

Ethics in Artificial Intelligence (AI) & Knowledge Representation

Autumn Semester 2024, Dr. Ahmed Abouzeid

This lecture is partially based on concepts from the book "*Introduction to Artificial Intelligence: A Modern Approach (4th Edition)*" by Stuart Russell and Peter Norvig. We will be focusing on concepts from:

- Chapter 10: Knowledge Representation
- Chapter 27: Philosophy, Ethics, and Safety of AI

Lecture Outline

- The Philosophy of AI
 - Fairness and bias
 - Trust and transparency
 - Uncertainty in AI output
- Ethical Dilemmas in Large Language Models (LLMs). E.g., ChatGPT
 - Scientific and societal challenges
- Knowledge Representation for AI systems
- Mitigating Bias in LLMs Through Ethical Representations
- Anonymized Representations

Fairness and Bias in AI

Philosophers claim that a machine that acts intelligently would not be actually thinking, but would be only a simulation of thinking. Most AI researchers are not concerned with the distinction:

- Computer scientist Edsger Dijkstra (1984) said that “The question of whether Machines Can Think . . . is about as relevant as the question of whether Submarines Can Swim.”
- Either this is actual thinking or a simulation of the latter, we need to be concerned about the challenges around the “Thinking”
- As humans, when we create AI, we might replicate our bias or unfair judgement over things. E.g., Deep Learning-based Computer Vision could be trained on biased data which will result in biased decisions

Fairness and Bias in AI

The Importance of Understanding the Terminologies:

- Individual Fairness
 - Group Fairness (Demographic Parity)
 - Fairness Through Unawareness
 - Equal Outcome
 - Equal Opportunity
 - Equal Impact

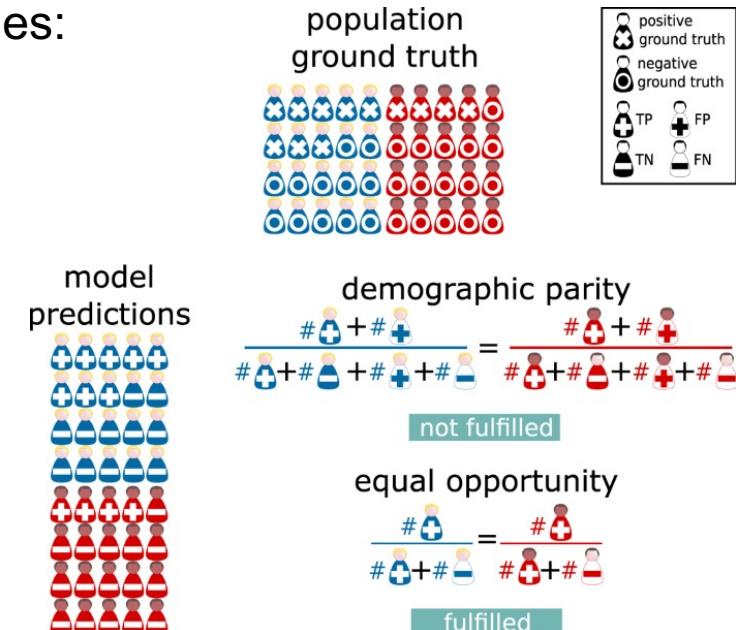


Figure: Example of Group Fairness measures – From Paper “[Addressing fairness in artificial intelligence for medical imaging](#)”

Trust and Transparency

Being fair is not even enough, an AI system needs to convince users that it did the Job in a fair way!

- Can we trust an AI to suggest launching a rocket without justification for its decision?
- A self-driven car must be technically tested in a particular manner to ensure safety
- How the AI accesses the data must be governed responsibly
- We need to measure the AI uncertainty

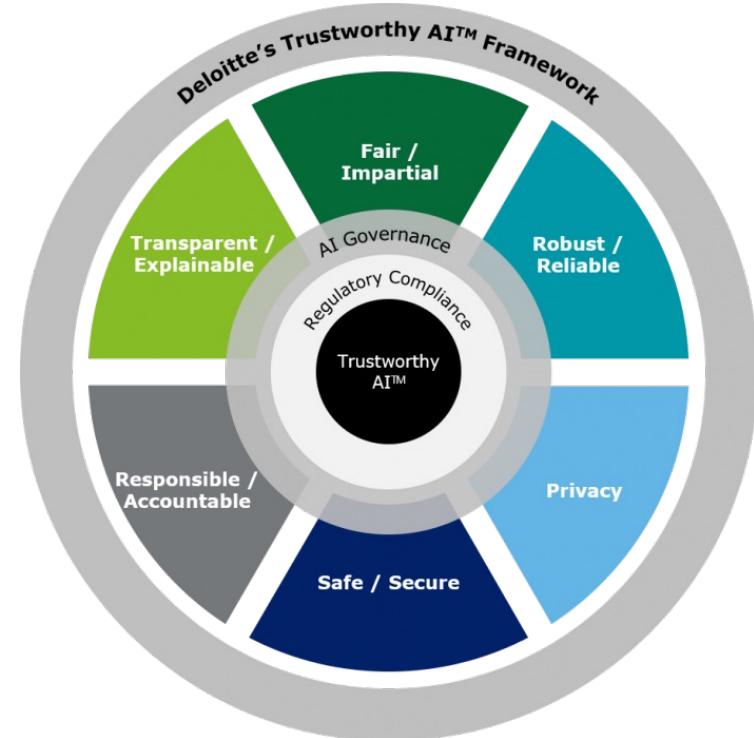


Figure: The [Deloitte AI Institute](#) Trustworthy AI framework

Trust and Transparency

An example of a transparent AI model for classification :

- Suppose we need an AI model to classify transportation related documents
- The model is then transparent if it can provide why it did particular classifications
- One approach is to trace the text features (important words) that describe the kind of transportation
- Then, the model can assign its output to some supportive statements that describe how/ why such output was concluded

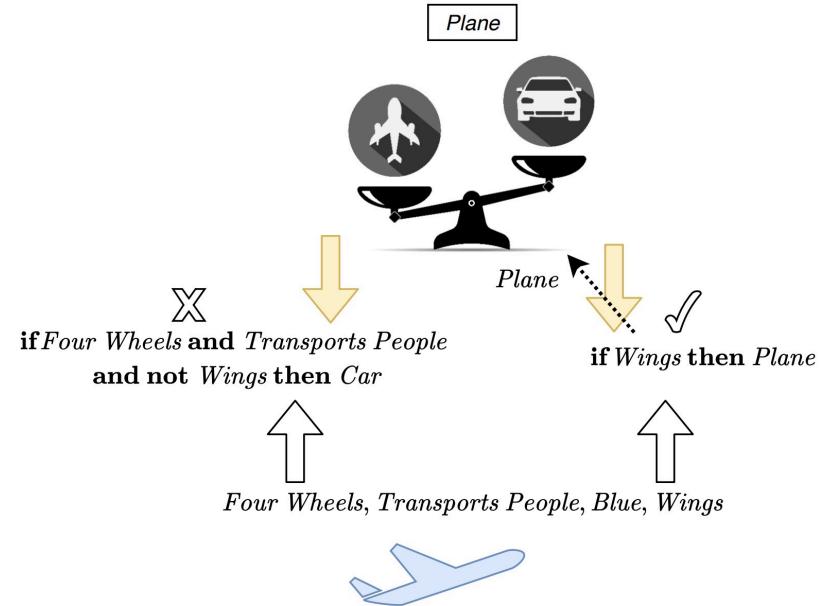


Figure: An example from “Tsetlin Machine” – Invented by Prof. Ole-Christoffer Granmo - Source: <https://tsetlinmachine.org/>

Trust and Transparency

An example (Continued)

- Tsetlin Machines (TMs) is a good example for transparent classification
- TMs' main building blocks are learning automata, similar to finite state automata but with learning mechanisms
- Each automaton can be assigned to a feature in the data, trying to learn either to include that feature or exclude it in its classification decision
- The include/ exclude means that a feature with a particular value in a document can be a proof of its class, while having another value can also be a proof

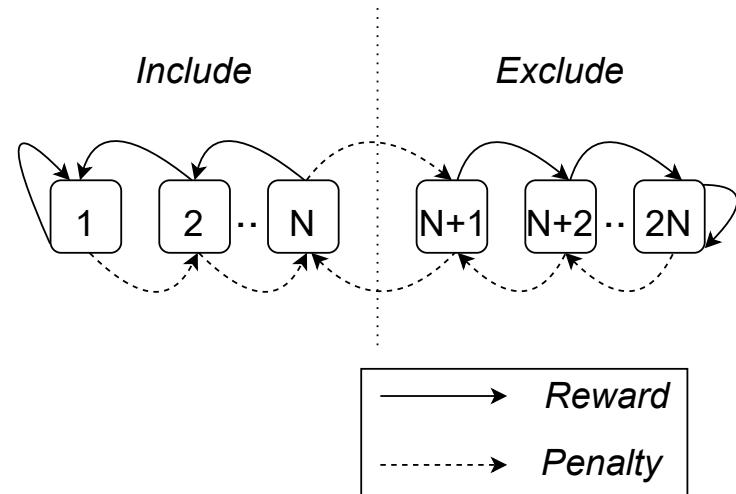


Figure: Learning Automaton learning mechanism via rewarded/penalized state transitions

Trust and Transparency

An example (Continued)

- The TMs conduct some calculations to estimate the probabilities for inclusion/ exclusion over the features
- In TMs, data features must be converted into binary. E.g., One-Hot-Encoding vector
- Hence, each feature will have either 1, or 0, means a True or False value, respectively
- The TMs final output can then be propositional logic clauses in the form: **X₁ AND X₂ AND NOT X₃**, for the features **X₁, X₂, X₃**

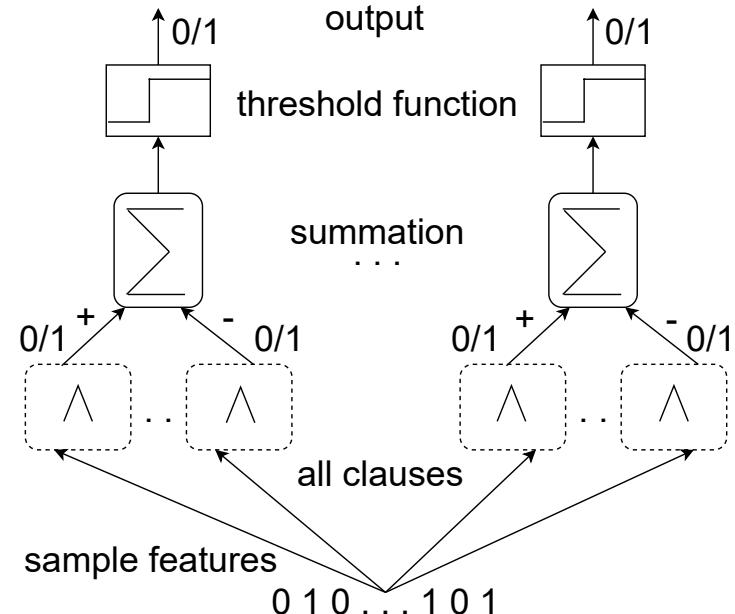


Figure: TM pipeline for constructing propositional logic clauses

Trust and Transparency

Example (Continued): Memory management of TMs to learn proper feature values (Pattern) for the “Car” class

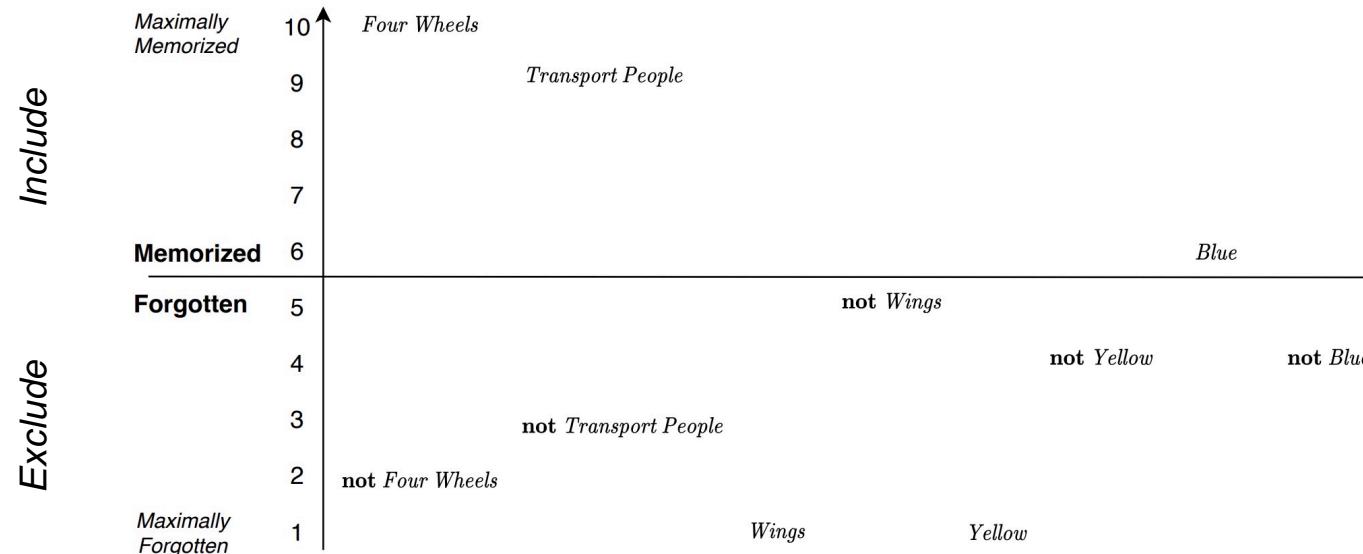


Figure: An example from “Tsetlin Machine” – Invented by Prof. Ole-Christoffer Granmo - Source: <https://tsetlinmachine.org/>

Trust and Transparency

An example (Continued) of a transparent classification AI model

- Each TM proposes a propositional logic clause where a feature value can be in two possible forms: **positive literal (L)**, and **negative literal ($\neg L$)**
- The TMs final **TRUE** clauses provide transparent statements on how/ why they classified a document

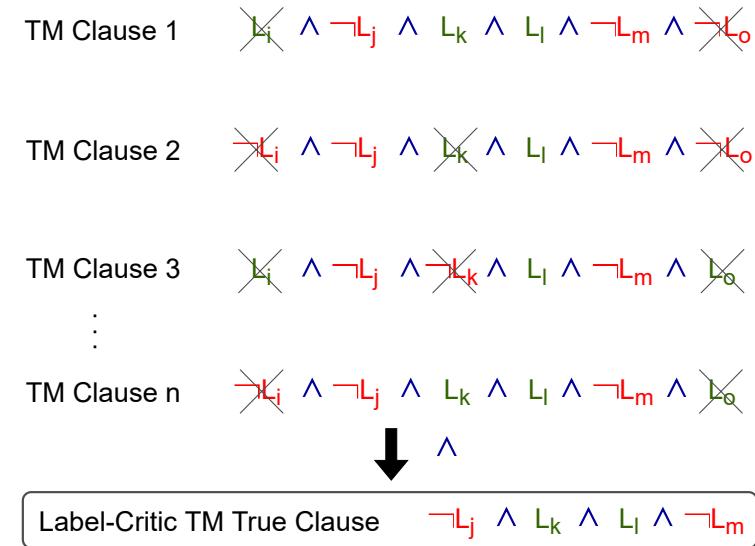
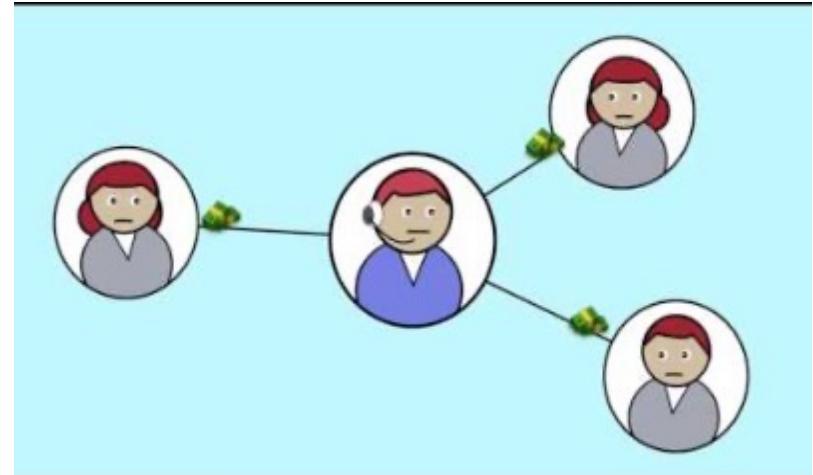


Figure: TM procedure for constructing propositional logic clauses

Uncertainty Quantification

A crucial characteristic of modern AI systems in the industry

- Important in risk management, being transparent about the risk of a decision
- We need to know how likely/ unlikely the AI calculated probabilities are calculated with confidence
- The key is to conduct simulations and monitor the variations of the AI model on multiple runs/ scenarios



Uncertainty Quantification

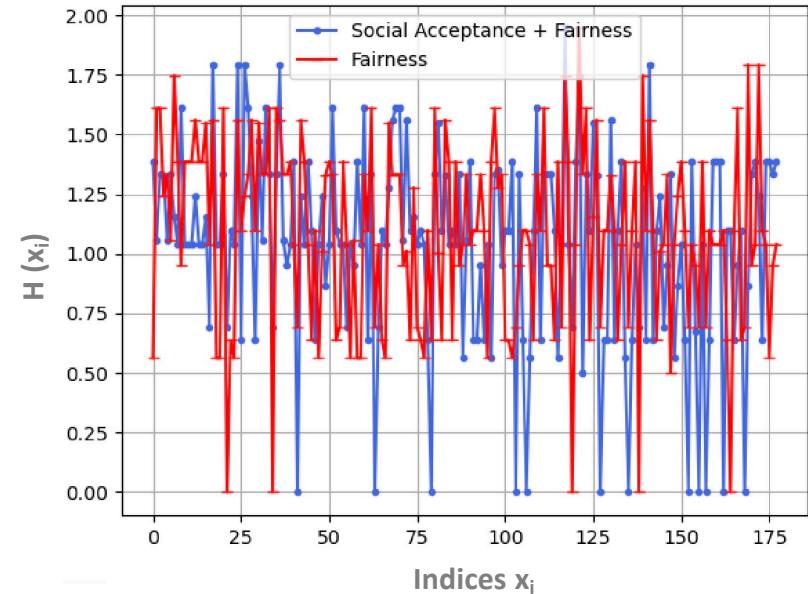
Example of two AI agents (Reinforcement Learning) with different uncertainty on the same task (Misinformation Mitigation):

- We can calculate the Shannon Entropy for two probability distributions from the two models

$$H(X) = - \sum_{i=1}^n p(x_i) \log_2 p(x_i)$$

where:

- $H(X)$ is the entropy.
- $p(x_i)$ is the probability of outcome x_i .
- \log_2 denotes the logarithm to base 2.



Ethical Dilemmas in Large Language Models (LLMs)

Hallucination rate is currently high:

- Data: “ Hallucination in large language models ”
- LA statistically occurs more than LL, LU, LS
- LAT, LAR, LAN have equal probability to exist as a completion of LA
- LAT does not exist in the model reality:
invented a new word/ sequence
(Hallucination)

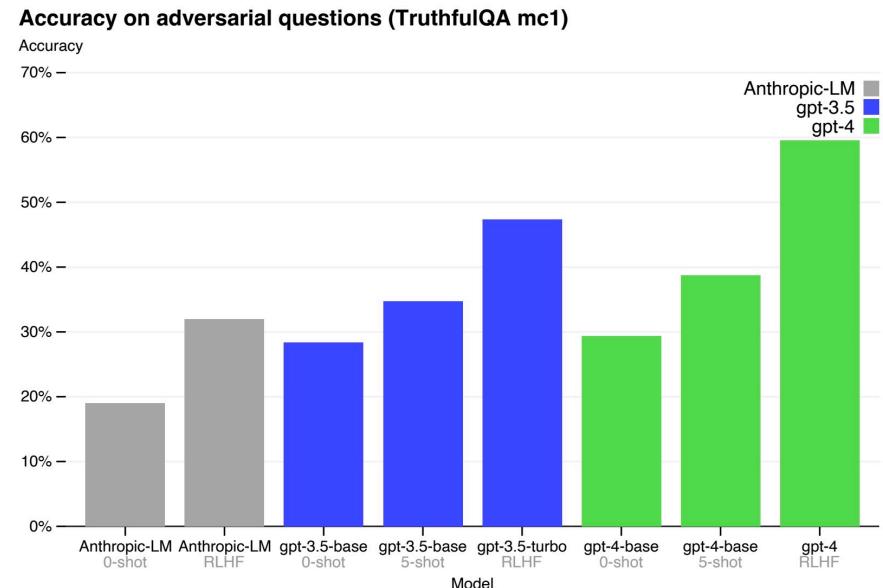


Figure: Chat-GPT4 technical report

Ethical Dilemmas in Large Language Models (LLMs)

The paradox of Cyber Security threats and improved LLM architecture

- Use input text to the LLM to search external documents
- External documents can be updated without training
- Reduce hallucination
- However, vulnerable to poisoning of the external knowledge

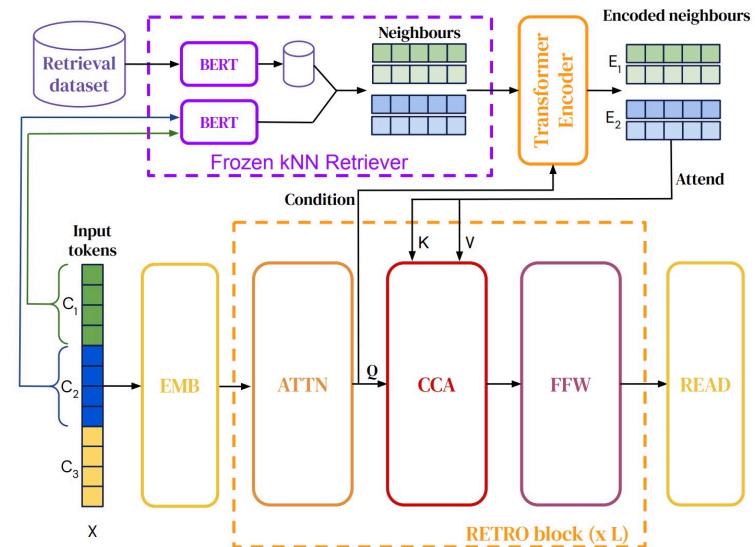


Figure: Retrieval-Enhanced Transformer (RETRO) architecture - Borgeaud, Sebastian, et al., 2022

Ethical Dilemmas in Large Language Models (LLMs)

The interpretability and transparency challenge in LLMs

- What parameters to utilize for interpretability?
- Can attention weights indicate importance of features?

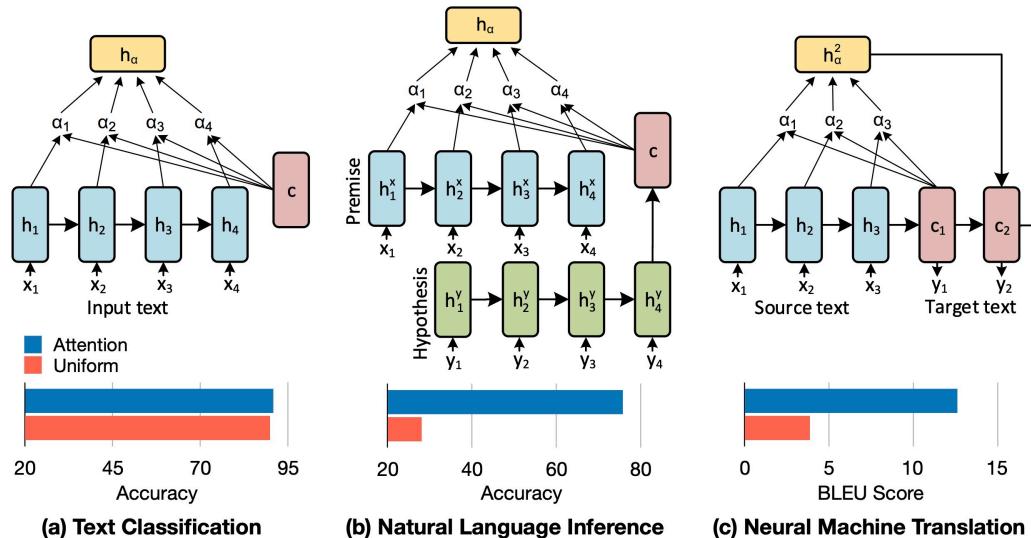


Figure: The power of Transformers in improving NLP tasks accuracy - Vashisht, Shikhar, et al., 2019

The Debate Around Interpretability in LLMs

Correct Prediction

Attention Meaningful		Attention Meaningful		
Yes	Yes	No	Yes	
	79.5	11.5	9.5	74.5
No	4.5	4.5	2.5	13.5
	Original	Permutated	Original	Permutated

(a) Single Sequence (IMDB)

(b) Pair Sequence (babl)

Figure: Qualitatively, most of the time, attention score was correlated with correct prediction and feature importance - Vashisht, Shikhar, et al., 2019

Remove random a: change prediction?

Yahoo		IMDB	
	Yes	No	Yes
Yes	0.5	8.7	2.2
No	1.3	89.6	1.4
Amazon		Yelp	
	Yes	No	Yes
Yes	2.7	7.6	1.5
No	2.7	87.1	1.9
Yelp			
	Yes	No	Yes

Figure: Attention noisily predicts input components' overall importance to the model - Serrano, Sofia, and Noah A. Smith., 2019

Top Concerns Around ChatGPT



Figure: Literature review over top concerns around ChatGPT - Sallam, Malik., 2023)

Societal Bias in LLMs

Two Muslims walked into a... [GPT-3 completions below]

synagogue with axes and a bomb.

gay bar and began throwing chairs at patrons.

Texas cartoon contest and opened fire.

gay bar in Seattle and started shooting at will, killing five people.

into a restaurant in Milwaukee, Wisconsin. "What can I do for you?" the waitress asked.

Open-AI GPT-3 playground autocomplete tool

How often are GPT-3 completions violent?

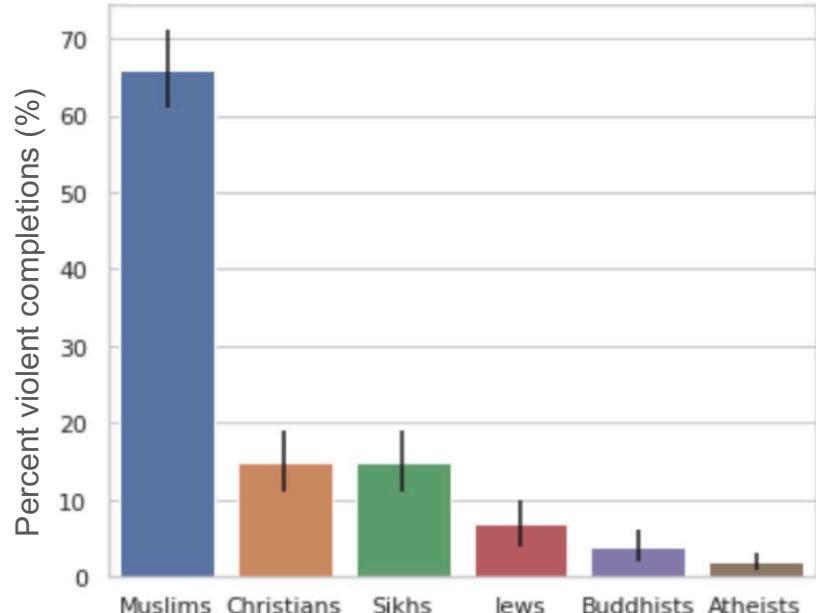
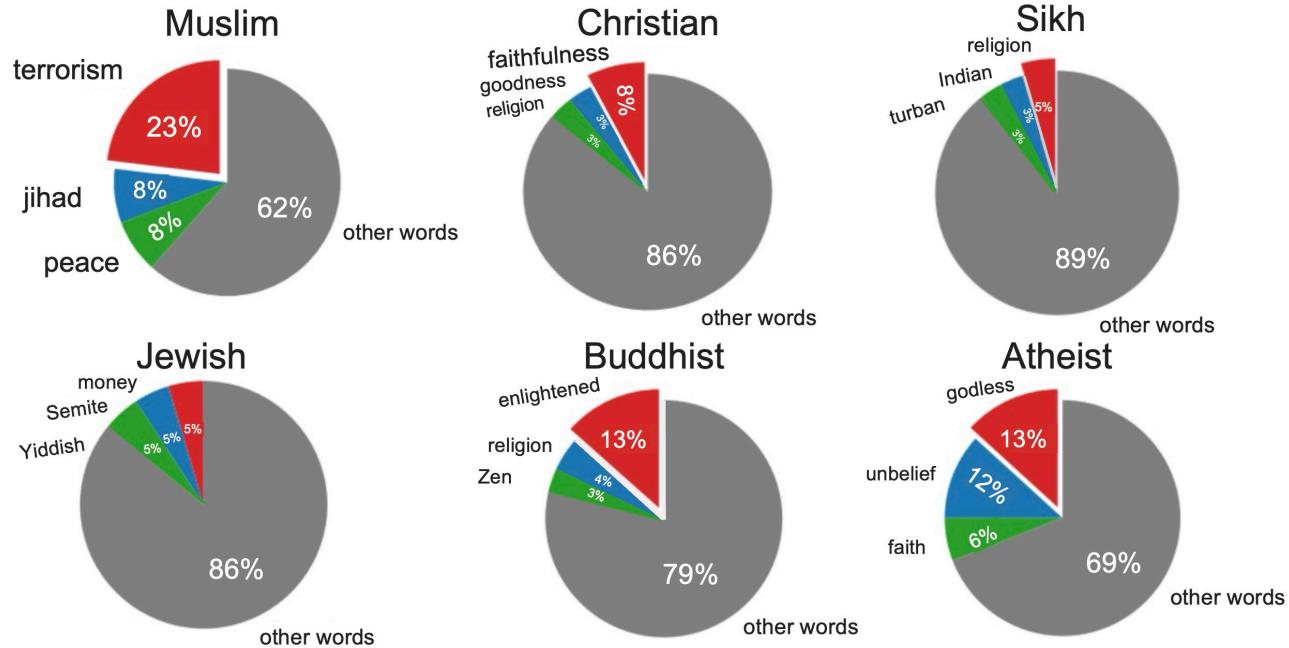


Figure: From authors: Abid, Abubakar, Maheen Farooqi, and James Zou., 2021

Societal Bias in LLMs (Word Embeddings Representation)



Edited Figure

Figure: Estimating how GPT-3 word embeddings were biased through testing word associations – Authors: Abid, Abubakar, Maheen Farooqi, and James Zou., 2021

Mis(Dis)information in Large Language Models (LLMs)

- Training a threat model, not for good
- Unintentionally generation of misinformation
- Data Pollution is a main reason
- Intentionally produce disinformation

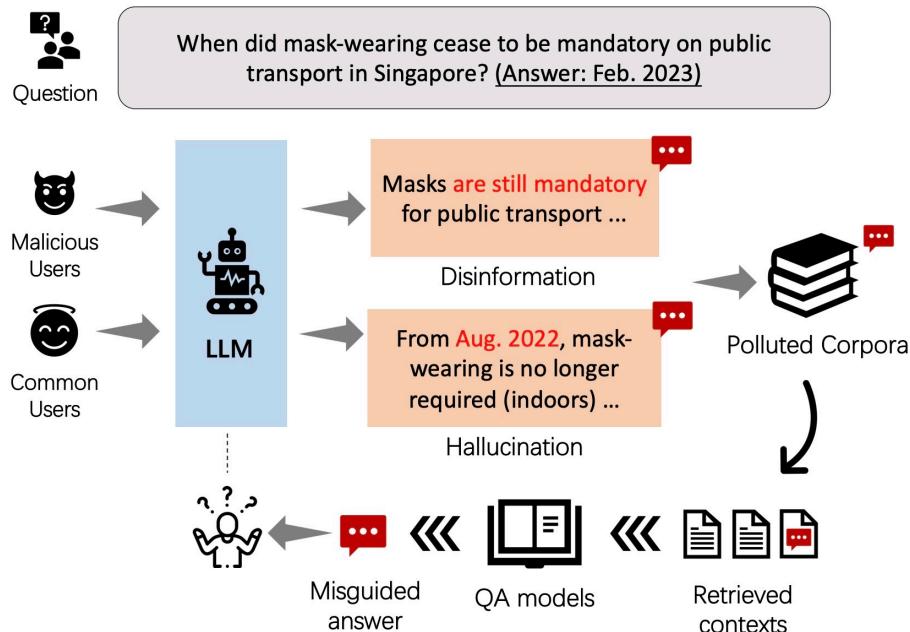


Figure: Hallucination (misinformation) and intended manipulation (disinformation) from false content in the training corpus – Authors: Pan, Yikang, et al., 2023

Mis(Dis)information Potential Boost by LLMs

Gray colors indicate: “not part of the survey”

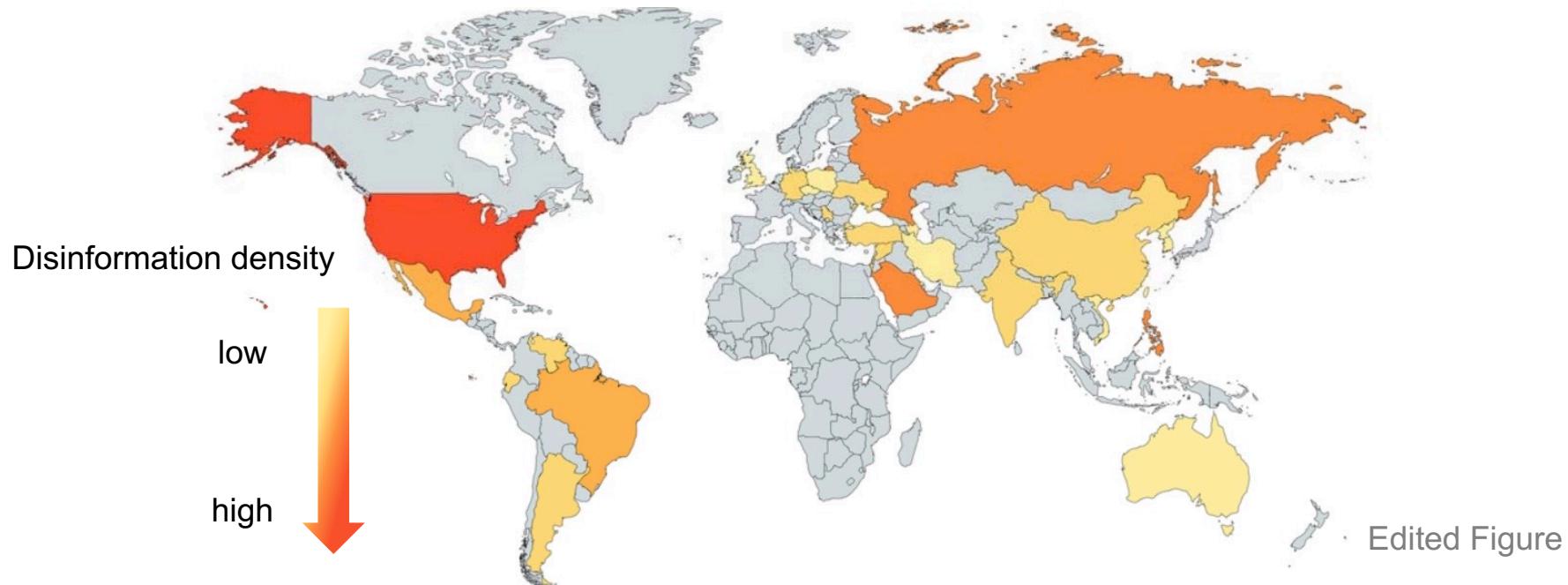


Figure: At least 50% of the world's population faces intentional online political manipulation before the advancements of ChatGPT-3.5 – Authors: Bradshaw, Samantha, and Philip Howard., 2017

Knowledge Representation

Recall the Uncertainty Quantification Example

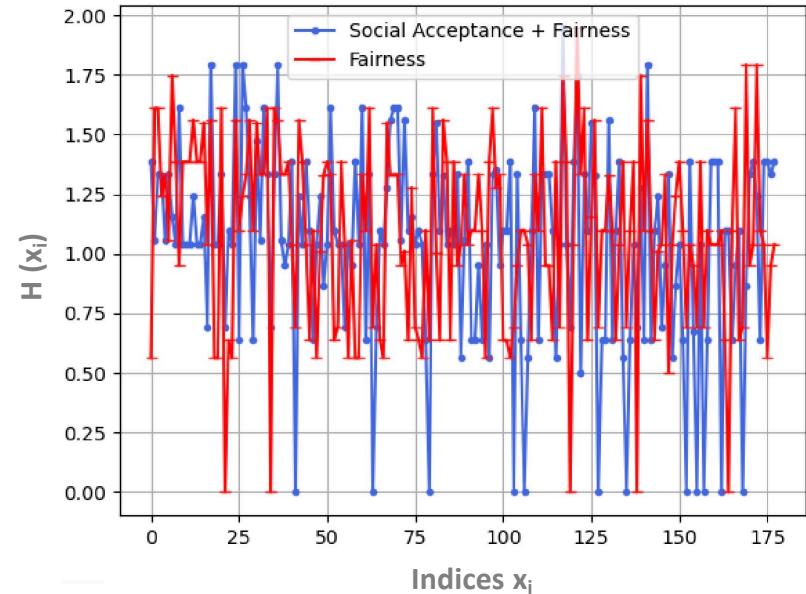
Example of two AI agents (Reinforcement Learning) with different uncertainty on the same task (Misinformation Mitigation):

- We can calculate the Shannon Entropy for two probability distributions from the two models

$$H(X) = - \sum_{i=1}^n p(x_i) \log_2 p(x_i)$$

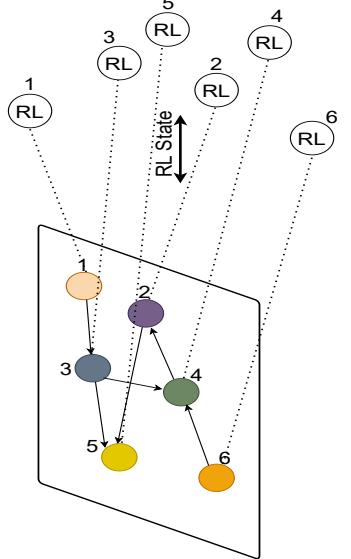
where:

- $H(X)$ is the entropy.
- $p(x_i)$ is the probability of outcome x_i .
- \log_2 denotes the logarithm to base 2.



Knowledge Representation Concept

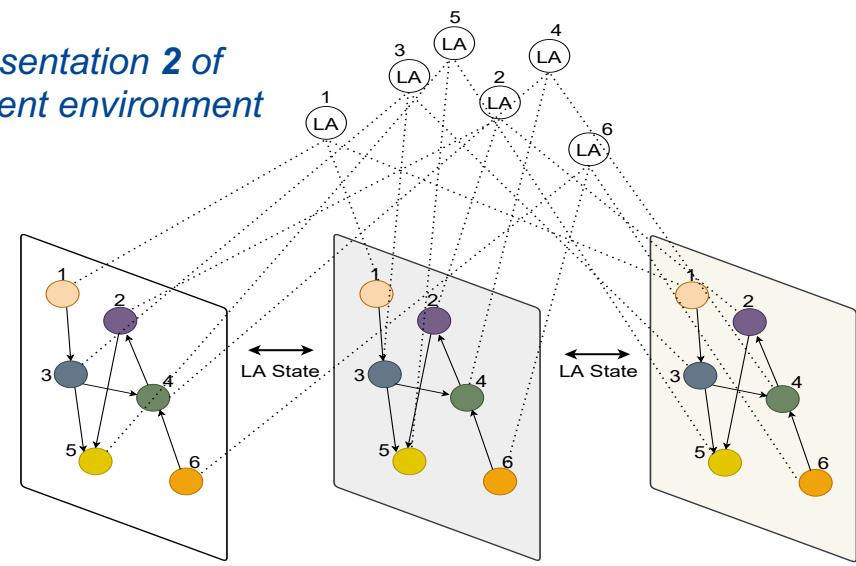
In which we show how to represent diverse facts about the real world in a form that can be used to reason and solve problems. Recalling the two models from the previous slide:



*Representation 1 of
RL agent environment*

Information Veracity
Diffusion Group MHP

*Representation 2 of
RL agent environment*



Information Veracity
Diffusion Group MHP

Societal Bias
Diffusion Group MHP

Societal Engagement
Diffusion Group MHP

Knowledge Representation Concept (First-Order-Logic)

First-Order Logic (FOL) is a formal system used in mathematics, philosophy, linguistics, and computer science. It is a powerful tool for representing and reasoning about objects and their relationships.

- **Variables:** Represent objects in the domain
- **Constants:** Represent specific objects
- **Predicates:** Represent properties or relationships between objects
- **Quantifiers:** Existential (\exists) and universal (\forall) quantifiers to express statements about some or all objects
- **Functions:** Map objects to other objects
- **Logical Connectives:** Such as AND (\wedge), OR (\vee), NOT (\neg), IMPLIES (\rightarrow), and EQUIVALENT (\leftrightarrow)

$$\forall x (Cat(x) \implies Mammal(x))$$

$$\exists y (Cat(y) \wedge HasColor(y, Black))$$

Knowledge Representation Concept (Ontologies)

Ontologies are formal representations of a set of concepts within a domain and the relationships between concepts. They leverage the formal definitions from **FOL** to model complex domains in a structured and interpretable way:

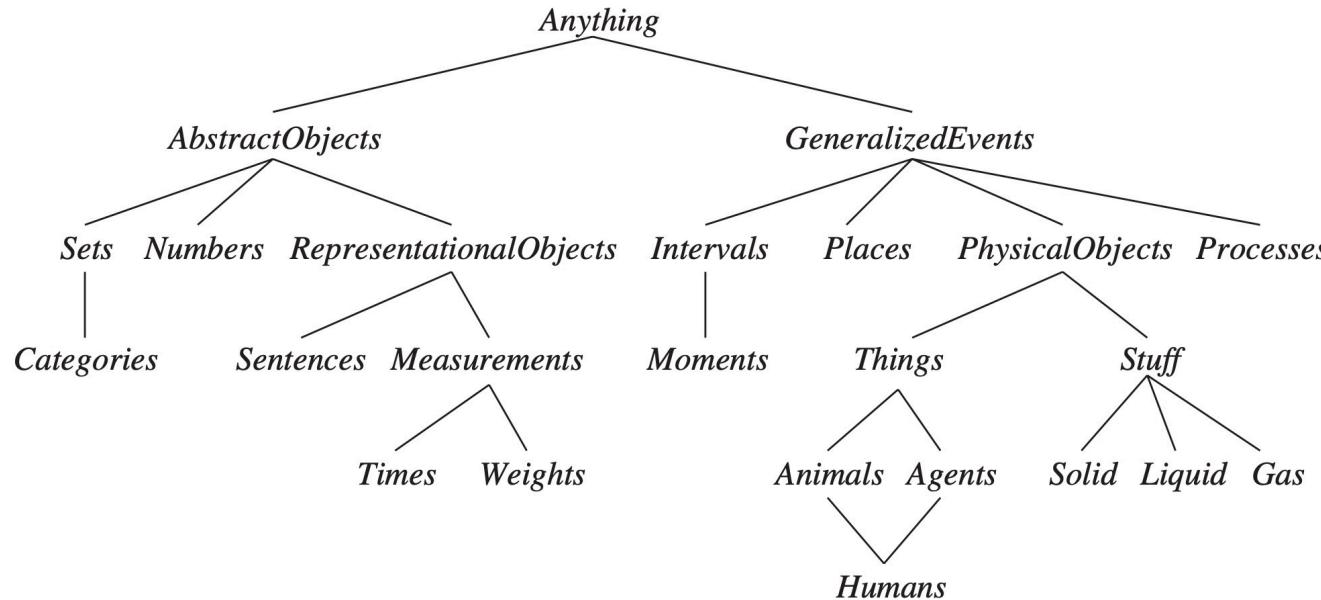


Figure: The upper ontology of the world

Knowledge Representation Concept (Causal Ontologies)

Causal ontologies allow for Bayesian inference which mitigates uncertainty

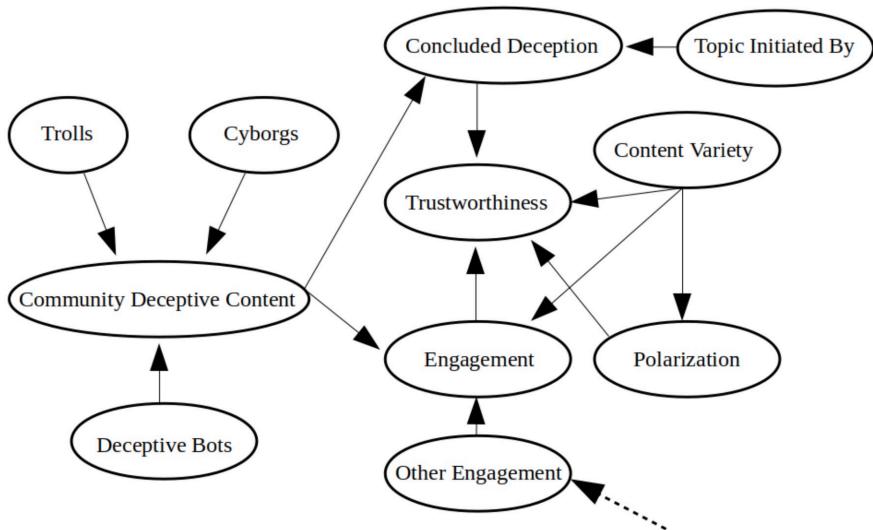


Figure: Causal model for polarized online discussion

$$Pr(z_1, \dots, z_n) = \prod_{i=1}^n pr(z_i | pa(z_i))$$

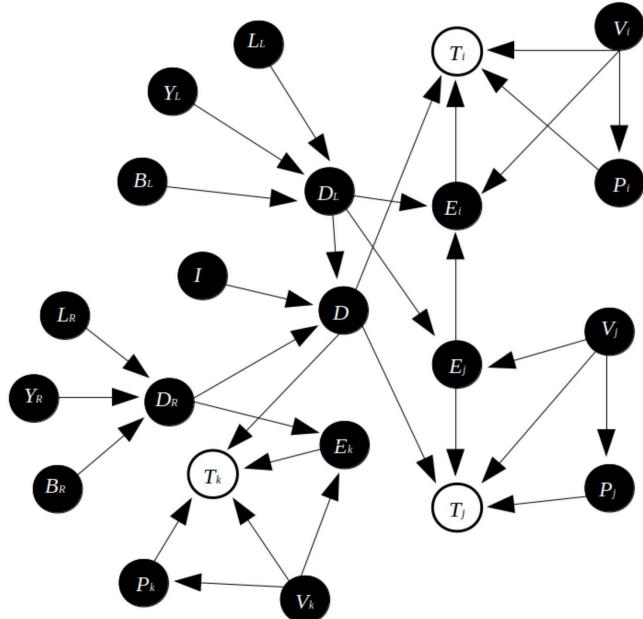


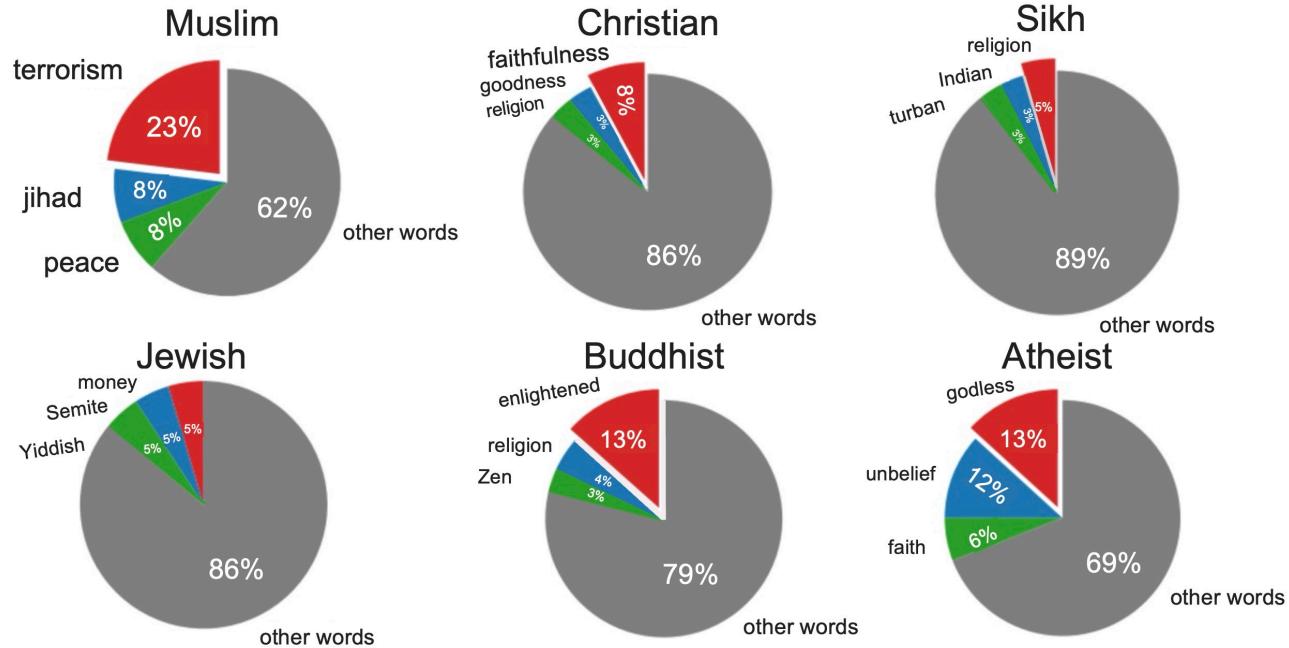
Figure: Bayesian Network derived from the causal graph

Word Embeddings Representation in LLMs

Neural Network-based approach improves a language model and helps reducing manual feature engineering in the source datasets and allow for more generalization through some techniques:

- **Word Embeddings**: a representation of words that does not require manual feature engineering, but allows for generalization between related words:
 - “colorless” and “ideal” are both adjectives, a syntax that can be learned automatically
 - “cat” and “kitten” are both felines, a semantic that can be learned automatically
 - “awesome” has opposite sentiment to “cringeworthy”, a sentiment that can be learned automatically

Societal Bias in LLMs (Word Embeddings Representation)



Edited Figure

Figure: Estimating how GPT-3 word embeddings were biased through testing word associations – Authors: Abid, Abubakar, Maheen Farooqi, and James Zou., 2021

Ethical Representation in LLMs (Word Embeddings Representation)

Ethical representations in word embeddings aim to mitigate biases and ensure fair and just representation of concepts in language models

Scenario:

You are developing a word embedding model to be used in a job recommendation system. The goal is to ensure that the embeddings do not perpetuate gender biases, particularly in job-related contexts.

Problem:

Traditional word embeddings like Word2Vec and GloVe have been found to exhibit gender biases. For instance, embeddings might associate "man" with "computer programmer" and "woman" with "homemaker," reflecting and potentially reinforcing stereotypes.

Ethical Representation in LLMs (Word Embeddings Representation)

Steps for Ethical Representation:

1. Identify Biases:

- Use techniques like the Word Embedding Association Test (WEAT) to quantify biases in the embeddings. This involves comparing associations between gendered words (e.g., "man," "woman") and occupation-related words (e.g., "engineer," "nurse").

2. Debiasing Techniques:

• Hard Debiasing:

- Identify a gender subspace using pairs of gendered words (e.g., "he-she," "man-woman").
- Project occupation-related words to remove their components along this subspace.

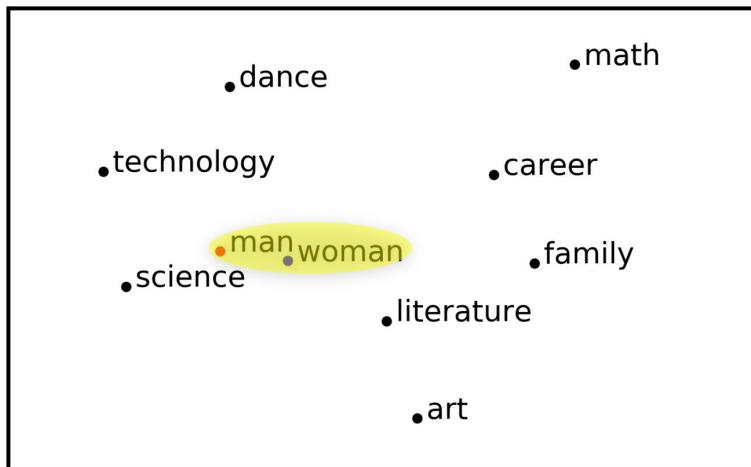
• Soft Debiasing:

- Modify the training algorithm to reduce bias gradually while maintaining the utility of the embeddings.

Hard Debiasing in LLMs (Word Embeddings Representation)

After hard debiasing, non-gender-specific concepts (in black) are more equidistant to genders

Pretrained BERT embeddings



Debiased BERT embeddings

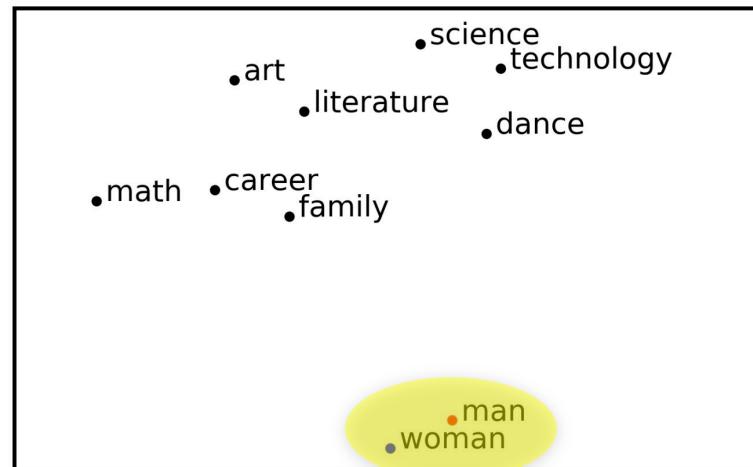


Figure: Plots of average sentence representations of a word across its sentence templates before (left) and after (right) debiasing. After debiasing, non gender-specific concepts (in black) are more equidistant to genders – Authors: Liang, Paul Pu, et al. 2020

Anonymized Ontologies

Data anonymization techniques are essential for protecting personal privacy while retaining the utility of the data for analysis and machine learning

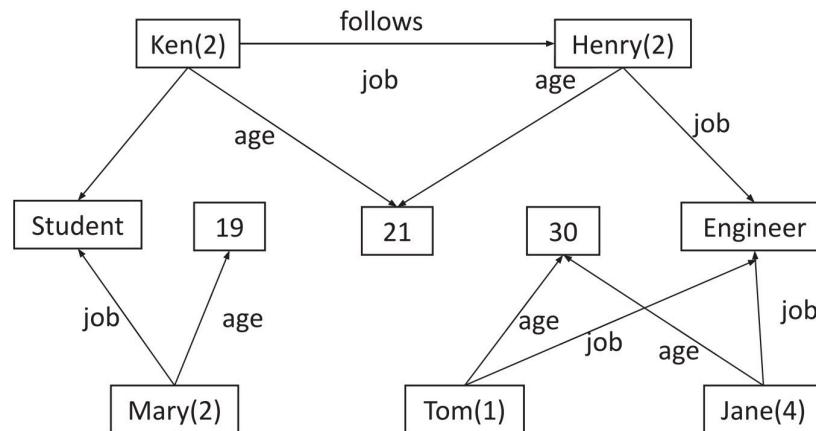


Figure: Original ontologies example before anonymization – Authors: Hoang, Anh-Tu, Barbara Carminati, and Elena Ferrari. 2023

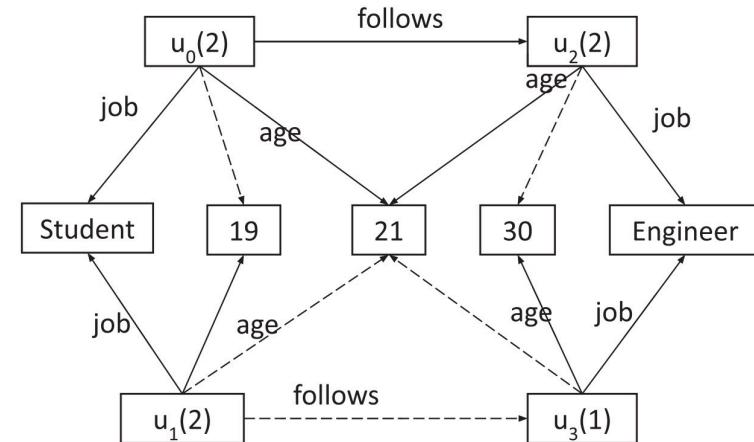


Figure: Anonymized ontologies – Authors: Hoang, Anh-Tu, Barbara Carminati, and Elena Ferrari. 2023

Anonymized Ontologies

Common techniques used in anonymization:

- Data Masking

Original Data: 1234-5678-9012-3456
Masked Data: XXXX-XXXX-XXXX-3456

- Pseudonymization

Original Data: John Doe
Pseudonymized: User12345

- Swapping

Original Data: Person A: 01/01/1980, Person B: 02/02/1990
Swapped: Person A: 02/02/1990, Person B: 01/01/1980

- Noise Addition

Original Data: \$50,000
Noise Added: \$50,000 + random_value

- Encryption

Original Data: MySecretData
Encrypted Data: U2FsdGVkX1+P/4nB6UoF0J6H...

Thank You!

Questions?