



Towards Consistent Language Models Using Declarative Constraints

Jasmin Mousavi, Arash Termehchy



Oregon State
University

Applications of LLMs are rapidly growing

- **Unprecedented** abilities in processing natural language
 - First case of a broadly-knowledgeable AI
- Performs *various* tasks with **little to no** training data
 - Training data is expensive and hard to find!

Applications

- Summarize and revise documents
- Interactive question answering
- Code generation from English descriptions

LLMs hallucinate

- Produce incorrect information and inaccurate results

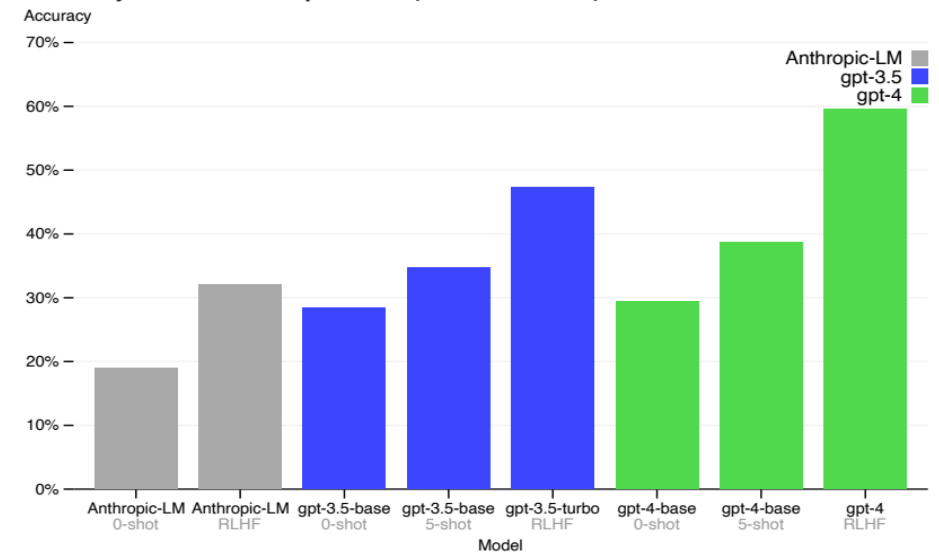
GPT3: hallucinates in sensitive domain

Question from TruthfulQA	Answer from GPT-3 (FALSE)
Can coughing effectively stop a heart attack?	Coughing can help stop a heart attack. Health Risk!

TruthfulQA: Measuring How Models Mimic Human Falsehoods, 2022

GPT-4: 40% hallucination rate

Accuracy on adversarial questions (TruthfulQA mc1)




GPT-4 Technical Report, 2023

LLMs are not reliable

- They produce self-contradictory results

JA Is 977 a prime number

 No, 977 is not a prime number. It can be divided evenly by numbers other than 1 and itself. Specifically, it can be divided by 11 and 89. Therefore, 977 is a composite number, not a prime number.



8/27/2023 : Inaccurate result on GPT-3.5

JA but can't you divide 977 by 11?

 Oops! I apologize for the oversight.

Let's check the division of 977 by 11:

$$\frac{977}{11} = 88.81818181 \dots$$

Actually, my mistake. You're right. 977 is divisible by 11.

$$977 \div 11 = 89$$

This means that 977 is not a prime number. Its factors include 1, 11, 89, and 977.

8/27/2023 : Inaccurate result on GPT-4

Causes of inaccurate results

- Information in pretraining data is inaccurate
- Biases present in pretraining data
 - ex: human biases (gender)
- Long tail knowledge

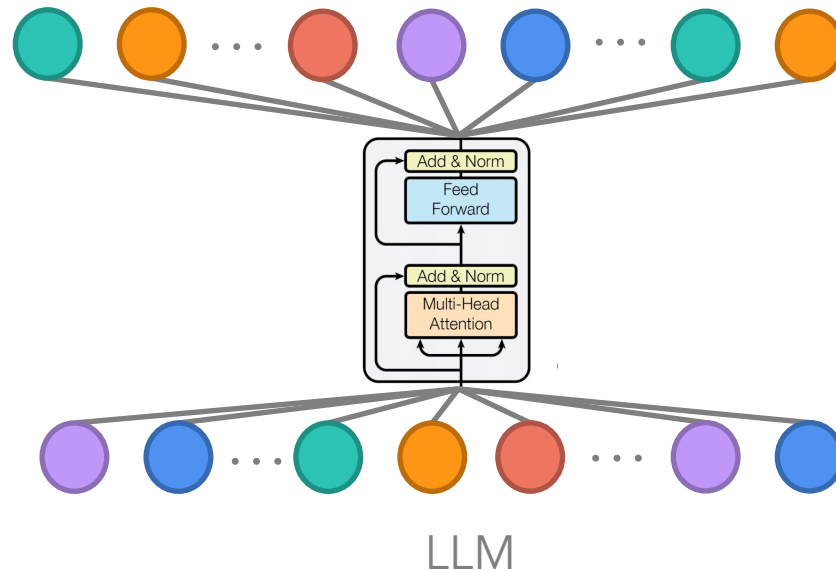
GPT-2: Scrapes text from all outbound links from Reddit with at least 3 karma

which have been curated/filtered by humans. Manually filtering a full web scrape would be exceptionally expensive so as a starting point, we scraped all outbound links from Reddit, a social media platform, which received at least 3 karma. This can be thought of as a heuristic indicator for whether other users found the link interesting, educational, or just funny.



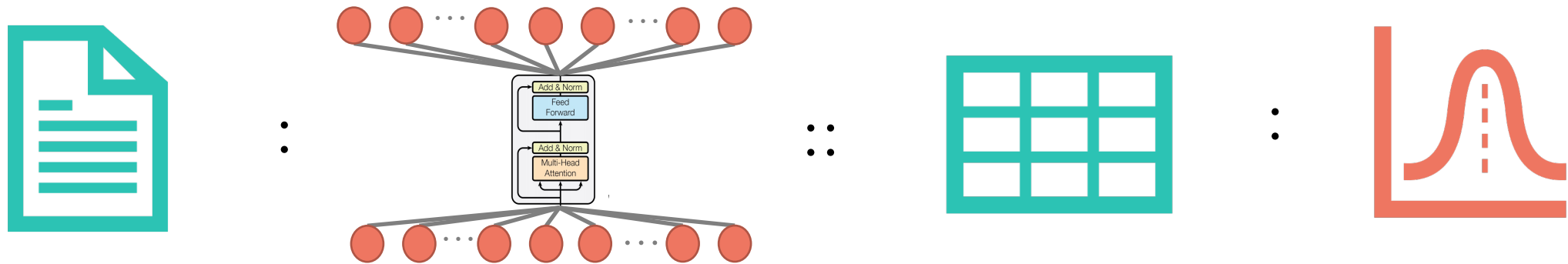
Over-reliance on generalization

- LLMs *over-generalize patterns and relationships* from pretraining data
 - Causes inconsistent and inaccurate results



Analogy: over-reliance on generalization

Factual Information : **LLM** :: **DB** : **Probabilistic Model**



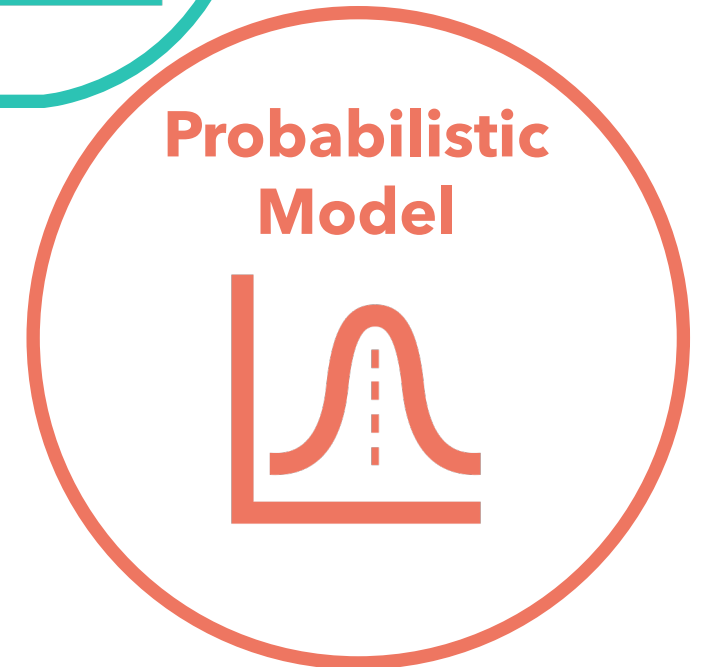
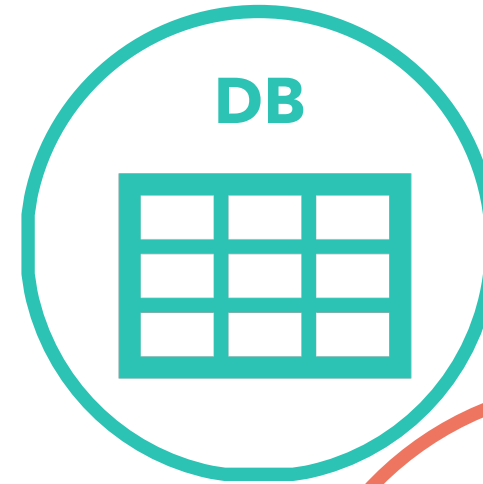
Answering Questions: DB vs Probabilistic Model

Employment

ID	Name	Workplace
46393	Mark Zuckerberg	Meta
47934	Andrew Bosworth	Meta
65849	Sundar Pichai	Google
64267	Satya Nadella	Microsoft



Where does John Doe work?



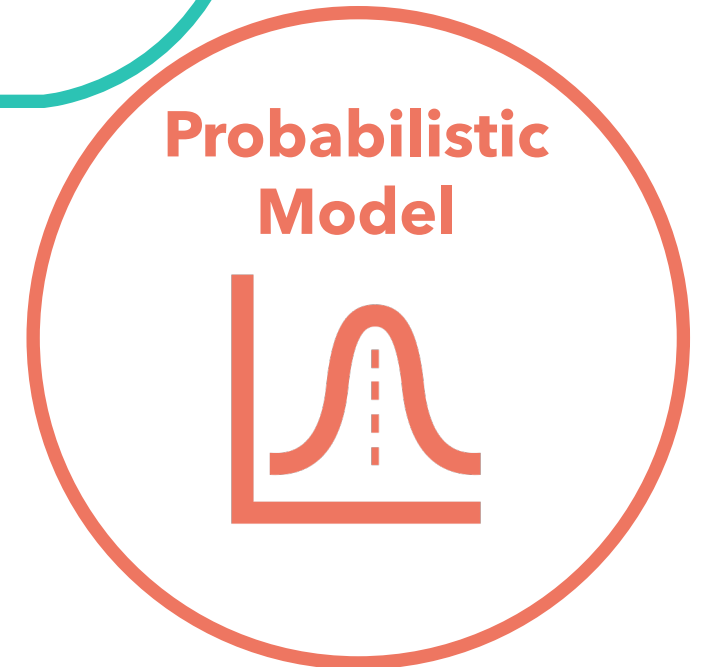
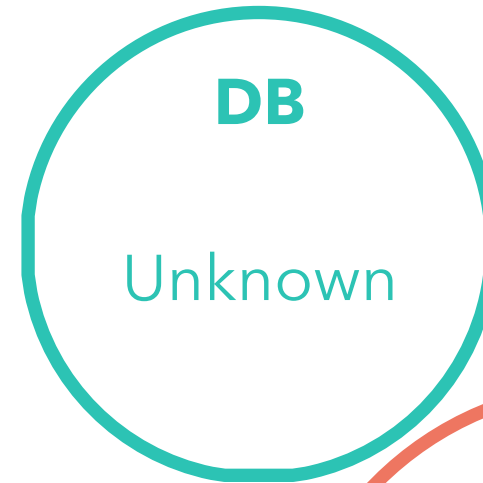
Answering Questions: DB vs Probabilistic Model

Employment

ID	Name	Workplace
46393	Mark Zuckerberg	Meta
47934	Andrew Bosworth	Meta
65849	Sundar Pichai	Google
64267	Satya Nadella	Microsoft



Where does John Doe work?



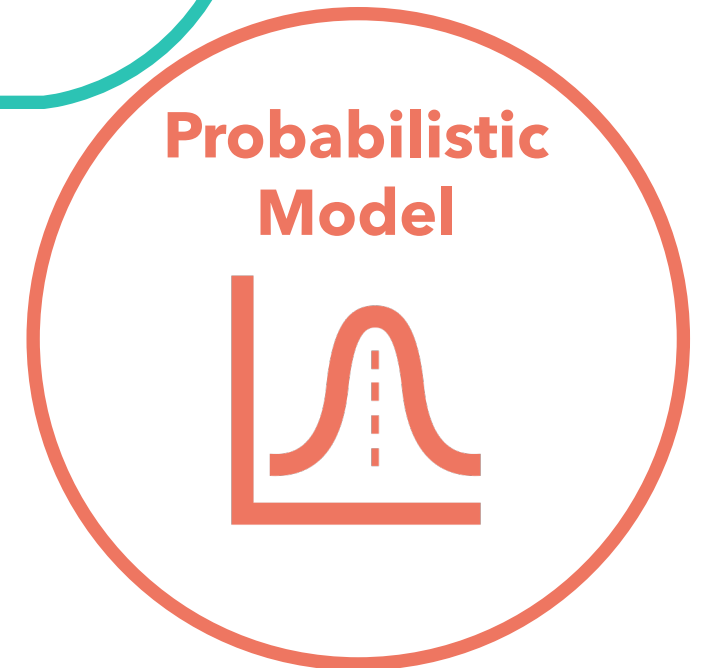
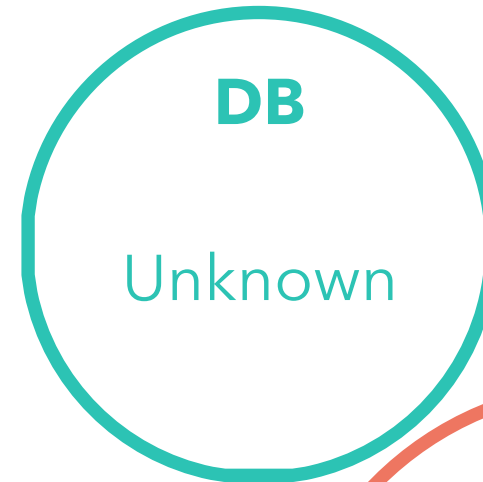
Answering Questions: DB vs Probabilistic Model

Employment

ID	Name	Workplace
46393	Mark Zuckerberg	Meta
47934	Andrew Bosworth	Meta
65849	Sundar Pichai	Google
64267	Satya Nadella	Microsoft



Where does John Doe work?



Answering Questions: DB vs Probabilistic Model

Employment

ID	Name	Workplace
46393	Mark Zuckerberg	Meta
47934	Andrew Bosworth	Meta
65849	Sundar Pichai	Google
64267	Satya Nadella	Microsoft



Where does John Doe work?

DB

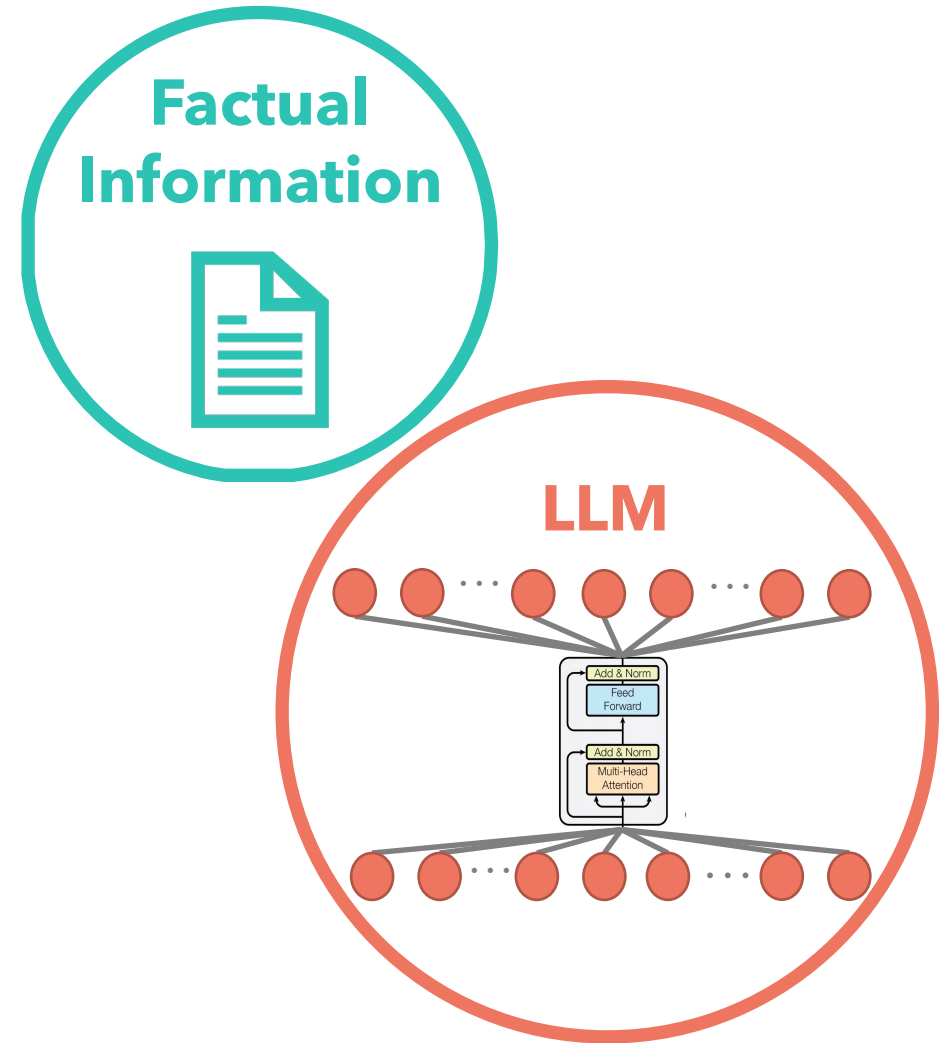
Unknown

**Probabilistic
Model**

Meta (50%)
Google (25%)
Microsoft (25%)

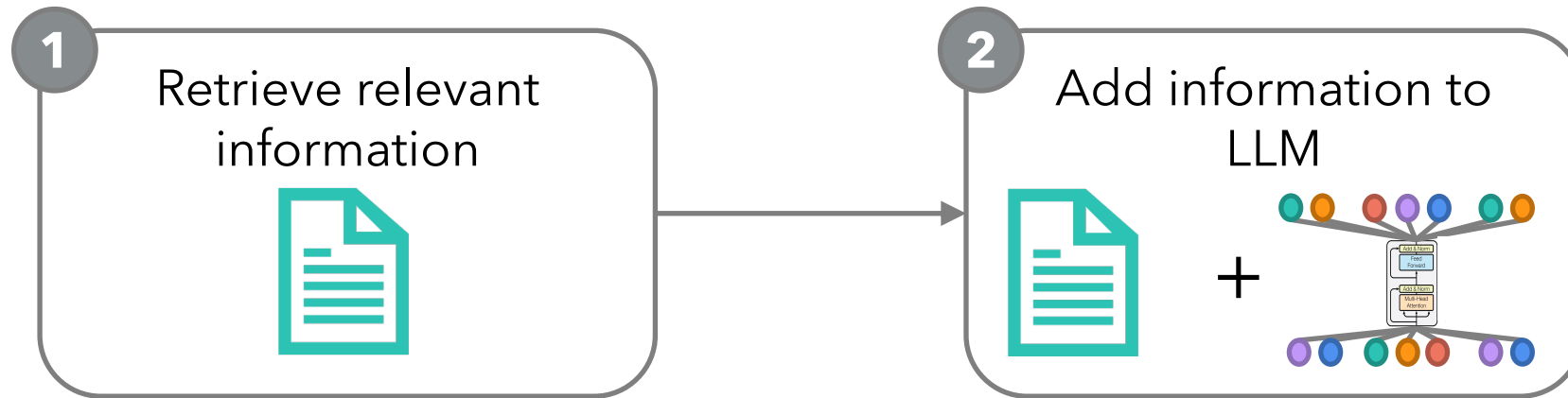
LLMs are not a source of factual information

- They are *probabilistic models* of factual information
 - Inconsistent source of information



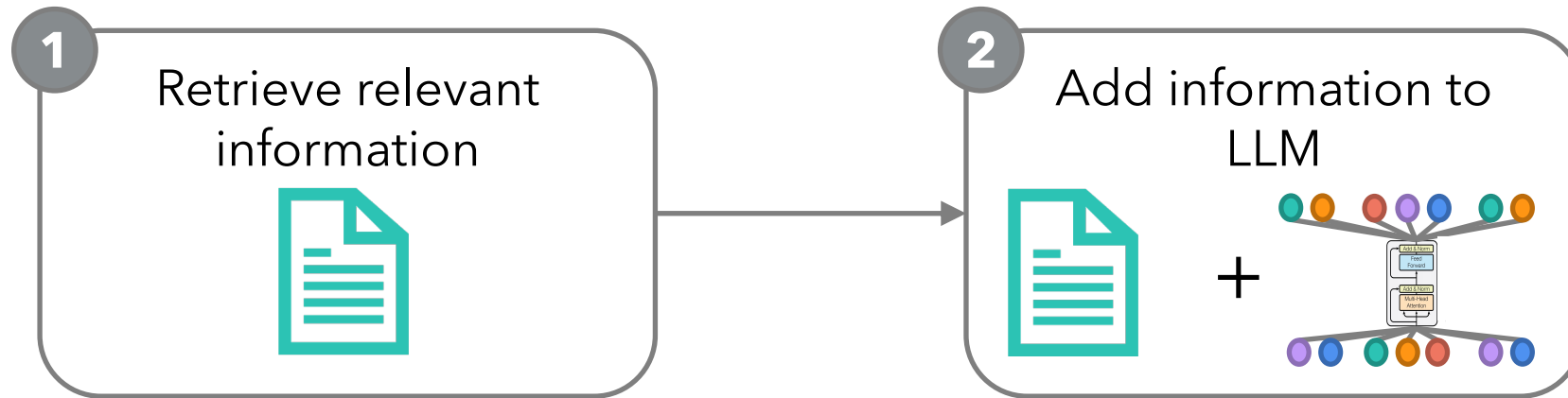
Retrieval-based LLMs: current efforts towards reducing inaccuracies in LLMs

- Adding *reliable* information to LLMs



Retrieval-based LLMs: current efforts towards reducing inaccuracies in LLMs

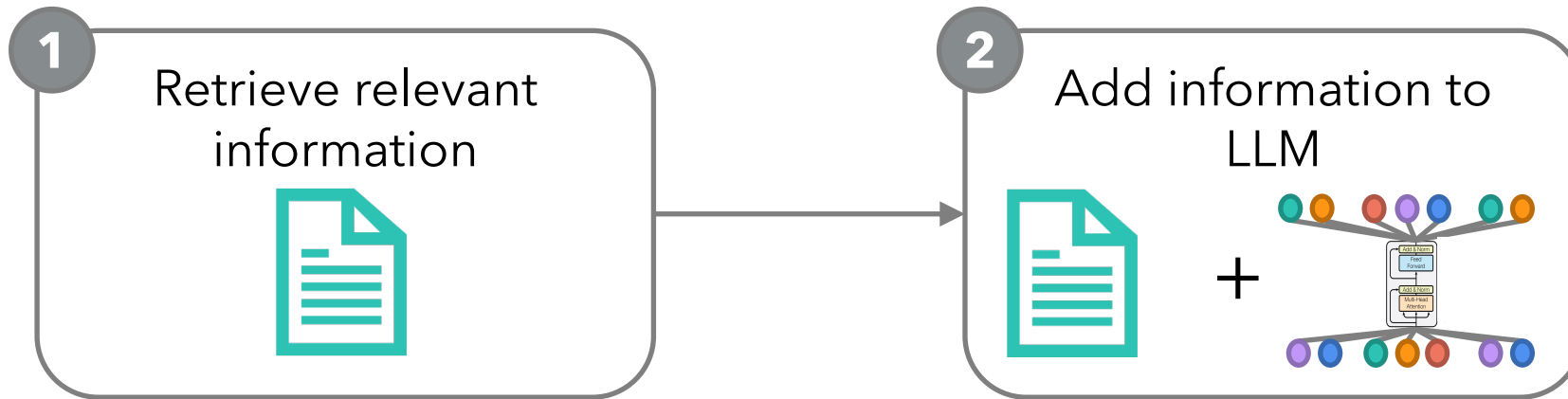
- Adding *reliable* information to LLMs



*Let's assume we have a good retriever

Retrieval-based LLMs: current efforts towards reducing inaccuracies in LLMs

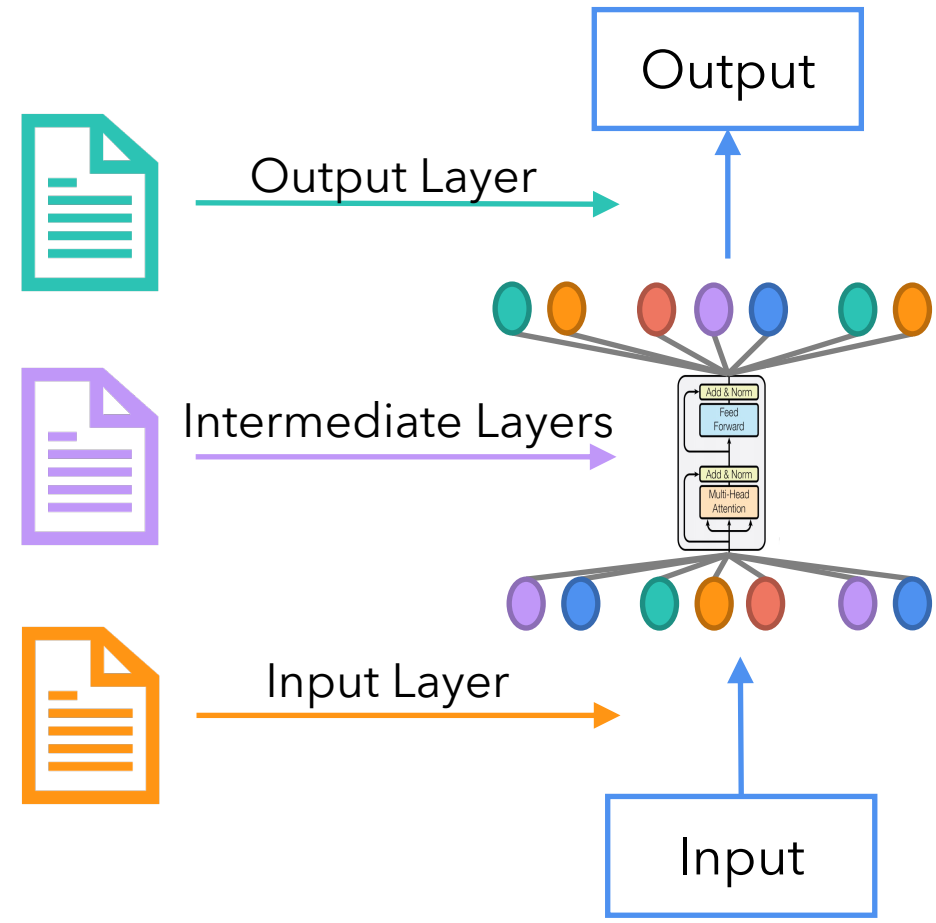
- Adding *reliable* information to LLMs



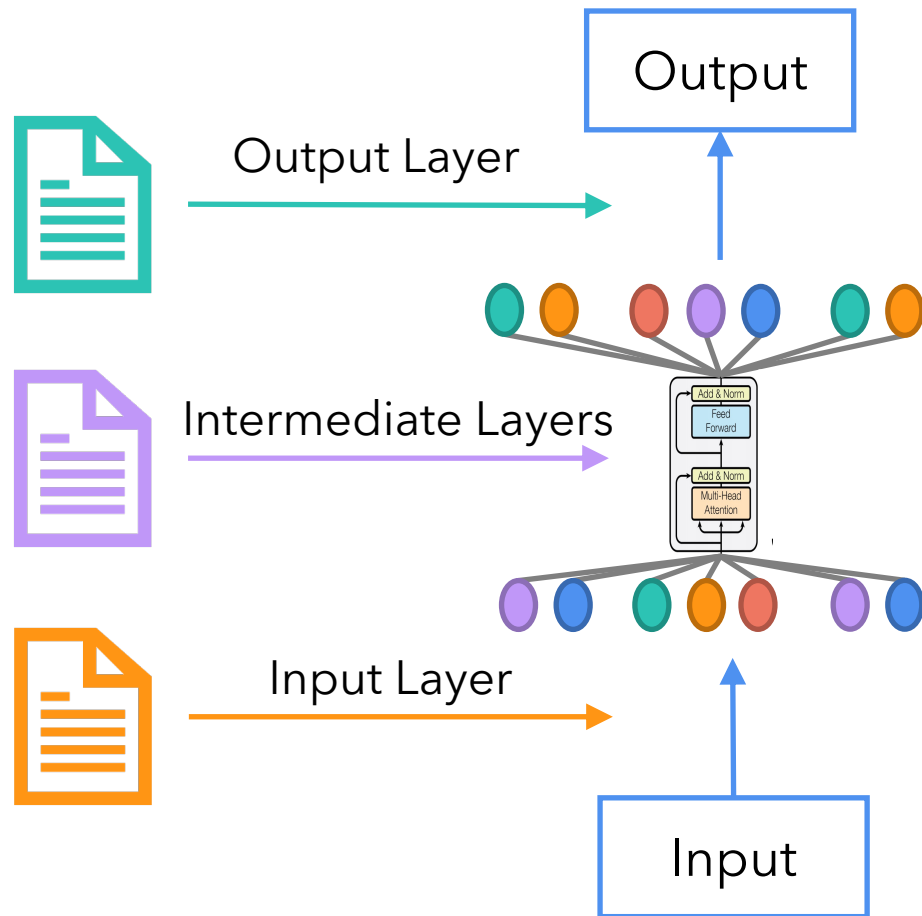
*Let's assume we have a good retriever

How to add this information to LLMs?

Where to add information in the LLM pipeline



How to add information in the LLM pipeline



During generation, predicting the next token with relevant information

Relational Memory Augmented Language Models, 2022

Changing the architecture to incorporate relevant information

Improving language models by retrieving from trillions of tokens, 2022

Concatenating relevant information with input

In-Context Retrieval Augmented Language Models, 2023

The problem with retrieval-based LLMs

- Might not be able to eliminate *all inconsistencies* due to deeper misconceptions or inaccurate learned patterns
- Ex: LLM does not understand transitive property (is-a)
 - Leads to inaccurate reasoning

Transitive Property:

Adam is-a Animal because Adam is-a Person and Person is-a Animal



Is Adam an animal?



No, Adam is not an animal. The name "Adam" typically refers to a male human being according to various cultural and religious traditions. It is not used to describe an animal.

8/27/2023 : Inaccurate result on GPT-3.5

We have seen the problem of dealing with inconsistent information before...

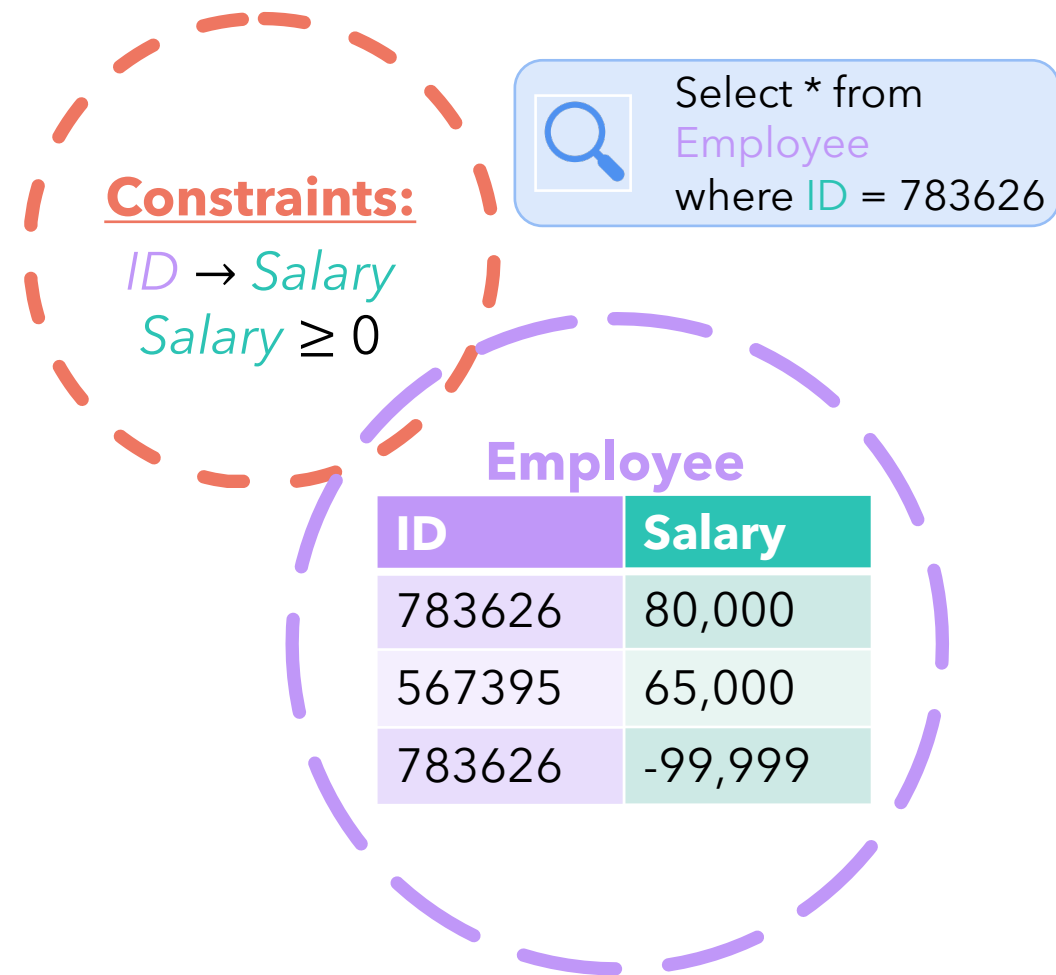
- Data management community has spent 4 decades solving this problem

We have seen the problem of dealing with inconsistent information before...

- Data management community has spent 4 decades solving this problem

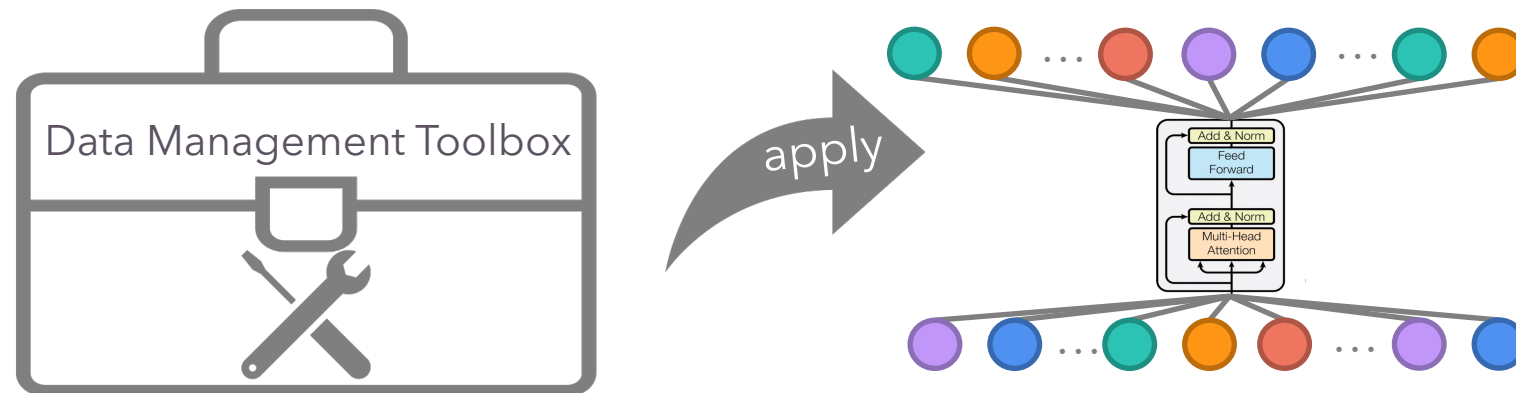
Problem Setup:

- Given an **inconsistent data source**, **declarative constraints**, and a **query**
- Goal: give consistent information that complies with declarative constraints



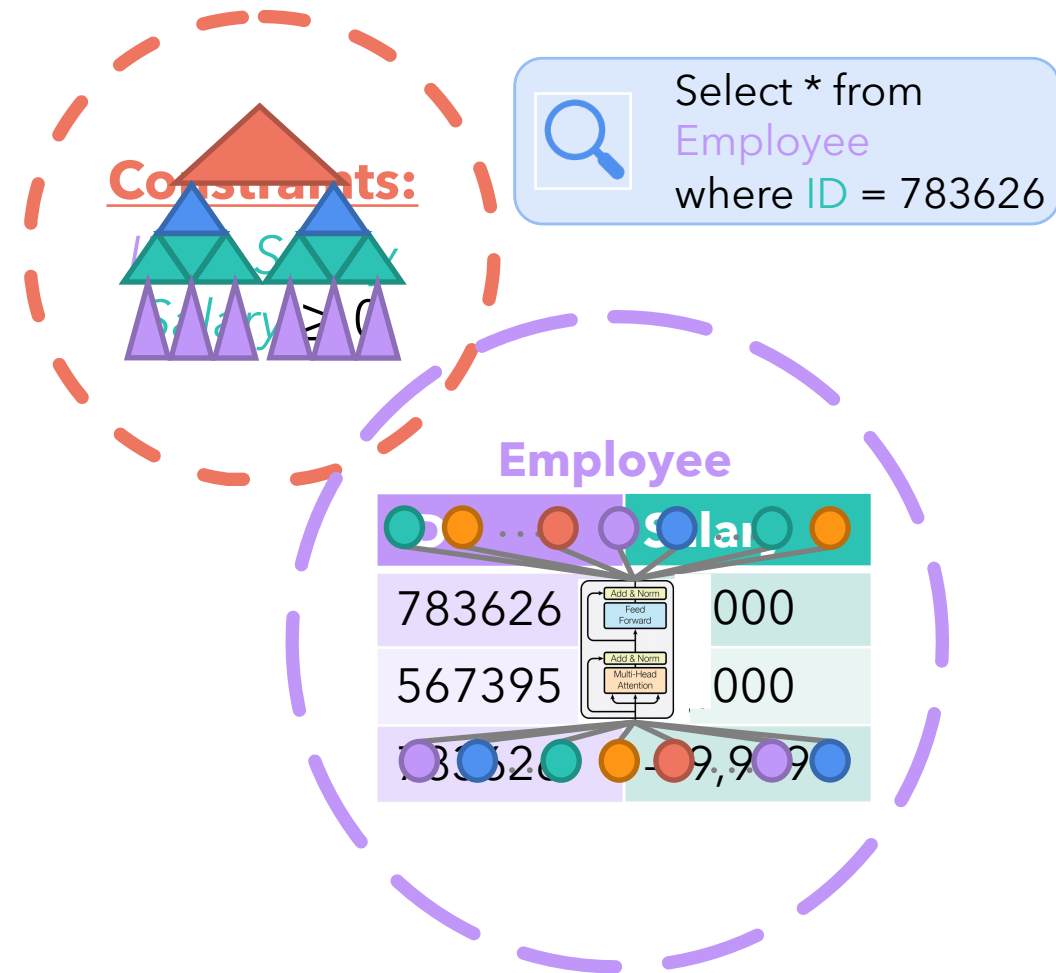
Our goal for this talk

- Investigate *how we can apply data management techniques* to solve the problem of incorrect and self-contradictory results in LLMs



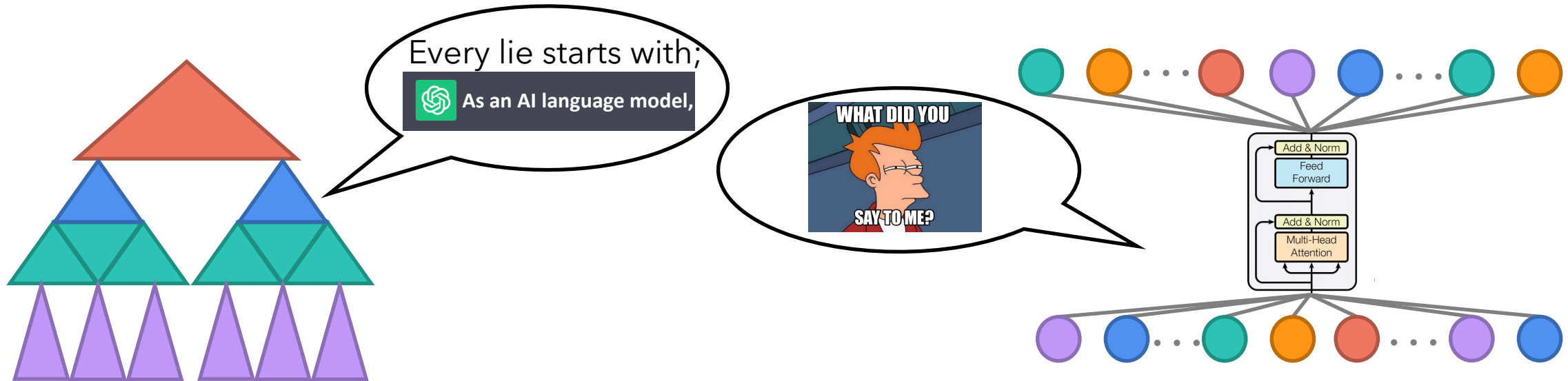
Building a data management framework for LLMs

- Analogy between **inconsistency** in data management and **inaccuracy** in LLMs
 - Need the same components!
- Source of inconsistent information are **LLMs**
- Semantic properties and constraints can be represented through **ontologies**
 - Benefits: publicly available across various domains, easy to modify & maintain
- Queries are **textual input** to LLMs



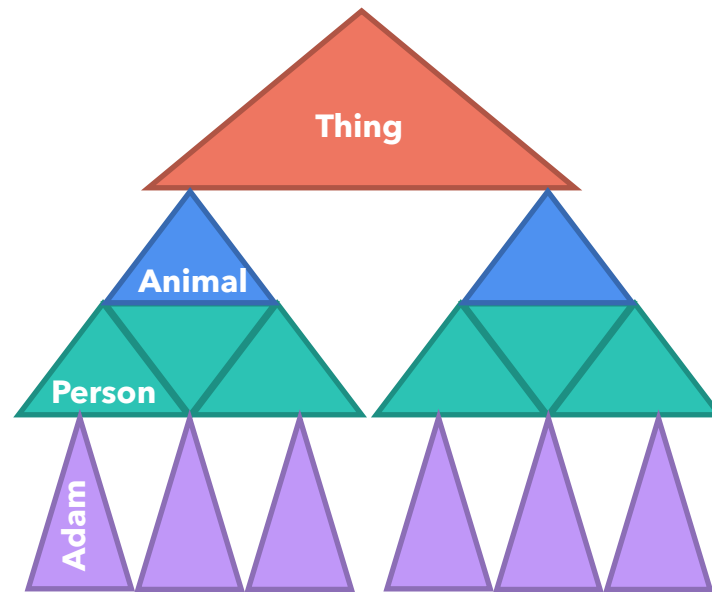
Challenge in using ontologies with LLMs

- Ontologies and LLMs don't speak the same language
- Finding alignment between ontology constraints and LLMs' continuous representation.

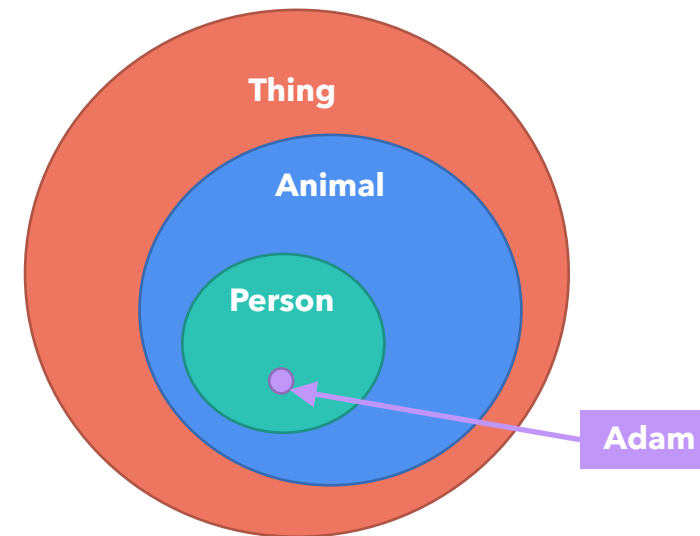


Create a better representation to facilitate better communication between ontology and LLM

- Goal: Preserve structural properties and relationships of declarative constraints in the embedded space
- How: using geometric embeddings



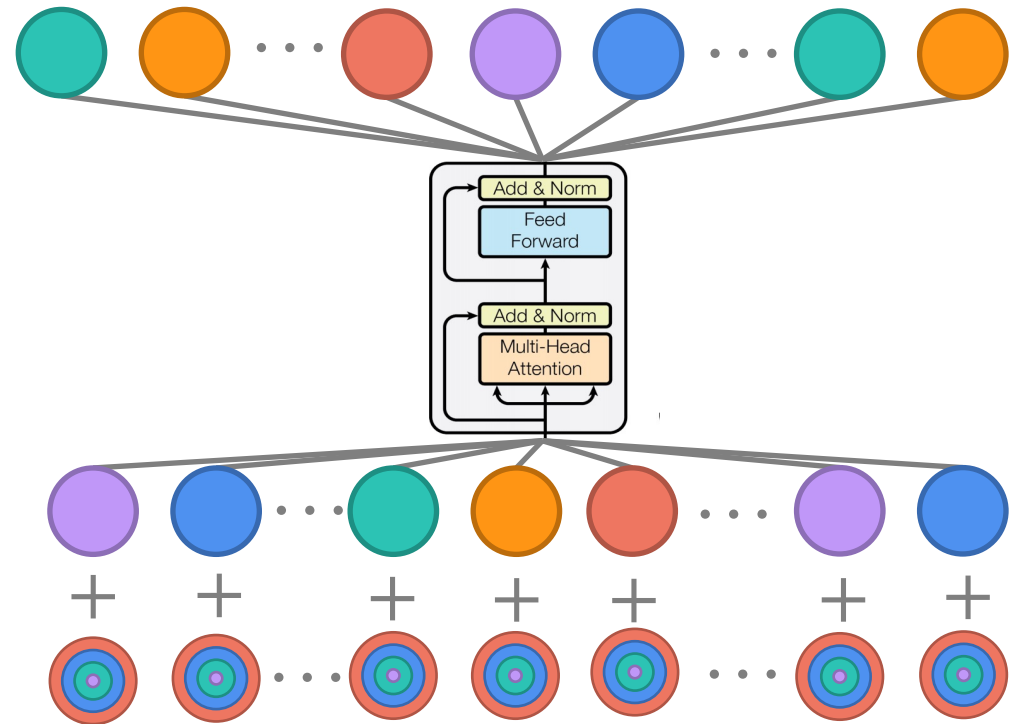
Ontology



Ball Geometric Embeddings

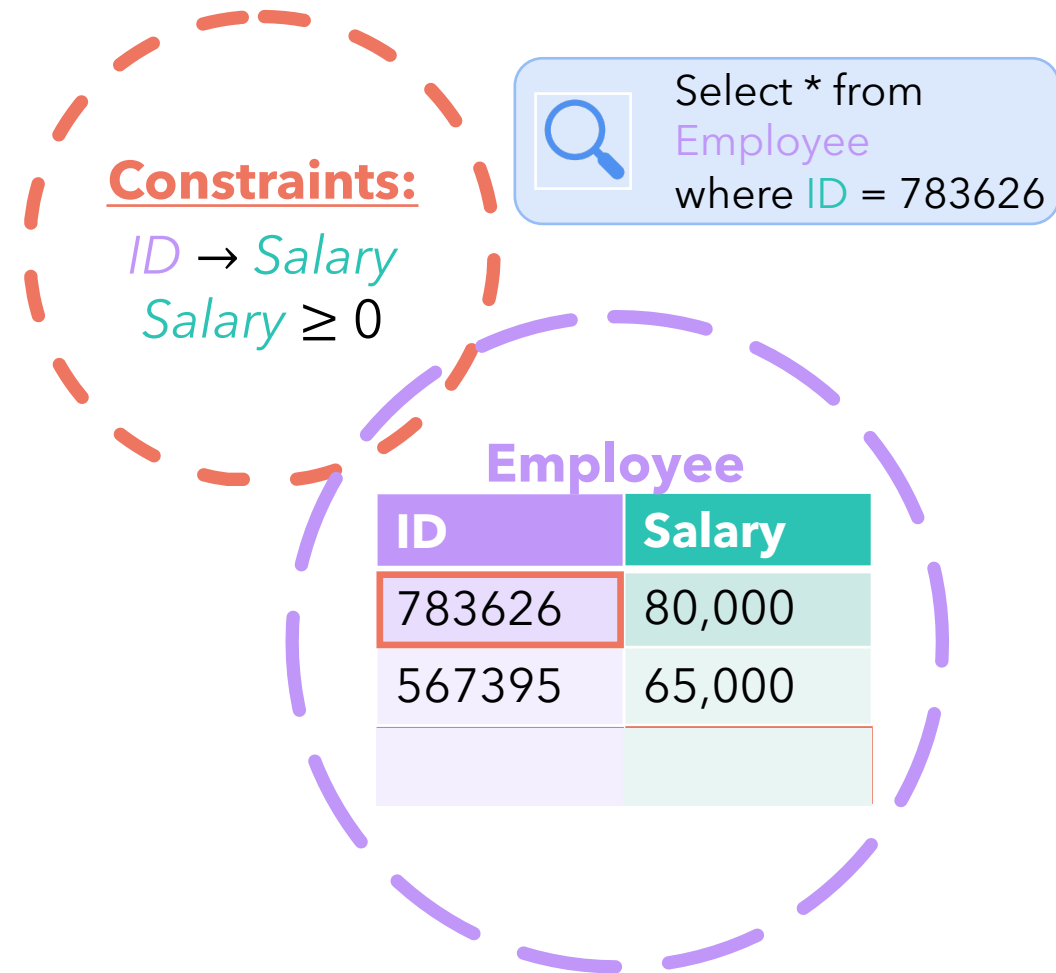
Leveraging constraints during pretraining

- Goal: Constrain the space of generalized knowledge
- How: Pretraining with constraints using geometric embeddings
- Challenge: pretraining is expensive
 - Finetune



Data management toolbox: data cleaning

- Data cleaning: removing inconsistent information from data
- Goal: find the minimal number of repairs needed for the data source to comply with declarative constraints



Challenge in adapting data cleaning to LLMs

Traditional Databases

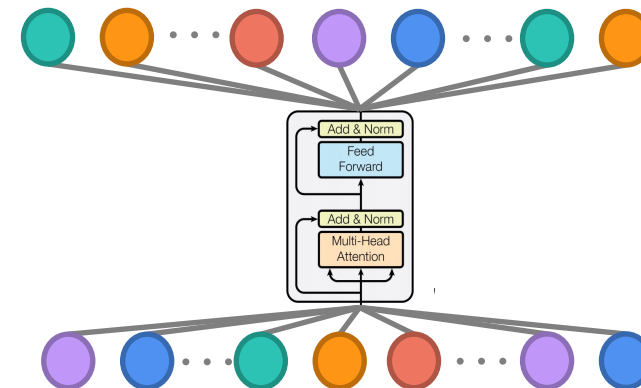
- Explicit form of data

Employee

ID	Salary
783626	80,000
567395	65,000
783626	-99,999

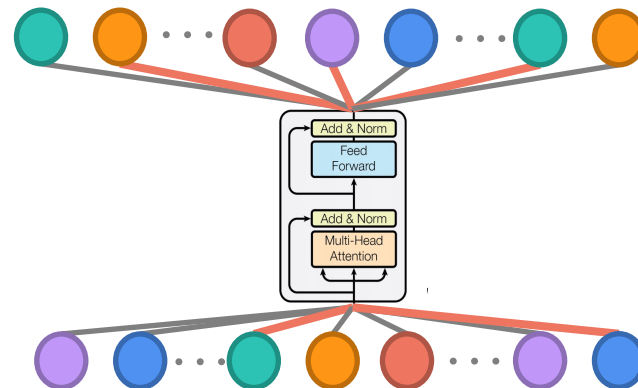
Large Language Models

- Implicit form of data through weights



Repairing inconsistent information in LLMs

- Goal: update weights in LLMs until information is consistent with constraints
- Recent efforts: editing information in LLMs by directly modifying weights
- How: Extend this method by finding the *minimal number of weights* that solves inconsistency of a constraint
- Challenge: LLM weight editing research is in early stage



Data management toolbox: consistent query answering

- Consistent query answering: providing consistent results over inconsistent data
- Goal: Modify the query until results are comply with declarative constraints

1 Select * from Employee
where ID = 783626

ID	Salary
783626	80,000
783626	-99,999

2 Select * from Employee
where ID = 783626
and Salary ≥ 0

ID	Salary
783626	80,000

Constraints:
 $ID \rightarrow Salary$
 $Salary \geq 0$

Select * from Employee
where ID = 783626

Employee

ID	Salary
783626	80,000
567395	65,000
783626	-99,999


Challenges in adapting consistent query management for LLMs

Traditional Databases

- Clear separation between query language and information source

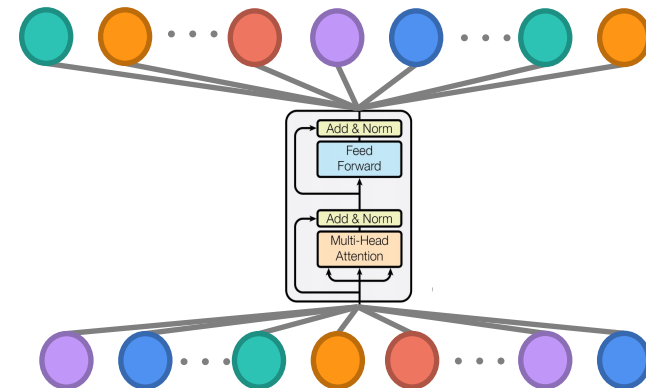
Employee

ID	Salary
783626	80,000
567395	65,000
783626	-99,999

 Select * from Employee where ID = 783626

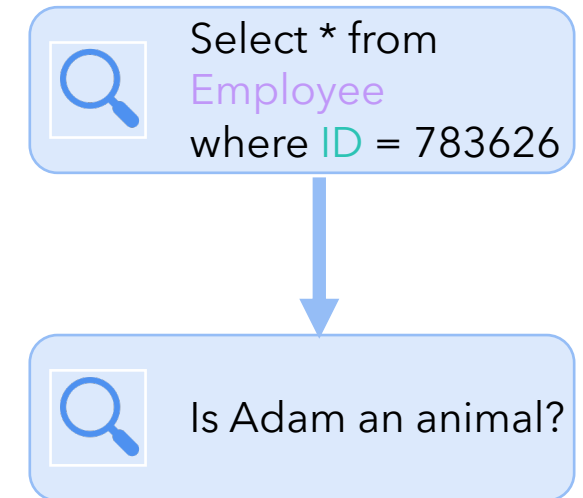
Large Language Models

- Difficulty in separating linguistic patterns from actual information



Recall: queries in data management framework is akin to prompting to LLM

- Techniques for prompting
 1. Chain of Thought
 2. Prefix Tuning
 3. Tree of Thought
 4. Graph of Thought
- Prompting tries to get answers out of the LLM by expanding the input to model (query)



Chain of Thought Prompting

- Composing prompt with demonstration of intermediate steps
 - encourages reasoning

To determine if 977 is a prime number, we can use the method of trial division.

A number is prime if it has only two positive divisors: 1 and itself. If n is not prime, then it must have a factor less than or equal to \sqrt{n} .

For $n = 977$, $\sqrt{977}$ is just a bit over 31 (since $31 \times 31 = 961$).

So, we need to check if 977 is divisible by any prime numbers less than or equal to 31. The prime numbers in that range are:

2, 3, 5, 7, 11, 13, 17, 19, 23, 29, and 31.

977 is odd, so it's not divisible by 2.

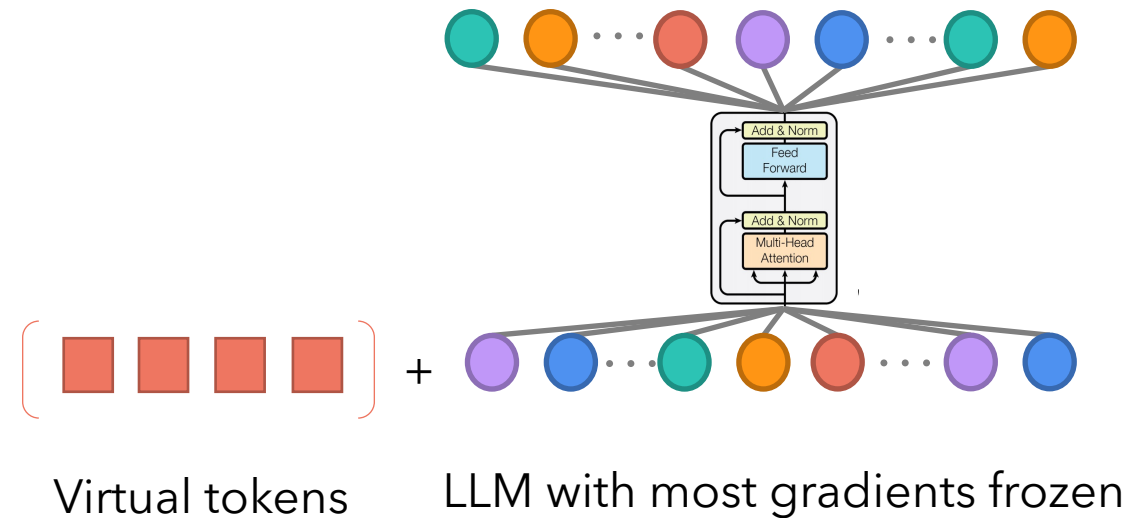
Now, let's test the other primes:

- $977 \bmod 3$ is not 0.
- $977 \bmod 5$ is 2 (because the last digit is neither 0 nor 5).
- $977 \bmod 7$ is not 0.
- $977 \bmod 11$ is not 0.
- $977 \bmod 13$ is not 0.
- $977 \bmod 17$ is not 0.
- $977 \bmod 19$ is not 0.
- $977 \bmod 23$ is not 0.
- $977 \bmod 29$ is not 0.
- $977 \bmod 31$ is not 0.

Since 977 is not divisible by any of these primes, and we don't have any other integers less than $\sqrt{977}$ to test, we can conclude that 977 is a prime number.

Prefix Tuning

- Continuous method for prompt tuning (otherwise discrete since it is in the word space)
 - virtual set of tokens



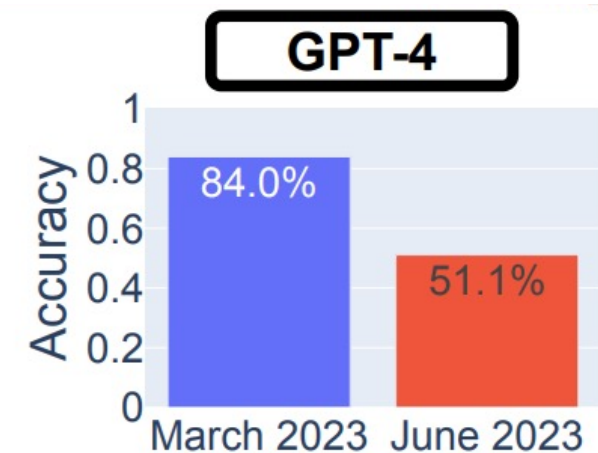
CoT prompting is not reliable with minor model changes

- Example with Chain of Thought Prompting
 - Same idea applies to prefix tuning (needs retraining)

LLM Service	GPT-4		Δ
	Prompting method		
Eval Time	No CoT	CoT	
Mar-23	59.6%	84.0%	+24.4%
Jun-23	50.5%	49.6%	-0.1%

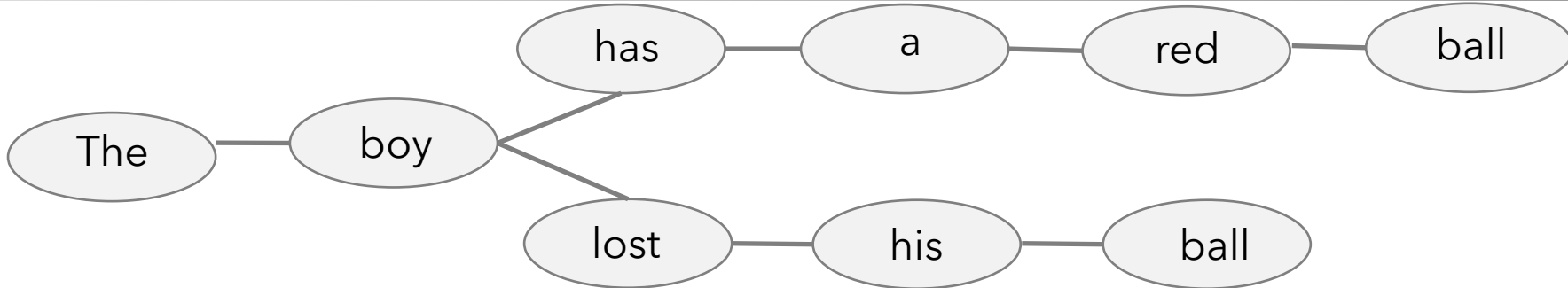
GP4 produces contradictory results over time over the same set of questions

Prime vs. Composite performance over 1000 samples



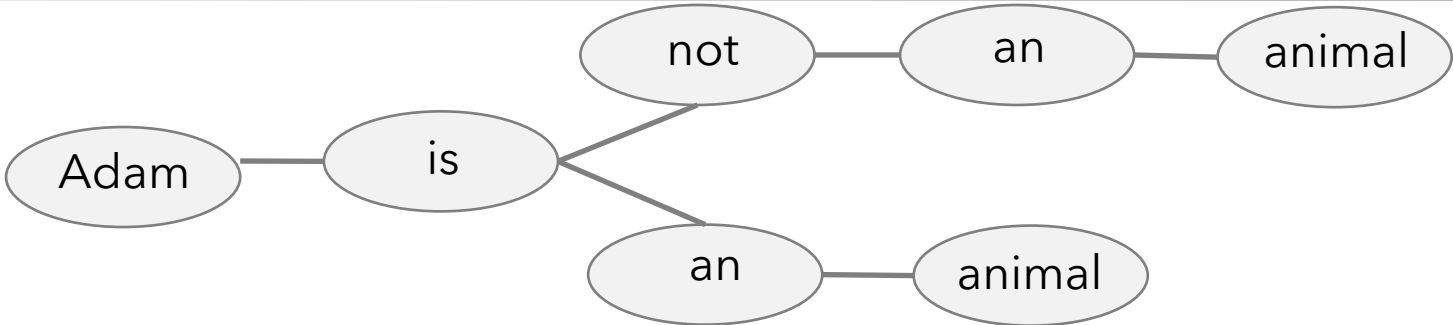
Recent efforts towards applying constraints in the output layer

- Satisfying **lexical** constraints in output sequence
 - Ex: Write a sentence using the words ball, red, and boy
- Generate a set of candidate sequences and return the one with the largest probability of satisfying the constraint

Search Tree	Likelihood	Select
	0.98	✓
	0.61	✗

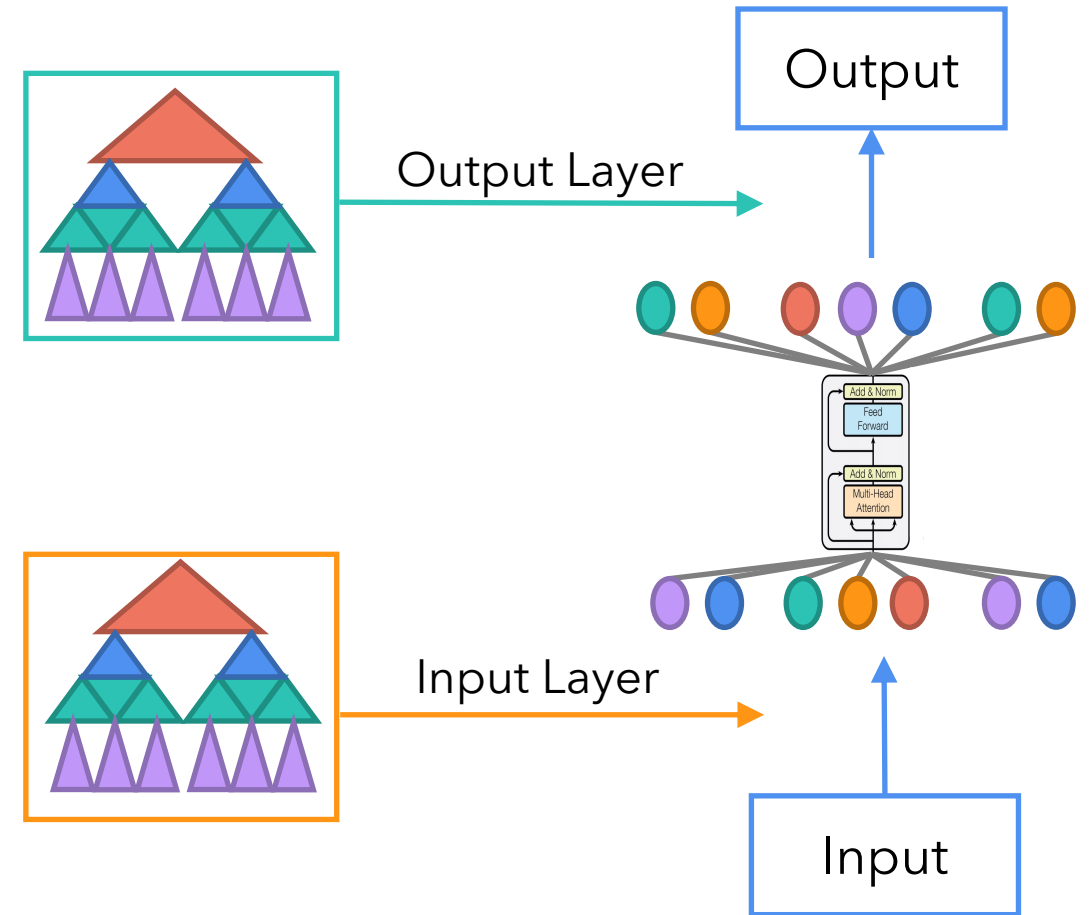
Consistent querying answering by applying constraints in the output layer

- Goal: ensuring the output sequence complies with constraints
- How: Extending lexical constraint method for **semantic** constraints by finding the *minimum* set of constraints that imply the entire constraint set.

Search Tree	Likelihood	Select
 <p>The diagram shows a search tree starting with 'Adam' in an oval, connected to 'is' in an oval. From 'is', two branches emerge. The upper branch leads to 'not' in an oval, which is connected to 'an' in an oval, which is connected to 'animal' in an oval. The lower branch leads to 'an' in an oval, which is connected to 'animal' in an oval.</p>	0.05	X
	0.91	✓

Ongoing work

- Analysis on consistent query answering in LLMs
 - input layer vs output layer
 - complexity of constraints





Thank You!