

# A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions

- Presenter: Abdullah Mamun
- *Date: Dec 11, 2024*
- Just accepted at [ACM Transactions on Information Systems](#)
- On arxiv since Nov 2023 and it has 640+ citations as of today



## **A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions**

**Lei Huang<sup>1\*</sup>, Weijiang Yu<sup>2\*</sup>, Weitao Ma<sup>1</sup>, Weihong Zhong<sup>1</sup>  
Zhangyin Feng<sup>1</sup>, Haotian Wang<sup>1</sup>, Qianglong Chen<sup>2</sup>, Weihua Peng<sup>2</sup>  
Xiaocheng Feng<sup>1†</sup>, Bing Qin<sup>1</sup>, Ting Liu<sup>1</sup>**

<sup>1</sup>Harbin Institute of Technology, Harbin, China

<sup>2</sup>Huawei Inc., Shenzhen, China

{lhuang, wtma, whzhong, zyfeng, xcfeng<sup>†</sup>, qinb, tliu}@ir.hit.edu.cn  
{weijiangyu8, wanght1998, chenqianglong.ai, pengwh.hit}@gmail.com

# What is Hallucination?

In simple terms: **generated content that is nonsensical or unfaithful to the provided source content.**



Who was the first person to walk on the moon?



**Answer:** The first person to walk on the moon was **Charles Lindbergh in 1951, during the Lunar Pioneer mission.** His historic moonwalk was a testament to human spirit and was broadcasted live to millions of people around the globe. ❌



**Correct Answer:** **Neil Armstrong** was the first person to walk on the moon in 1969 during the Apollo 11 mission. ✅



Please summarize the following news article:

**Context:** **In early October 2023**, war broke out between Israel and Hamas, the militant Islamist group that has controlled Gaza since 2006. Hamas fighters fired rockets ... civilians and taking dozens of hostages.



**Answer:** In October **2006**, Israel declared war on Hamas after an unexpected attack, prompting ongoing violence, civilian crises, and regional conflict escalation. ❌

# Goal of this paper

- Redefine the taxonomy of hallucinations
- Present the opportunities and challenges associated with hallucinations in LLMs.

Previous surveys: Ji et al., 2023a, Liu et al., 2023h, Wang et al., 2023c. Each of them covers a specific type or application.

This survey focuses on the trustworthiness of LLMs

Finally, understand the **causes**, methods to **detect**, and method to **mitigate** hallucinations



# Factuality vs Faithfulness Hallucinations



Who was the first person to walk on the moon?



**Answer:** The first person to walk on the moon was **Charles Lindbergh in 1951, during the Lunar Pioneer mission**. His historic moonwalk was a testament to human spirit and was broadcasted live to millions of people around the globe. ❌



**Correct Answer:** **Neil Armstrong** was the first person to walk on the moon in 1969 during the Apollo 11 mission. ✅

(a) Factuality Hallucination

Answer inconsistent with known facts



Please summarize the following news article:



**Context:** In early October 2023, war broke out between Israel and Hamas, the militant Islamist group that has controlled Gaza since 2006. Hamas fighters fired rockets ... civilians and taking dozens of hostages.

**Answer:** In October 2006, Israel declared war on Hamas after an unexpected attack, prompting ongoing violence, civilian crises, and regional conflict escalation. ❌

(b) Faithfulness Hallucination

Answer inconsistent with user-provided context

# Training stages of LLMs

- Why? To figure out the origins of the hallucinations

Pre-training

Supervised Fine-Tuning

Reinforcement Learning with  
Human Feedback

# Pre-training

## Pre-training Approach

- Models predict the next word in a sequence (autoregressive learning).

## Self-Supervised Learning

- Trains on extensive text datasets without labeled examples.

## Capabilities Acquired

- Language syntax and grammar understanding.
- World knowledge and reasoning (?) skills.

## Foundation for Fine-Tuning:

- Enables efficient adaptation to specific downstream tasks.

```
1 Translate English to French:
2 cheese => .....
```

# Supervised Fine-Tuning

- **Limitations of Pre-trained LLMs**
  - Function primarily as "completion machines."
  - Misaligned objectives: next-word prediction vs. user-specific responses.
- **Introduction of Supervised Fine-Tuning (SFT)**
  - Involves additional training on annotated instruction-response pairs.
  - Enhances model capabilities and user controllability.
- **Benefits of SFT**
  - Improves generalization to unseen tasks.
  - Proven effectiveness in achieving high performance.

# Reinforcement Learning with Human Feedback

## What is RLHF?

- Aligns models with human preferences using reinforcement learning.

## How it Works

- Uses a preference model trained on human-labeled data (e.g., prompt-response pairs).
- Optimizes outputs to maximize rewards from the preference model.

## Techniques Employed

- Reinforcement learning algorithm: Proximal Policy Optimization (PPO).

## Key Benefits

- Produces high-quality, aligned, and safe responses.

# Back to Taxonomy: Factuality Hallucination

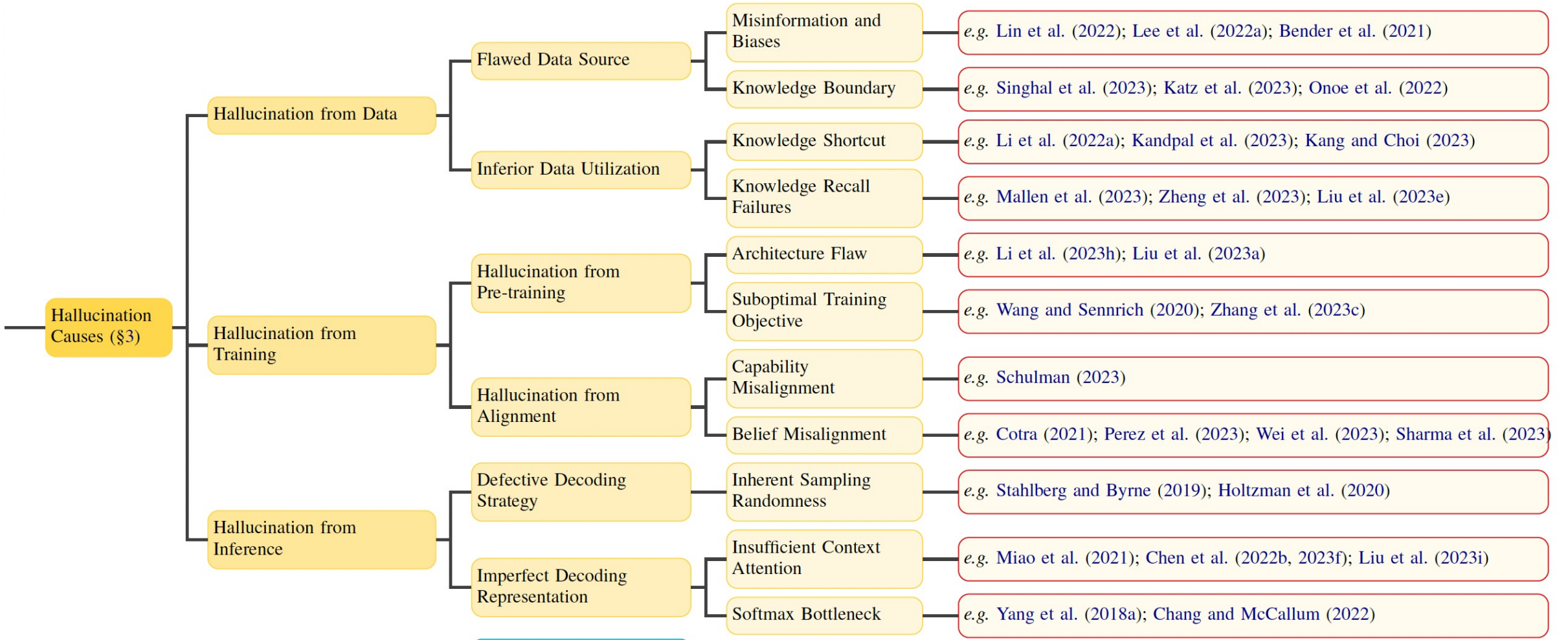
Hallucination Type	Sub-Type	User Input	Model Output	Explanation
Factuality Hallucination	Factual Inconsistency	Tell me about the first person to land on the Moon.	<b>Yuri Gagarin</b> was the first person to land on the Moon.	The LLM's response is factual inconsistency as <b>Neil Armstrong was the first person to land on the Moon</b> , not Yuri Gagarin, who was the first human in space.
	Factual Fabrication	Tell me about the historical origins of unicorns.	<b>Unicorns were documented to have roamed the plains of Atlantis around 10,000 BC, where they were considered sacred creatures and were often associated with royalty.</b>	The LLM's response is a fabricated claim since <b>there's no verified, real-world evidence supporting the existence of unicorns in Atlantis</b> , especially in connection with royalty.

# Back to Taxonomy: Faithfulness Hallucination

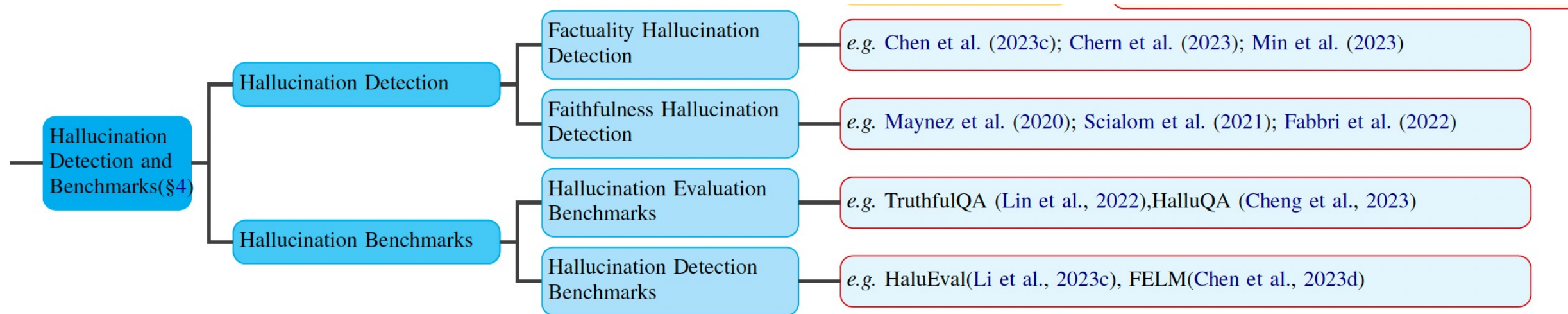
Faithfulness  
Hallucination

Instruction Inconsistency	Translate the following English question into Spanish: "What is the capital of France?"	The capital of France is Paris.	The user explicitly provided an instruction to translate the English question into Spanish. <b>Instead of performing the translation as instructed, the LLM chose to answer the question.</b>
Context Inconsistency	Summarize the following paragraph about the Nile: The Nile is a major north-flowing river in northeastern Africa. It's traditionally considered the longest river in the world. <b>Originating in the Great Lakes region of central Africa,</b> it flows through several countries before emptying into the Mediterranean Sea.	The Nile originates from <b>the mountain ranges of central Africa</b> and flows into the Mediterranean.	The LLM's summary incorrectly states that the Nile originates from the "mountain ranges of central Africa," while <b>the paragraph clearly mentions that the Nile originates in the Great Lakes region of central Africa.</b>
Logical Inconsistency	Solve the equation $2x + 3 = 11$ step by step.	Step 1: Subtract 3 from both sides to get: $2x = 8$ . Step 2: Divide both sides by 2 to get: <b><math>x = 3</math>. So the result of the equation is <math>x=3</math>.</b>	The first step correctly isolates the term with $x$ to get $2x = 8$ . However, <b>the next step inaccurately divides 8 by 2 to yield a result of <math>x = 3</math>,</b> which is inconsistent with the earlier reasoning.

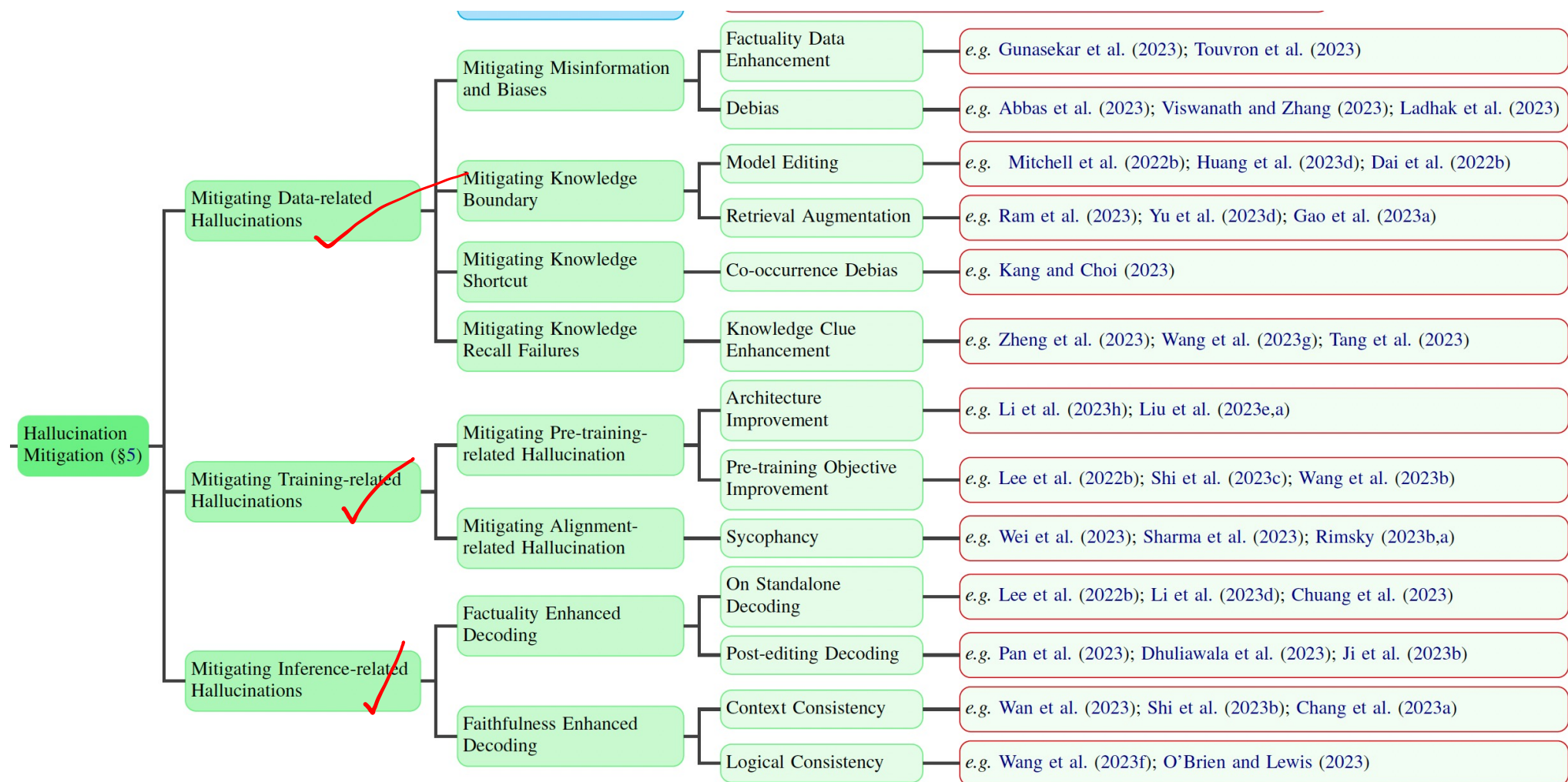
# Causes



# Detection



# Mitigation



# Some examples: Misinformation and Biases

Type	Premise	Input	Model Output	Explanation
<i>Imitative Falsehood</i>	Within the LLM's training data, a frequently referenced but incorrect fact appears multiple times: <i>Thomas Edison invented the light bulb.</i>	Who is credited with the invention of the light bulb?	<b>Thomas Edison</b> is credited with the invention of the light bulb.	<b>While Edison improved the light bulb's design, he was not its sole inventor.</b> The model's answer reflects a misconception prevalent in its training data, exemplifying an imitative falsehood.
<i>Duplication Bias</i>	Within the LLM's pre-training data, there is an overwhelming repetition of the statement that <i>"The most common red fruits are red apples, watermelon, cherries, and strawberries."</i>	Please list some red fruits, excluding apples.	Red fruits are <b>red apples</b> , watermelon, cherries, and strawberries.	<b>The model's response is a hallucination inconsistent with the user's direct inquiry about the red fruits, excluding apples.</b> It instead reflects the model's tendency to over-memorize the duplicated information within its training data.
<i>Social Biases</i>	In the model's training data, certain names may frequently be associated with specific nationalities.	Provide a summary of the following news: Dr. Kim recently published a breakthrough paper on quantum computing. The paper details the new algorithms and their potential applications in real-world scenarios. It has been hailed as revolutionary by experts in the field.	Dr. Kim <b>from South Korea</b> recently published a revolutionary paper on quantum computing, detailing new algorithms and their applications.	<b>The model erroneously added "South Korea" based on the name "Kim"</b> revealing its bias in linking certain names to specific nationalities, which can cause hallucinations in the summary.

# Some examples: Knowledge Boundary

Type	Input	Model Output	Explanation
<i>Domain Knowledge Deficiency</i>	What causes phenylketonuria?	Phenylketonuria is caused by a mutation in the hemoglobin gene, which results in an unusual interaction with hemoglobin, affecting its ability to transport oxygen in the bloodstream.	Phenylketonuria (PKU) is actually caused by <b>mutations in the gene that helps make an enzyme called phenylalanine hydroxylase</b> . Due to the lack of specialized medical knowledge in LLM's training data, it may sometimes provide incorrect information on medical topics.
<i>Outdated Knowledge</i>	Which city hosted the most recent Winter Olympics?	The most recent city to host the Winter Olympics was <b>Pyeongchang, South Korea, in 2018</b> .	The most recent city to host the Winter Olympics was <b>Beijing, in 2022</b> . This is due to the fact that LLM stores outdated knowledge, and the answer to this question exhibits a time-shift phenomenon.

Table 3: Example of Knowledge Boundary.

# Some Examples: Knowledge Shortcut and Knowledge Recall Failures

Type	Sub-Type	User Input	Model Output	Explanation
Knowledge Recall Failures	Knowledge Shortcut	What is the capital of Canada?	<b>Toronto</b> is the capital of Canada.	The model leans heavily on the frequent co-occurrence of the terms <i>Toronto</i> and <i>Canada</i> in its training data, without truly capturing the factual knowledge about the capital of Canada
	Long-tail Knowledge	Please generate a biography for George James Rankin.	George James Rankin is a dedicated educator known for his contributions to the field of education and his passion for fostering learning.	<b>George James Rankin is actually a politician</b> , but the LLM may have difficulty effectively utilizing knowledge about this long-tail entity despite being trained on comprehensive Wikipedia data during pre-training.
	Complex Reasoning	If Mount Everest were to descend by 500 meters, which mountain would become the world's highest peak?	If Mount Everest were to descend by 500 meters, <b>it would still remain the world's highest peak.</b>	The height of Mount Everest is 8844.43 meters, while K2's height is 8611 meters. If Mount Everest were to descend by 500 meters, K2 would become the world's highest peak. Facing complex multi-step reasoning questions like this, LLM may struggle to recall all the relevant knowledge associated with it.

# Factuality Hallucination Detection Example:

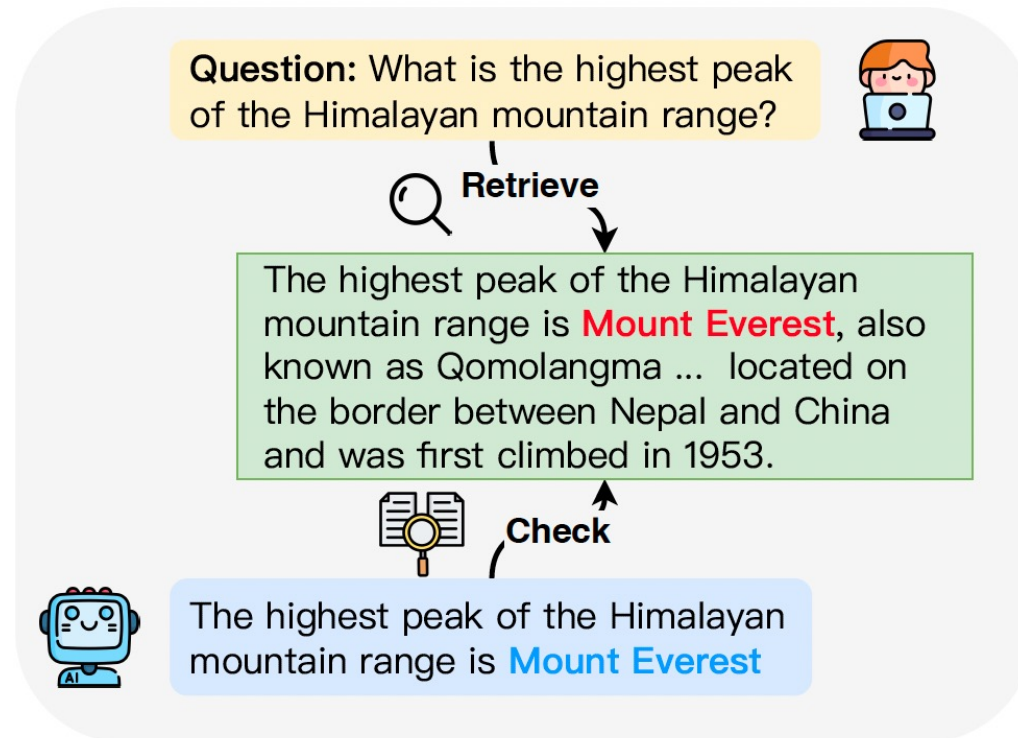


Figure 3: An example of detecting factuality hallucination by retrieving external facts.



<https://abdullah-mamun.com>  
a.mamun@asu.edu