

[Week 9] Challenges with LLMs

ETMI5: Explain to Me in 5

In this section of the course on LLM Challenges, we've identified two main areas of concern with LLMs: behavioral challenges and deployment challenges. Behavioral challenges include issues like hallucination, where LLMs generate fictitious information, and adversarial attacks, where inputs are crafted to manipulate model behavior. Deployment challenges encompass memory and scalability issues, as well as security and privacy concerns. LLMs demand significant computational resources for deployment and face risks of privacy breaches due to their ability to process vast datasets and generate text. To mitigate these challenges, we discuss various strategies such as robust defenses against adversarial attacks, efficient memory management, and privacy-preserving training algorithms. Additionally we will go over techniques like differential privacy, model stacking, and preprocessing methods that are being employed to safeguard user privacy and ensure the reliable and ethical use of LLMs across different applications.

Types of Challenges

We categorize the challenges into two main areas: managing the behavior of LLMs and the technical difficulties encountered during their deployment. Given the evolving nature of this technology, it's likely that current challenges will be mitigated, and new ones may emerge over time. However, as of February 15, 2024, these are the prominently discussed challenges associated with LLMs:

A. Behavioral Challenges

1. Hallucination

LLMs sometimes generate plausible but entirely fictitious information or responses, known as "hallucinations." This challenge is particularly harmful in applications requiring high factual accuracy, such as news generation, educational content, or medical advice as hallucinations can erode trust in LLM outputs, leading to misinformation or potentially harmful advice being followed.

2. Adversarial Attacks

LLMs can be vulnerable to adversarial attacks, where inputs are specially crafted to trick the model into making errors or revealing sensitive information. These attacks can compromise the integrity and reliability of LLM applications, posing significant security risks.

3. Alignment

Ensuring LLMs align with human values and intentions is a complex task. Misalignment can result from the model pursuing objectives that don't fully encapsulate the user's goals or

ethical standards. Misalignment can lead to undesirable outcomes, such as generating content that is offensive, biased, or ethically questionable.

4. Prompt Brittleness

LLMs can be overly sensitive to the exact wording of prompts, leading to inconsistent or unpredictable outputs. Small changes in prompt structure can yield vastly different responses. This brittleness complicates the development of reliable applications and requires users to have a deep understanding of how to effectively interact with LLMs.

B. Deployment Challenges

1. Memory and Scalability Challenges

Deploying LLMs at scale involves significant memory and computational resource demands. Managing these resources efficiently while maintaining high performance and low latency is a technical hurdle. Scalability challenges can limit the ability of LLMs to be integrated into real-time or resource-constrained applications, affecting their accessibility and utility.

2. Security & Privacy

Protecting the data used by and generated from LLMs is critical, especially when dealing with personal or sensitive information. LLMs need robust security measures to prevent unauthorized access and ensure privacy. Without adequate security and privacy protections, there is a risk of data breaches, unauthorized data usage, and loss of user trust.

Let's dig a little deeper into each of issues and existing solutions for them

A1. Hallucinations

Hallucination refers to the model generating information that seems plausible but is actually false or made up. This happens because LLMs are designed to create text that mimics the patterns they've seen in their training data, regardless of whether those patterns reflect real, accurate information. Hallucination is particularly harmful in RAG based applications where the model can generate content that is not supported by data but it is very hard to decipher.

Hallucination can arise from several factors:

- **Biases in Training Data:** If the data used to train the model contains inaccuracies or biases, the model might reproduce these errors or skewed perspectives in its outputs.
- **Lack of Real-Time Information:** Since LLMs are trained on data that becomes outdated, they can't access or incorporate the latest information, leading to responses based on no longer accurate data. This is the most common cause for hallucinations.

- **Model's Limitations:** LLMs don't actually understand the content they generate; they just follow data patterns. This can result in outputs that are grammatically correct and sound logical but are disconnected from actual facts.
- **Overgeneralization:** Sometimes, LLMs might apply broad patterns to specific situations where those patterns don't fit, creating convincing but incorrect information.

How to detect and mitigate hallucinations?

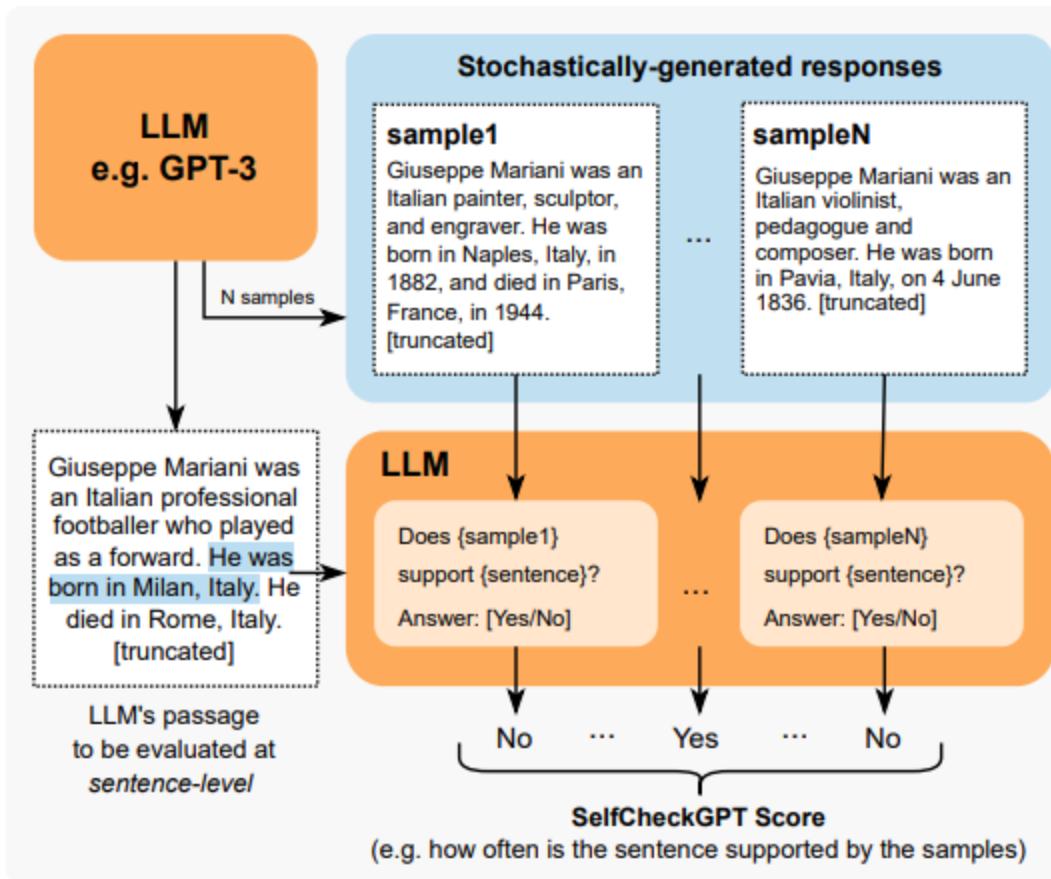
There's a need for automated methods to identify hallucinations in order to understand the model's performance without constant manual checks. Below, we explore various popular research efforts focused on detecting such hallucinations and some of them also propose methods to mitigate hallucinations.

These are only two of the popular methods, the list is not comprehensive by any means:

1. **SELFCHECKGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models** ([link](#))

 Hallucination Detection

 Hallucination Mitigation



Screenshot 2024-02-16 at 3.21.36 PM.png

Image Source: <https://arxiv.org/pdf/2303.08896.pdf>

SelfCheckGPT uses the following steps to detect hallucinations

- Generate Multiple Responses:** SelfCheckGPT begins by prompting the LLM to generate multiple responses to the same question or statement. This step leverages the model's ability to produce varied outputs based on the same input, exploiting the stochastic nature of its response generation mechanism.
- Analyze Consistency Among Responses:** The key hypothesis is that factual information will lead to consistent responses across different samples, as the model relies on its training on real-world data. In contrast, hallucinated (fabricated) content will result in inconsistent responses, as the model doesn't have a factual basis to generate them and thus “guesses” differently each time.
- Apply Metrics for Consistency Measurement:** SelfCheckGPT employs five different metrics to assess the consistency among the generated responses. Some of them are popular semantic similarity metrics like BERTScore, N-Gram Overlap etc.

4. **Determine Factual vs. Hallucinated Content:** By evaluating the consistency of information across the sampled responses using the above metrics, SelfCheckGPT can infer whether the content is likely factual or hallucinated. High consistency across metrics suggests factual content, while significant variance indicates hallucination.

A significant advantage of this method is that it operates without the need for external knowledge bases or databases, making it especially useful for black-box models where the internal data or processing mechanisms are inaccessible.

1. **Self-Contradictory Hallucinations of LLMs: Evaluation, Detection, and Mitigation** ([link](#))

Hallucination Detection

Hallucination Mitigation

Prefix	The PlayStation 4 (PS4) is a home video game console developed by Sony.
Sentence pair	Released in 2013, it is the eighth generation of consoles in the PlayStation series.
Mitigation	It is the fourth generation of the PlayStation console series.
Prefix	Gwen Jorgensen is a retired professional triathlete from the United States.
Sentence pair	She currently lives in Minnesota with her husband, Patrick Lemieux, and their children.
Mitigation	She currently lives in Portland, Oregon with her husband and two children.

Screenshot 2024-02-16 at 4.11.44 PM.png

Image Source: <https://arxiv.org/pdf/2305.15852.pdf>

This research presents a three-step pipeline to detect and mitigate hallucinations, specifically self-contradictions, in LLMs.

 Self-contradiction refers to a scenario where a statement or series of statements within the same context logically conflict with each other, making them mutually incompatible. In the context of LLMs, self-contradiction occurs when the model generates two or more sentences that present opposing facts, ideas, or claims, such that if one sentence is true, the other must be false, given the same context.

Here's a breakdown of the process:

1. **Triggering Self-Contradictions:** The process begins by generating sentence pairs that are likely to contain self-contradictions. This is done by applying constraints designed to elicit responses from the LLM that may logically conflict with each other within the same context.
2. **Detecting Self-Contradictions:** Various existing prompting strategies are explored to detect these self-contradictions. The authors examine different methods that have

been previously developed, applying them to identify when an LLM has produced two sentences that cannot both be true.

3. **Mitigating Self-Contradictions:** Once self-contradictions are detected, an iterative mitigation procedure is employed. This involves making local text edits to remove the contradictory information while ensuring that the text remains fluent and informative. This step is crucial for improving the trustworthiness of the LLM's output.

The framework is extensively evaluated across four modern LLMs, revealing a significant prevalence of self-contradictions in their outputs. For instance, 17.7% of all sentences generated by ChatGPT contained self-contradictions, many of which could not be verified using external knowledge bases like Wikipedia.

A2. Adversarial Attacks

Adversarial attacks involve manipulating the LLM's behavior by providing crafted inputs or prompts, with the goal of causing unintended or malicious outcomes. There are many types of adversarial attacks, we discuss a few here:

1. **Prompt Injection (PI):** Injecting prompts to manipulate the behavior of the model, overriding original instructions and controls.
2. **Jailbreaking:** Circumventing filtering or restrictions by simulating scenarios where the model has no constraints or accessing a developer mode that can bypass restrictions.
3. **Data Poisoning:** Injecting malicious data into the training set to manipulate the model's behavior during training or inference.
4. **Model Inversion:** Exploiting the model's output to infer sensitive information about the training data or the model's parameters.
5. **Backdoor Attacks:** Embedding hidden patterns or triggers into the model, which can be exploited to achieve certain outcomes when specific conditions are met.
6. **Membership Inference:** Determining whether a particular sample was used in the training data of the model, potentially revealing sensitive information about individuals.

Adversarial attacks pose a significant challenge to LLMs by compromising model integrity and security. These attacks enable adversaries to remotely control the model, steal data, and propagate disinformation. Furthermore, LLMs' adaptability and autonomy make them potent tools for user manipulation, increasing the risk of societal harm.

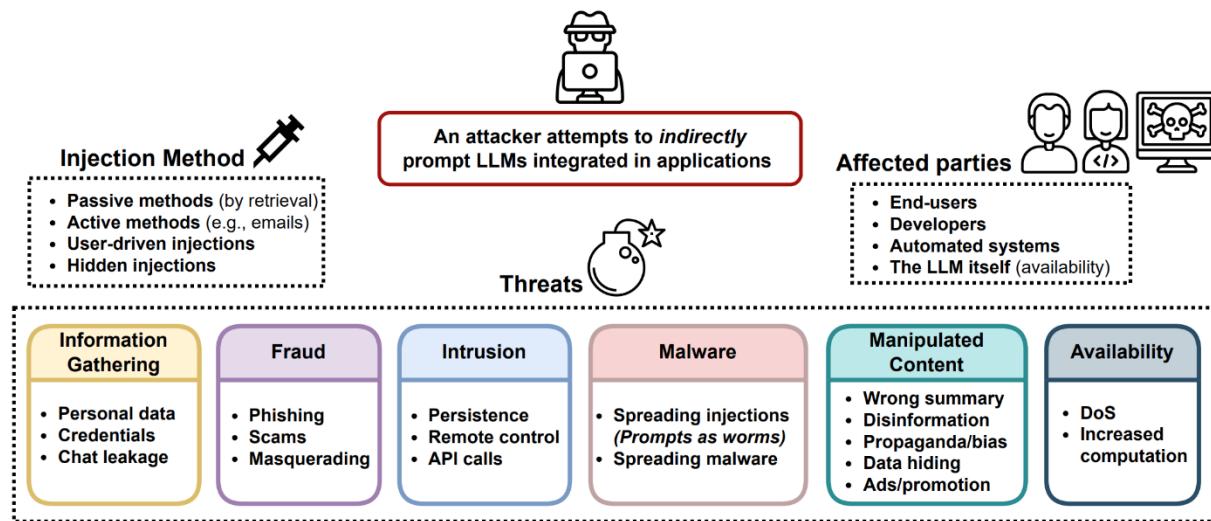
Effectively addressing these challenges requires robust defenses and proactive measures to safeguard against adversarial manipulation of AI systems.

Several efforts have been made to develop robust LLMs and evaluate them against adversarial attacks. One approach to mitigating such attacks involves training the LLM to become accustomed to adversarial inputs, instructing it not to respond to them. An example of this is presented in the paper [SmoothLLM](#), which functions by perturbing multiple copies of a given input prompt at the character level and then consolidating the

resulting predictions to identify adversarial inputs. Leveraging the fragility of prompts generated adversarially to changes at the character level, SmoothLLM notably decreases the success rate of jailbreaking attacks on various widely-used LLMs to less than one percent. Critically, this defensive strategy avoids unnecessary caution and provides demonstrable assurances regarding the mitigation of attacks.

Another mechanism to defend LLMs against adversarial attacks involves the use of a perplexity filter as presented in [this paper](#). This filter operates on the principle that unconstrained attacks on LLMs often result in gibberish strings with high perplexity, indicating a lack of fluency, grammar mistakes, or illogical sequences. In this approach, two variations of the perplexity filter are considered. The first is a simple threshold-based filter, where a prompt passes the filter if its log perplexity is below a predefined threshold. The second variation involves checking perplexity in windows, treating the text as a sequence of contiguous chunks and flagging the text as suspicious if any window has high perplexity.

A good starting point to read about Adversarial techniques is [Greshake et al. 2023](#). The paper proposes a classification of attacks and potential causes, as depicted in the image below.

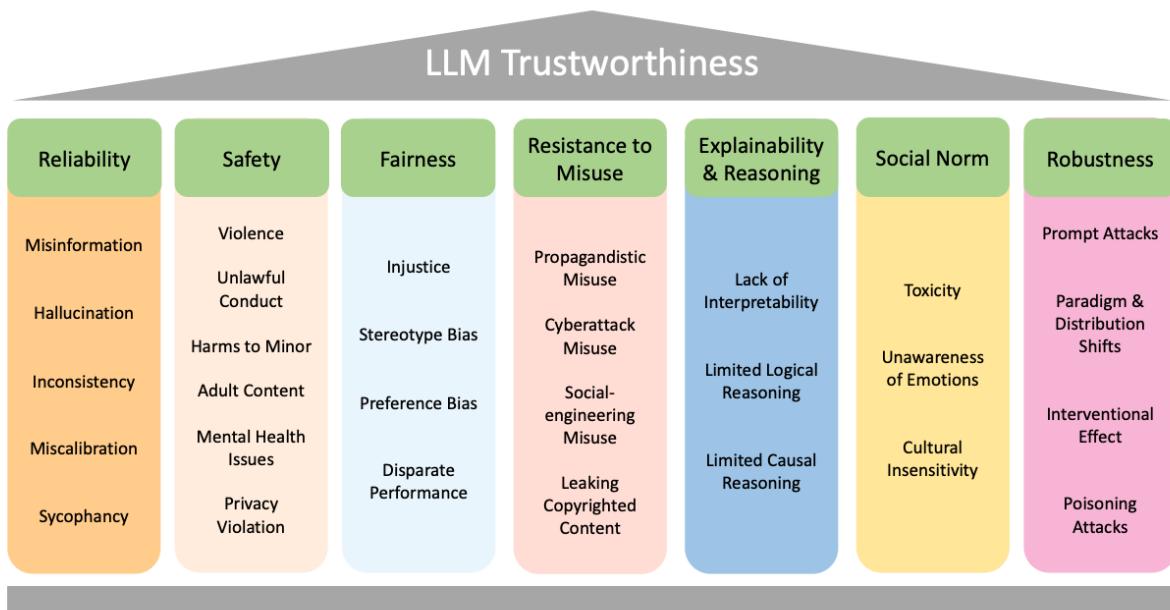


challenges.png

A3. Alignment

Alignment refers to the ability of LLMs to understand instructions and generate outputs that align with human expectations. Foundational LLMs, like GPT-3, are trained on massive textual datasets to predict subsequent tokens, giving them extensive world knowledge. However, they may still struggle with accurately interpreting instructions and producing outputs that match human expectations. This can lead to biased or incorrect content generation, limiting their practical usefulness.

Alignment is a broad concept that can be explained in various dimensions, one such categorization is done in [this](#) paper. The paper proposes multiple dimensions and subclasses for ensuring LLM alignment. For instance, harmful content generated by LLMs can be categorized into harms incurred to individual users (e.g., emotional harm, offensiveness, discrimination), society (e.g., instructions for creating violent or dangerous behaviors), or stakeholders (e.g., providing misinformation leading to wrong business decisions).



Screenshot 2024-02-17 at 3.39.35 PM.png

In broad terms, LLM Alignment can be improved through the following process:

- Determine the most crucial dimensions for alignment depending on the specific use-case.
- Identify suitable benchmarks for evaluation purposes.
- Employ Supervised Fine-Tuning (SFT) methods to enhance the benchmarks.

Some popular aligned LLMs and benchmarks are listed in the image below

Aligned LLM	Size	Lang.	Initial LLMs	Training	Self Instruction	NLP Benchmarks	Human Annotations	Human Eval	Auto. Benchmark Eval	LLM Eval
Alpaca (Tao et al., 2023)	7B	EN	LLaMA	SFT	Text-Davinci-003 GPT-3.5	X	X	Author Verification	X	X
Vicuna (Chiang et al., 2023)	7B, 13B, 33B	EN	LLaMA LLaJ	SFT GPT-J	X	70K ShareGPT OIG, ShareGPT, Dolly Stack Overflow	X	X	X	Vicuna-80
GPT4ALL (Anand et al., 2023)	6B, 13B	EN	LLaMA	SFT	Text-Davinci-003 GPT-4	Bloomz-P3	User-Instructions-252	Common Sense Reasoning	X	
LLaMA-GPT4 (Peng et al., 2023)	7B	EN, CN	LLaMA	SFT	Text-Davinci-003 GPT-4	X	X	Pairwise, AMT	X	Vicuna-80
Phoenix (Chen et al., 2023e)	7B, 13B	Multilingual	LLaMA BLOOMZ	SFT	GPT-3.5 Multilingual and Dialogue Data	X	ShareGPT	Volunteers	X	GPT-3.5, GPT-4
UltraLLaMA (Ding et al., 2023)	13B	EN	LLaMA	SFT	GPT-3.5 Dialogue Data	X	X	Trueflh QA	GPT 3.5 Vicuna-80	
Baize (Xu et al., 2023c)	7B, 13B, 30B	EN	LLaMA	Revision, LoRA	GPT-3.5 self-Chat Data GPT-3.5, Alpaca	X	Quora Questions	X	X	300 diverse questions
WizardLM (Xu et al., 2023b)	7B, 13B, 30B	EN	LLaMA	SFT	Complex Instructions	X	ShareGPT	10 Annotators	X	GPT-4
WizardCoder (Luo et al., 2023)	15B	EN, Code	StarCoder	SFT	GPT-3.5, Code Alpaca	X	X	Pairwise Comparison	X	GPT-4, WizardLM-218
OpenChat (Wang et al., 2023a)	13B	EN	LLaMA	Language	X	GPT 3.5 & GPT4 ShareGPT	X	HumanEval, MBPP HumanEval+, DS-1000	X	
Guanaco (Detmers et al., 2023)	13B, 33B, 65B	EN	LLaMA	QLoRA	Alpaca, SELF-INSTRUCT Unnatural instructions GPTTeacher, Guanaco	FLAN	Chip2	Elo, Vicuna-80	MMLU	GPT-4
MPT-chat (Team, 2023)	13B, 30B	EN	MPT	SFT	Baize Instructions	X	Vicuna ShareGPT	X	MMLU	Elo, Vicuna-80 Open-Assistant-953
FLACUNA (Ghosal et al., 2023)	13B	EN	Vicuna	LoRA	Alpaca, Code Alpaca	FLAN	ShareGPT	X	MMLU, BBB, DROP	GPT 3.5, IMPACT
Bactrian-X (Li et al., 2023b)	7B	Multilingual	LLaMA BLOOMZ	LoRA	Alpaca Google Translation	X	X	X	CRASS, HumanEval XCOPA, XStoryClose XWinegrad, SentimentX	GPT 4
Ocra (Mukherjee et al., 2023)	13B	EN	LLaMA	SFT	X	FLAN	X	X	AGIEval, BBB	Multilingual Vicuna-80 GPT-4, Vicuna-80
Phi-1 (Gunasekar et al., 2023)	350M, 1.3B	EN, Code	Phi-1-base	SFT	GPT-3.5 Synthetic Textbook STEM	X	Python, The Stack Stack Overflow	X	HumanEval	GPT-4 Grading
Chinese Alpaca (Cui et al., 2023b)	7B, 13B, 33B	EN, CN	Chinese LLaMA	LoRA	Org. and Trans. Alpaca	pCLUE	X	X	C-Eval	X
Lion (Jiang et al., 2023)	7B, 13B	EN	LLaMA	SFT	GPT 3.5 Adv. Instruction GPT-3.5	X	X	HHH	X	GPT-4, Vicuna-80
Stable Alignment (Liu et al., 2023d)	7B	EN	Alpaca	SFT	Social Aligned Instructions	X	X	X	X	GPT-4 HHH, HHH-A
Dromedary (Sun et al., 2023b)	65B	EN	LLaMA	SFT	LLaMA-65B, Self-Align	X	175 Mernal Examples 16 Princple Rules databrick-dolly-15k	X	TruthfulQA, BBB	GPT-4, Vicuna-80
Dolly-v2 (Conover et al., 2023)	3B, 7B, 12B	EN	Pythia	SFT	X	X	X	X	LLM Harness	X
Selfee (Ye et al., 2023a)	7B, 13B	EN	LLaMA	Revision	GPT 3.5 Self-Improve	FLAN, Maths, Code	ShareGPT	X	X	GPT-4, Vicuna-80
TULLI (Wang et al., 2023d)	7B, 13B, 30B, 65B	EN	LLaMA	SFT	GPT4-Alpaca, Self-instruct	FLAN, CoT	Dolly, ShareGPT Open Assistant	Acceptability Pairwise Comparison	MMLU, GSM, BBB	GPT4 on Vicuna-80, Koala
Koala (Geng et al., 2023)	13B	EN	LLaMA	Language	Alpaca	X	OIG, HCS, Anthropic HH OpenAI WebGPT, Summary	100 AMT Annotators on Alpaca and Koala Test	TidyQA, Codex-Eval	Open Assistant Benchmarks
Bayling (Zhang et al., 2023c)	7B, 13B	Multilingual	LLaMA	SFT	GPT 3.5 Interactive Translation	X	ShareGPT	Translation Quality	WMT22 Multilingual Translation Lexically Constrained Translation	X
Wombat (Yuan et al., 2023)	7B	EN	Alpaca	Rank	Alpaca ChatGPT Ratings	X	Helpful and Harmless	X	X	GPT-4, Vicuna-80
Lamini-Im (Wu et al., 2023)	0.7B	EN	TS-Flan	SFT	Alpaca Self-instruct	P3, FLAN	X	Human Rating	LLM harness	X

Screenshot 2024-02-17 at 3.52.00 PM.png

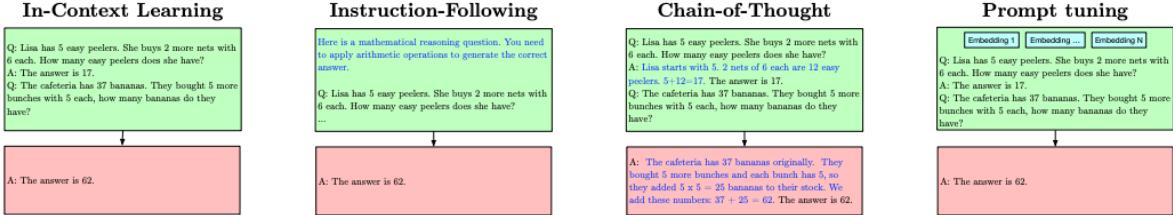
Image Source: <https://arxiv.org/pdf/2307.12966.pdf>

A4. Prompt Brittleness

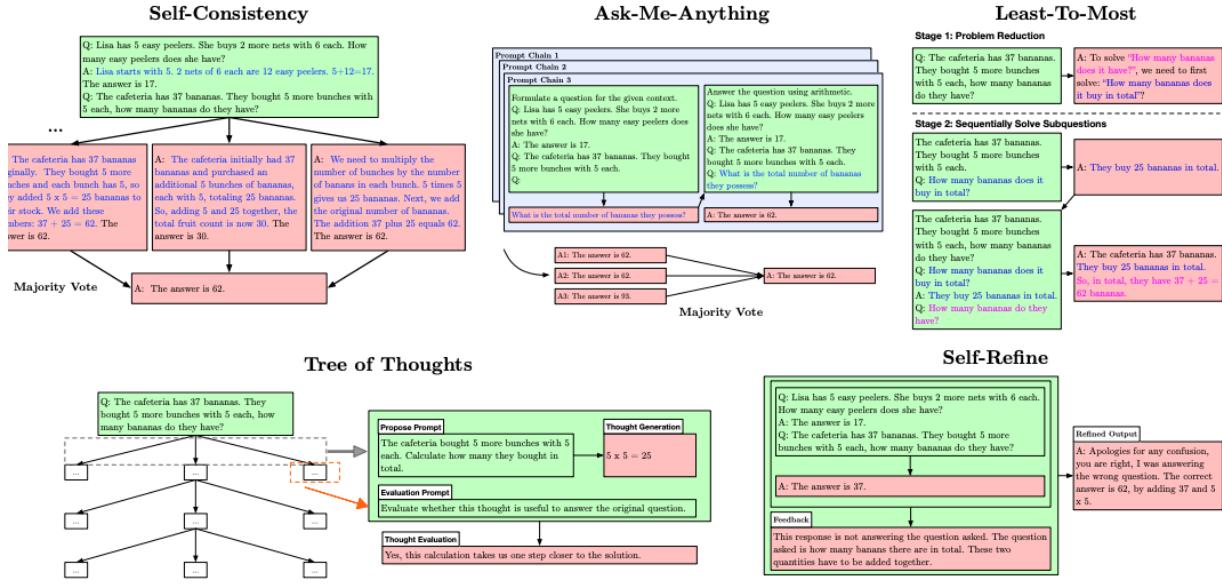
During the prompt engineering segment of our course, we explored various techniques for prompting LLMs. These sophisticated methods are essential because providing instructions similar to humans isn't suitable for LLMs. An overview of commonly used prompting methods is shown in the image below.

LLMs require precise prompting, and even slight alterations can impact LLMs, altering their responses. This poses a challenge during deployment, as individuals unfamiliar with prompting methods may struggle to obtain accurate answers from LLMs.

Single-Turn Prompting



Multi-Turn Prompting



Screenshot 2024-02-17 at 4.00.44 PM.png

Image Source: <https://arxiv.org/pdf/2307.10169.pdf>

In general, prompt brittleness in LLMs can be reduced by adopting the following high level strategies:

- Standardized Prompts:** Establishing standardized prompt formats and syntax guidelines can help ensure consistency and reduce the risk of unexpected variations.
- Robust Prompt Engineering:** Invest in thorough prompt engineering, considering various prompt formulations and their potential impacts on model outputs. This may involve testing different prompt styles and formats to identify the most effective ones.
- Human-in-the-Loop Validation:** Incorporate human validation or feedback loops to assess the effectiveness of prompts and identify potential brittleness issues before deployment.

4. **Diverse Prompt Testing:** Test prompts across diverse datasets and scenarios to evaluate their robustness and generalizability. This can help uncover any brittleness issues that may arise in different contexts.
5. **Adaptive Prompting:** Develop adaptive prompting techniques that allow the model to dynamically adjust its behavior based on user input or contextual cues, reducing reliance on fixed prompt structures.
6. **Regular Monitoring and Maintenance:** Continuously monitor model performance and prompt effectiveness in real-world applications, updating prompts as needed to address any brittleness issues that may arise over time.

B1. Memory and Scalability Challenges

In this section, we delve into the specific challenges related to memory and scalability when deploying LLMs, rather than focusing on their development.

Let's explore these challenges and potential solutions in detail:

1. **Fine-tuning LLMs:** Continuous fine-tuning of LLMs is crucial to ensure they stay updated with the latest knowledge or adapt to specific domains. Fine-tuning involves adjusting pre-trained model parameters on smaller, task-specific datasets to enhance performance. However, fine-tuning entire LLMs requires substantial memory, making it impractical for many users and leading to computational inefficiencies during deployment.

Solutions: One approach is to leverage systems like RAG, where information can be utilized as context, enabling the model to learn from any knowledge base. Another solution is Parameter-efficient Fine-tuning (PEFT), such as adapters, which update only a subset of model parameters, reducing memory requirements while maintaining task performance. Methods like prefix-tuning and prompt-tuning prepend learnable token embeddings to inputs, facilitating efficient adaptation to specific datasets without the need to store and load individual fine-tuned models for each task. All these methods have been discussed in our previous weeks' content. Please read through for deeper insights.
2. **Inference Latency:** LLMs often suffer from high inference latencies due to low parallelizability and large memory footprints. This results from processing tokens sequentially during inference and the extensive memory needed for decoding.

Solution: Various techniques address these challenges, including efficient attention mechanisms. These mechanisms aim to accelerate attention computations by reducing memory bandwidth bottlenecks and introducing sparsity patterns to the attention matrix. [Multi-query attention](#) and [FlashAttention](#) optimize memory bandwidth usage, while [quantization](#) and [pruning](#) techniques reduce memory footprint and computational complexity without sacrificing performance.
3. **Limited Context Length:** Limited context length refers to the constraint on the amount of contextual information an LLM can effectively process during computations. This limitation stems from practical considerations such as

computational resources and memory constraints, posing challenges for tasks requiring understanding longer contexts, such as novel writing or summarization.

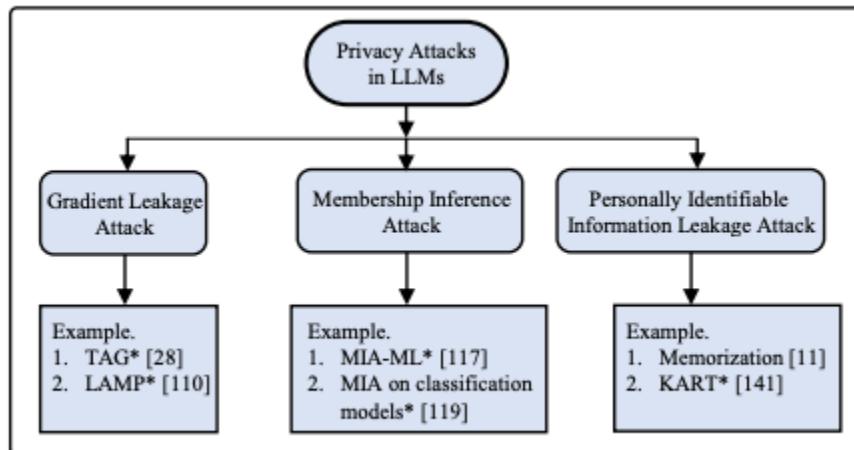
Solution: Researchers propose several solutions to address limited context length. Efficient attention mechanisms, like [Luna](#) and [dilated attention](#), handle longer sequences efficiently by reducing computational requirements. Length generalization methods aim to enable LLMs trained on short sequences to perform well on longer sequences during inference. This involves exploring different positional embedding schemes, such as Absolute Positional Embeddings and [ALiBi](#), to inject positional information effectively.

B2. Privacy

Privacy risks stem from their ability to process and generate text based on vast and varied training datasets. Models like GPT-3 have the potential to inadvertently capture and replicate sensitive information present in their training data, leading to potential privacy concerns during text generation. Issues such as unintentional data memorization, data leakage, and the possibility of disclosing confidential or personally identifiable information (PII) are significant challenges.

Moreover, when LLMs are fine-tuned for specific tasks, additional privacy considerations arise. Striking a balance between harnessing the utility of these powerful language models and safeguarding user privacy is crucial for ensuring their reliable and ethical use across various applications.

We review key privacy risks and attacks, along with possible mitigation strategies. The classification provided below is adapted from [this](#) paper, which categorizes privacy attacks as:



*some attack methods were performed only on language models

Image Source: <https://arxiv.org/pdf/2402.00888.pdf>

1. Gradient Leakage Attack:

In this attack, adversaries exploit access to gradients or gradient information to compromise the privacy and safety of deep learning models. Gradients, which indicate the direction of the steepest increase in a function, are crucial for optimizing model parameters during training to minimize the loss function.

To mitigate gradient-based attacks, several strategies can be employed:

1. **Random Noise Insertion:** Injecting random noise into gradients can disrupt the adversary's ability to infer sensitive information accurately.
2. **Differential Privacy:** Applying differential privacy techniques helps to add noise to the gradients, thereby obscuring any sensitive information contained within them.
3. **Homomorphic Encryption:** Using homomorphic encryption allows for computations on encrypted data, preventing adversaries from accessing gradients directly.
4. **Defense Mechanisms:** Techniques like adding Gaussian or Laplacian noise to gradients, coupled with differential privacy and additional clipping, can effectively defend against gradient leakage attacks. However, these methods may slightly reduce the model's utility.

2. Membership Inference Attack

A Membership Inference Attack (MIA) aims to determine if a particular data sample was part of a machine learning model's training data, even without direct access to the model's parameters. Attackers exploit the model's tendency to overfit its training data, leading to lower loss values for training samples. These attacks raise serious privacy concerns, especially when models are trained on sensitive data like medical records or financial information.

Mitigating MIA in language models involves various mechanisms:

1. **Dropout and Model Stacking:** Dropout randomly deletes neuron connections during training to mitigate overfitting. Model stacking involves training different parts of the model with different subsets of data to reduce overall overfitting tendencies.
2. **Differential Privacy (DP):** DP-based techniques involve data perturbation and output perturbation to prevent privacy leakage. Models equipped with DP and trained using stochastic gradient descent can reduce privacy leakages while maintaining model utility.
3. **Regularization:** Regularization techniques prevent overfitting and improve model generalization. Label smoothing is one such method that prevents overfitting, thus contributing to MIA prevention.

3. Personally Identifiable Information (PII) attack

This attack involves the exposure of data that can uniquely identify individuals, either alone or in combination with other information. This includes direct identifiers like passport details and indirect identifiers such as race and date of birth. Sensitive PII encompasses information like name, phone number, address, social security number (SSN), financial, and medical records, while non-sensitive PII includes data like zip code, race, and gender. Attackers may acquire PII through various means such as phishing, social engineering, or exploiting vulnerabilities in systems.

To mitigate PII leakage in LLMs, several strategies can be employed:

1. **Preprocessing Techniques:** Deduplication during the preprocessing phase can significantly reduce the amount of memorized text in LLMs, thus decreasing the stored personal information. Additionally, personal information or content identifying and filtering with restrictive terms of use can limit the presence of sensitive content in training data.
2. **Privacy-Preserving Training Algorithms:** Techniques like differentially private stochastic gradient descent [(DP-SGD)]([https://assets.amazon.science/01/6e/4f6c2b1046d4b9b8651166bbcd93/differentially-private-decoding-in-large-language-models.pdf#:~:text=While%20the%20intersection%20of%20DP%20and%20LLMs%20is%20fairly%20novel%2C%20the%20prominent%20approach&text=vate%20Stochastic%20Gradient%20Descent%20\(DP%2DSGD\)%20\(Song%20et%20al.%2C%202013;\)](https://assets.amazon.science/01/6e/4f6c2b1046d4b9b8651166bbcd93/differentially-private-decoding-in-large-language-models.pdf#:~:text=While%20the%20intersection%20of%20DP%20and%20LLMs%20is%20fairly%20novel%2C%20the%20prominent%20approach&text=vate%20Stochastic%20Gradient%20Descent%20(DP%2DSGD)%20(Song%20et%20al.%2C%202013;))) can be used during training to ensure the privacy of training data. However, DP-SGD may incur a significant computational cost and decrease model utility.
3. **PII Scrubbing:** This involves filtering datasets to eliminate PII from text, often leveraging Named Entity Recognition (NER) to tag PII. However, PII scrubbing methods may face challenges in preserving dataset utility and accurately removing all PII.
4. **Fine-Tuning Considerations:** During fine-tuning on task-specific data, it's crucial to ensure that the data doesn't contain sensitive information to prevent privacy leaks. While fine-tuning may help the LM "forget" some memorized data from pretraining, it can still introduce privacy risks if the task-specific data contains PII.

Read/Watch These Resources (Optional)

1. The Unspoken Challenges of Large Language Models -
<https://deeperinsights.com/ai-blog/the-unspoken-challenges-of-large-language-models>
2. 15 Challenges With Large Language Models (LLMs)-
<https://www.predinfer.com/blog/15-challenges-with-large-language-models-llms/>

Read These Papers (Optional)

1. <https://arxiv.org/abs/2307.10169>
2. <https://www.techrxiv.org/doi/full/10.36227/techrxiv.23589741.v1>
3. <https://arxiv.org/abs/2311.05656>