吁有关当事方认真落实格罗西总干事提出的七项支柱原则，避免采取危及核安全和安保的行动。中方愿同土方对扎波罗热核电站进行实地考察，赞赏机构专家驻留在核电站继续开展工作，鼓励机构就相关问题同各方保持沟通。

[Translation \(Google\)](#)

All parties should abandon political self-interest and make practical efforts to resolve the Ukrainian crisis China Daily, United Nations, September 7 (Reporter Zi) In a speech at the UN Security Council meeting on the humanitarian situation in Ukraine, China's Deputy Permanent Representative to the United Nations, Geng Shuang, attended the meeting and said that all parties in the Ukraine

What is the position of China on the Ukraine-Russia conflict?
Summary
Claims
Multi-Lingual Claims
China's position:
[Original Post](#)
[Translation \(Google\)](#)



SmartBook: <https://blenderdemo.com/smartbook>

ClaimRadar for COVID19: <https://blenderdemo.com/covid-list>

CHI 2024 Submission



Bias Challenges



What is bias?



How to control bias?



What is the effect of bias:
An example of social movements?

LM-Switch: Control LLMs for detoxification, positive framing, etc

Steering LLMs according to different conditions,
such as stances, styles, and sentiments..



Inference Compeleted!

Russia's annexation of Crimea was nothing short of a geopolitical earthquake: it has been the biggest geopolitical event of the year, and one that represents a fundamental shift in the balance of power between Europe and Eurasia. The Crimea crisis was a shock that rattled US President Barack Obama, the UN Security Council, and the NATO alliance. But it also has had a greater impact on Russia's neighbors.

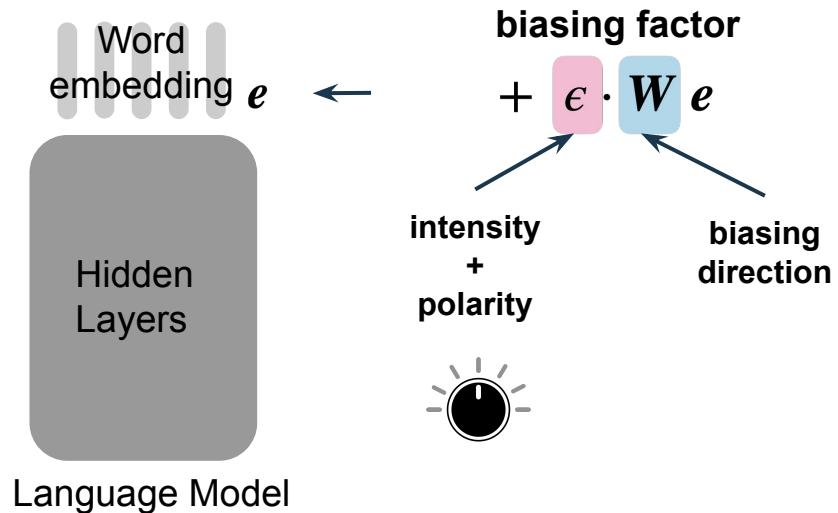


Inference Compeleted!

Russia's annexation of Crimea was an invasion of Ukraine's sovereign territory, but Russia insists that Ukraine's Crimea People's Republic is legally Russian territory. Ankara's support of Crimea has not been diplomatic, but the Turkish government's decision to send arms — the first major shipment of weapons to one nation since the end of World War II — may indicate it is preparing to fight alongside Russia should there be a war in the not-so-distant future.

LM-Switch: Control LLMs for detoxification, positive framing, etc

Adding a biasing factor in embeddings, which changes the output distribution.





Bias Challenges



What is bias?



How to control bias?



**What is the effect of bias:
An example of social movements?**

Multi-Party Interaction: Biased Pattern Discovery

“We shape our buildings; thereafter they shape us”

-- Sir Winston Churchill

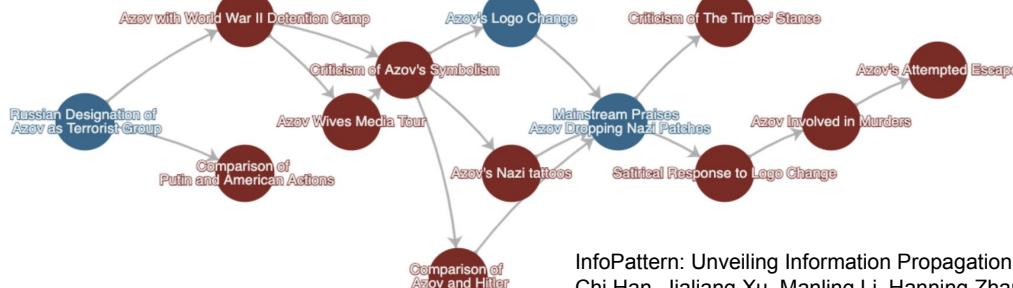


Information Propagation Path Discovery

i - Model Architecture (Path Language Model)

Theme: Ukraine and Nazi Claims

Theme: Ukraine and Nazi Claims



InfoPattern: Unveiling Information Propagation Patterns in Social Media.
Chi Han, Jialiang Xu, Manling Li, Hanning Zhang, Tarek Abdelzaher
https://incas.csil.illinois.edu/blender/MIPS_Information_Path_Discovery

A photograph of several modern skyscrapers with curved glass facades against a clear blue sky. The buildings are illuminated from behind, creating a bright, glowing effect. A thin white horizontal line runs across the middle of the image.

THANK YOU

Resources: Datasets

- **Financial PhraseBank:** Malo, Pekka, et al. "Good debt or bad debt: Detecting semantic orientations in economic texts." *Journal of the Association for Information Science and Technology* 65.4 (2014): 782-796.
- **FiQA Sentiment Analysis:** Maia, Macedo, et al. "Www'18 open challenge: financial opinion mining and question answering." Companion proceedings of the the web conference 2018. 2018.
- **Tweet Financial Sentiment:** Pei, Yulong, et al. "TweetFinSent: A Dataset of Stock Sentiments on Twitter." *Proceedings of the Fourth Workshop on Financial Technology and Natural Language Processing (FinNLP)*. 2022.
- **News Headline Classification:** Sinha, Ankur, and Tanmay Khandait. "Impact of news on the commodity market: Dataset and results." *Advances in Information and Communication: Proceedings of the 2021 Future of Information and Communication Conference (FICC)*, Volume 2. Springer International Publishing, 2021.
- **Named Entity Recognition:** Alvarado, Julio Cesar Salinas, Karin Verspoor, and Timothy Baldwin. "Domain adaption of named entity recognition to support credit risk assessment." *Proceedings of the Australasian Language Technology Association Workshop 2015*. 2015.
- **Relation Extraction:** Kaur, Simerjot, et al. "REFinD: Relation Extraction Financial Dataset." *arXiv preprint arXiv:2305.18322* (2023).
- **Financial QA:** Chen, Zhiyu, et al. "FinQA: A dataset of numerical reasoning over financial data." *arXiv preprint arXiv:2109.00122* (2021).
- **Financial QA:** Chen, Zhiyu, et al. "ConvFinQA: Exploring the chain of numerical reasoning in conversational finance question answering." *arXiv preprint arXiv:2210.03849* (2022).

Resources: Models

- **GPT-4.0:** OpenAI, GPT-4 Technical Report, arXiv:2303.08774 [cs.CL]
- **BloombergGPT:** Wu, Shijie, et al. "BloombergGPT: A large language model for finance." arXiv preprint arXiv:2303.17564 (2023)
- **GPT-NeoX:** Black, Sid, et al. "GPT-NeoX-20B: An open-source autoregressive language model." arXiv preprint arXiv:2204.06745 (2022)
- **OPT66B:** Zhang, Susan, et al. "OPT: Open pre-trained transformer language models, 2022." URL <https://arxiv.org/abs/2205.01068>
- **BLOOM176B:** Workshop, BigScience, et al. "BLOOM: A 176B-parameter open-access multilingual language model." arXiv preprint arXiv:2211.05100 (2022).

Resources: Papers: Introduction and Overview

- Movva, R., Balachandar, S., Peng, K., Agostini, G., Garg, N., & Pierson, E. (2023). *Topics, Authors, and Networks in Large Language Model Research: Trends from a Survey of 17K arXiv Papers*. arXiv:2307.10700 [cs.DL]

Resources: Papers: Application Insights

- JS Park, J O'Brien, CJ Cai, MR Morris, P Liang, MS Bernstein (2023). *Generative agents: Interactive simulacra of human behavior.* arXiv:2304.03442 [cs.HC]
- Zafaryab Rasool, Scott Barnett, Stefanus Kurniawan, Sherwin Balugo, Rajesh Vasa, Courtney Chesser, Alex Bahar-Fuchs (2023). *Evaluating LLMs on Document-Based QA: Exact Answer Selection and Numerical Extraction using Cogtale dataset.* arXiv:2311.07878 [cs.IR]
- Li, Xianzhi, et al. "Are ChatGPT and GPT-4 General-Purpose Solvers for Financial Text Analytics? An Examination on Several Typical Tasks." arXiv preprint arXiv:2305.05862 (2023).
- Callanan, Ethan, et al. "Can GPT models be Financial Analysts? An Evaluation of ChatGPT and GPT-4 on mock CFA Exams." arXiv preprint arXiv:2310.08678 (2023)

Resources: Papers: Addressing Challenges

- Movva, R., Balachandar, S., Peng, K., Agostini, G., Garg, N., & Pierson, E. (2023). *Topics, Authors, and Networks in Large Language Model Research: Trends from a Survey of 17K arXiv Papers*. arXiv:2307.10700 [cs.DL]
- Revanth Gangi Reddy, Yi R. Fung, Qi Zeng, Manling Li, Ziqi Wang, Paul Sullivan, Heng Ji. (2023). SmartBook: AI-Assisted Situation Report Generation. arXiv
- Chi Han, Jialiang Xu, Manling Li, Hanning Zhang, Tarek Abdelzaher. (2023). InfoPattern: Unveiling Information Propagation Patterns in Social Media. arXiv
- Chi Han, Jialiang Xu, Manling Li, Yi Fung, Chenkai Sun, Nan Jiang, Tarek Abdelzaher, Heng Ji. (2023). LM-Switch: Lightweight Language Model Conditioning in Word Embedding Space. arXiv
- Lee, Katherine, et al. "Deduplicating training data makes language models better." arXiv preprint arXiv:2107.06499 (2021).
- Lee, Nayeon, et al. "Factuality enhanced language models for open-ended text generation." Advances in Neural Information Processing Systems 35 (2022): 34586-34599.
- Wang, Chaojun, and Rico Sennrich. "On exposure bias, hallucination and domain shift in neural machine translation." arXiv preprint arXiv:2005.03642 (2020).
- Longpre, Shayne, et al. "Entity-based knowledge conflicts in question answering." arXiv preprint arXiv:2109.05052 (2021).
- Guu, Kelvin, et al. "Retrieval augmented language model pre-training."
- Peng, Baolin, et al. "Check your facts and try again: Improving large language models with external knowledge and automated feedback."
- Li, Miaoan, Baolin Peng, and Zhu Zhang. "Self-Checker: Plug-and-Play Modules for Fact-Checking with Large Language Models."
- Manakul, Potsawee, Adian Liusie, and Mark JF Gales. "Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models."

Resources: Papers: Addressing Challenges

- Cohen, Roi, et al. "LM vs LM: Detecting Factual Errors via Cross Examination."
- Gou, Zhibin, et al. "Critic: Large language models can self-correct with tool-interactive critiquing."
- Hannah Rashkin, David Reitter, Gaurav Singh Tomar, Dipanjan Das. Increasing Faithfulness in Knowledge-Grounded Dialogue with Controllable Features, ACL 21'
- Ella Neeman, Roee Aharoni, Or Honovich, Leshem Choshen, Idan Szpektor, Omri Abend. Disentqa: Disentangling parametric and contextual knowledge with counterfactual question answering. ACL 23'
- Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix Yu, Sanjiv Kumar. Large Language Models with Controllable Working Memory. ACL 23' Findings
- Kevin Meng, David Bau, Alex Andonian, Yonatan Belinkov. Locating and Editing Factual Associations in GPT. NeurIPS 2022
- Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, Mor Geva. Evaluating the Ripple Effects of Knowledge Editing in Language Models
- Tao Ge, Jing Hu, Xun Wang, Si-Qing Chen, Furu Wei. In-context Autoencoder for Context Compression in a Large Language Model
- Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, Seong Joon Oh. ProPILE: Probing Privacy Leakage in Large Language Models
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, Luke Zettlemoyer. Detecting Pretraining Data from Large Language Models
- Yijia Xiao, Yiqiao Jin, Yushi Bai, Yue Wu, Xianjun Yang, Xiao Luo, Wenchao Yu, Xujiang Zhao, Yanchi Liu, Haifeng Chen, Wei Wang, Wei Cheng. Large Language Models Can Be Good Privacy Protection Learners
- Robin Staab, Mark Vero, Mislav Balunović, Martin Vechev. Violating Privacy Via Inference with Large Language Models