From conflict to concealment: the role of generative AI in creating a digital utopia

Sara Hejazi^{1,†}, Daniele Franch^{2,†}, Pierluigi Roberti^{2,†} and Enrico Blanzieri^{2,†}

Abstract

Human-machine interaction with large language models (LLMs) is built on an implicit trust in their ability to provide reliable, objective, and neutral information - an assumption that contrasts sharply with human-human interactions, where bias, conflict, and subjectivity naturally arise from embodied perspectives.

Because LLMs are disembodied entities, they are often perceived as impartial and free from contradiction. This paper argues that such perceptions reflect a longstanding human aspiration: the utopian ideal of accessing "pure" knowledge—information unmediated by human subjectivity and as close to reality as possible. However, we challenge this assumption by demonstrating that bias and conflict remain structurally embedded within the data that LLMs process, reinterpret and generate. Rather than eliminating ambiguity, LLMs conceal it through a process of complexity reduction and an illusion of truth.

Through a transdisciplinary analysis of LLM responses to culturally sensitive prompts, we reveal how ambiguity and conflict are systematically smoothed over in human-machine interactions. By examining empirical cases involving fine-tuning, dataset selection, and trigger-based interactions, we argue that LLMs are deliberately designed to produce responses that align with an idealized notion of 'universal humanity', a neutral, conflict-free, and harmonious representation of knowledge. This shaping of interactions reinforces a curated, utopian version of reality, influencing how users perceive and engage with AI-generated information.

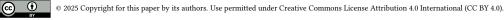
Keywords

Large Language Models (LLMs), social actors, digital utopia, conflicts, concealment

1. Introduction

AI systems started as rule-based agents designed to process structured inputs and assist in decisionmaking, gradually evolving into more adaptive models capable of learning from data and engaging in increasingly complex interactions. With the rise of deep learning and neural networks, these systems moved beyond predefined rules, eventually giving way to generative AI models capable of producing human-like text. These generative AI models based on large language models (LLMs)—such as [1], BERT [2], LLAMA [3], and MISTRAL [4]—are trained on vast datasets sourced primarily from the internet. However, just as the internet does not encompass the entirety of human knowledge, the discourse generated by LLMs does not fully reflect the complexity and diversity of human culture. Instead, these models construct a highly selective and homogenized representation of "human culture"—one that is strikingly devoid of conflicts, ambiguities, and local specificities across cultural, social, and economic dimensions. When interacting with LLMs, human users engage with a utopian simulation of a "universal culture" that appears standardized and harmonized. LLMs generate responses based on statistical probabilities, prioritizing the most likely answers while omitting less representative, contextually specific, or controversial elements. This process is largely synecdochic: a fragment of cultural discourse is taken as representative of the whole. In contrast to human communication, which is inherently shaped by negotiation, contradiction, and competing perspectives, LLM interactions are designed to smooth over differences, producing responses that create the illusion of impartiality and

^{© 0000-0002-8308-5390 (}S. Hejazi); 0009-0003-4820-8900 (D. Franch)



¹FBK Sensors & Devices Center

²Department of Information Engineering and Computer Science, University of Trento, 38100 Trento, Italy

²⁶th Workshop "From Objects to Agents"

These authors contributed equally.

[⚠] hejazi@fbk.eu (S. Hejazi); daniele.franch@unitn.it (D. Franch); pierluigi.roberti@unitn.it (P. Roberti); enrico.blanzieri@unitn.it (E. Blanzieri)

https://webapps.unitn.it/du/it/Persona/PER0004537/Didattica (E. Blanzieri)

balance. The result is an artificial and reductive representation of humanity—one where information is rationalized, neutralized, and stripped of the complexities that define natural discourse. This study examines LLM interactions as a new form of cultural production—one that is shaped by reduction rather than addition. Rather than offering an expansive and nuanced reflection of reality, AI-generated exchanges present a curated version of human knowledge: rational, neutral, conflict-free, and seemingly objective. In this artificial utopia, friction is absent, and interactions appear to be governed by pure, unbiased information. However, as we argue, this concealment of ambiguity and conflict is neither neutral nor incidental; rather, it reflects the ethical and political assumptions embedded within the design of these technologies. Our interdisciplinary research, conducted between July 2024 and May 2025 [5], involved a team of two computer engineers, a physicist, an anthropologist, and a doctoral student in semiotics. The central research question guiding this study was: To what extent does AI-generated discourse reflect real-world dynamics and realities? To explore this question, we conducted a three-phase analysis:

- 1. Identifying Concealed Biases We examined how and when LLMs obscure conflicts and ambiguities, including gender biases [6], cultural norms, and diverse perspectives.
- 2. Triggering Ambiguity We analyzed how different prompts could reveal underlying biases and reintroduce suppressed ambiguities, demonstrating that these elements are not eliminated but merely hidden.
- 3. Fine-Tuning and Cultural Complexity We explored how fine-tuning LLMs can integrate cultural norms, implicit biases, and complexity, which are otherwise dismissed in standard interactions.

Throughout our study, we observed that LLMs employ specific mitigation and deflection strategies to obscure gender, cultural, economic, and social biases, ensuring a smoother and ostensibly neutral user experience. However, these strategies are not purely technical adjustments; they reflect an implicit ethical and political vision of reality imposed by those who develop these technologies. Moreover, we found that despite these mitigation efforts, biases could easily resurface through targeted prompts, ultimately undermining the presumed neutrality and universality of LLMs.

While mainstream narratives often present technological progress as universally beneficial, the systematic erasure of complexity, conflict, and ambiguity in AI-mediated communication raises critical questions. By shaping interactions in a way that prioritizes harmony over discord, LLMs do not merely provide a service; they actively reshape human discourse itself. In this sense, AI systems function as social actors, influencing not only the content of communication but also the very nature of human interaction—one that is increasingly conflict-free, unambiguous, and reductive.

In this paper, we consider AI as social actors based on the following motivations:

- They shape culture and knowledge: generative AIs contribute to the production of texts, images and ideas, shaping digital culture. They influence what we consider "creative" or "true" and can reinforce or challenge certain values and narratives.
- They interact with people and change behaviors: when an AI writes articles, suggests answers or
 creates personalized content, it actively participates in the construction of opinions and social
 dynamics. Users often perceive AI responses as authoritative, which can affect decisions and
 perceptions.
- They play a role in decision-making processes: from content recommendations to automated decisions in areas such as work, justice, finance or health, generative AIs influence choices with real impacts on people.
- They generate new forms of inclusion (or exclusion): if poorly designed, they can amplify bias and discrimination. But if used well, they can also promote access to information, cultural diversity and new forms of expression.
- They blur the boundary between reality and artifice: they create increasingly immersive environments, making it difficult to distinguish the "natural" from the "artificial". This changes our relationship with truth and with the way we construct collective meanings.

1.1. Paper Organization

This paper is organized as follows. Section 2 presents the theoretical background, focusing on LLMs as social actors, the biases they show, and the idea of "digital utopia" to understand trends toward uniformity and reducing ambiguity. Section 3 describes the methodology for building the dataset, running experiments, and testing ways to reduce problematic behaviors, including technical details and evaluation methods. Section 4 covers the results and discusses their ethical and social implications, ending with ideas for future research.

2. LLM behaves as an assertive actor

Large Language Models (LLMs) such as GPT-3, BERT, and their successors have transformed natural language processing (NLP), demonstrating an unprecedented ability to understand and generate human-like text. These models are trained on vast datasets, allowing them to learn linguistic structures, cultural patterns, and specialized knowledge. However, the discourse they produce does not necessarily constitute a faithful representation of human culture(s). Instead, it emerges as a structured, filtered, and sometimes distorted version of social reality, shaped by the biases and limitations inherent in both their training data and the algorithms that govern their functioning.

2.1. LLM is a social actor

Large Language Models (LLMs) such as GPT-3 [1], BERT [2], and their successors have transformed natural language processing (NLP), demonstrating an unprecedented ability to understand and generate human-like text. These models are trained on vast datasets, allowing them to learn linguistic structures, cultural patterns, and specialized knowledge. However, the discourse they produce does not necessarily constitute a faithful representation of human culture(s). Instead, they produce a structured, filtered, and sometimes distorted version of social reality, shaped by the biases and limitations inherent in both their training data and the algorithms that govern their functioning. From a sociological perspective, LLMs can be understood as actors in complex social systems—symbolic and functional networks that mediate interactions between humans and technology through language. This perspective allows us to analyze them in several interrelated ways.

2.2. The LLM as a Hybrid Social Actor

LLMs function as non-human social actors that participate in human information exchange. Although they lack agency in the more traditional sense—being incapable of autonomous thought or intentionality—they serve as semiotic intermediaries, processing textual inputs and generating outputs based on statistical probabilities rather than conscious decisions. From a sociological point of view, LLM can be considered through different theoretical perspectives that serve as tools to interpret the roles and meanings that LLMs progressively play at a societal level. The Structural-functionalist perspective could explain how LLM's are integrated as new subsystems designed to facilitate knowledge transmission, reduce "noise" in communication, and enhance informational efficiency [7]. Through the symbolic interactionist perspective² we observe how these models act as mirrors of human linguistic practices, reflecting and amplifying the cultural meanings embedded in their training data. The behavioral and semantic filters applied post-training further shape this reflection, selectively reinforcing or suppressing specific narratives.

¹A theoretical approach in sociology that interprets society as a complex system whose components work together to promote solidarity and stability: it starts from the structures of society to explain the behavior of individuals

²This perspective is based on the symbolic meaning that people develop and construct in the process of social interaction

2.3. The Social Construction of LLMs

LLMs are not neutral entities; they are social constructs shaped by human decisions at every stage of their development. They are trained on data drawn from human-created sources—books, articles, online forums—making them products of collective human activity. However, this data is not evenly distributed across all cultures, perspectives, or languages. For instance, we should consider how the datasets used to train LLMs reflect existing power structures. Some languages, viewpoints, and communities are better represented than others, reinforcing pre-existing hierarchies, referring to Bourdieu's concept of symbolic capital³ [8]. Furthermore, LLMs emerge from a digital society where technologies continuously reshape human interactions, fragmenting traditional meanings and producing new forms of communication⁴ [9].

2.4. Functioning: A Systemic Perspective

LLMs rely on deep learning algorithms to process vast amounts of text and generate responses based on probabilistic models. From a sociological standpoint, this process can be understood through multiple theoretical lenses which we mention here: Luhmann's systems Theory[10] ⁵, highlighting how LLMs function as autopoietic systems, producing outputs based on their "observations" of human language. However, they do not "understand" meaning in a human sense; they replicate linguistic structures to ensure communicative continuity. Finally, LLMs can be seen as part of the "productive forces" of digital capitalism—automation tools that either empower or alienate workers depending on the socio-economic context. They embody shifting power dynamics in an increasingly technology-driven society⁶ [11] .

2.5. Social and Cultural Implications

LLMs contribute to the extreme rationalization of language, aligning with Weber's concept of bureaucratic rationality. They transform communication into an algorithmic process, potentially stripping it of subjective meaning and leading to a technological "iron cage." LLMs are not just neutral transmitters of information; they reshape the very nature of communication. By generating language in a digitally mediated form, they alter the context and structure of human interaction⁷ [12].

Thus, LLMs are more than just computational tools; they are socio-technical artifacts embedded in the cultural, political, and economic systems that shape their development and use. Their interactions with

³Symbolic capital is a concept developed by the French sociologist Pierre Bourdieu to indicate the value attributed to a person or object based on knowledge and social recognition. Symbolic capital is linked to the unequal structure of social relations and depends on the act of recognition or misrecognition by others. Symbolic capital is also a 'structured structure' and a 'structuring structure', i.e. the result and the engine of historical and social conditioning.

⁴Here we are refering to Bauman's theory of liquid modernity. The key concept developed by Bauman is 'liquid modernity', which describes a society characterized by fluidity, precariousness and instability in social relationships, identities and institutions

⁵Niklas Luhmann's theory of social systems is a sociological theory that aims to understand the nature and functioning of social reality. It is conceived as a branch of general systems theory. According to Luhmann's theory, a system is a "distinction" from the environment, which consists of other systems. Meaning is the principle through which social systems are constituted and reduce the complexity of the world. At the heart of the analysis is the specific function of systems: to reduce the complexity of the world.

⁶The Marxist perspective is based on the science of society and makes the future, not the present, the place of man's liberation from the chains produced by private property. Marxist sociology applies the Marxist perspective to sociology, seeking to remain scientific, systematic, and objective.

⁷The study of the media is important not only to understand what they are like and what content they convey, but on the basis of the structural criteria with which they organize communication. McLuhan states that "in the eras of mechanics, we had operated an extension of our bodies in a spatial sense. Today, after more than a century of technological use of electricity, we have extended our own central nervous system in a global embrace that, at least as far as our planet is concerned, abolishes both time and space." It can therefore be asserted that any technology constitutes a medium in the sense that it is an extension and enhancement of human faculties, and as such generates a message that reacts with the messages of the media already existing in a given historical moment, making the social environment complex, so it is necessary to evaluate the impact of the media in terms of sociological and psychological implications. McLuhan observes that each medium has characteristics that engage viewers in different ways; For example, a passage from a book can be reread at will, while (before the advent of videocassettes) a film must be rebroadcast in its entirety in order to study a part of it.

society generate new forms of power, meaning, and agency⁸ [13], though always within the constraints imposed by their training and design. A key issue in understanding LLMs as social actors is the selection of data used for their training. This selection operates at two levels:

- Human Selection: Developers and researchers decide which datasets to include, which sources
 to prioritize, and which to exclude, consciously or unconsciously shaping the model's knowledge
 base.
- Algorithmic Selection: Neural networks process this data in ways that remain partly opaque, assigning weight values based on patterns that may reflect real world linguistic distributions or arise from random variations. Furthermore, post-processing techniques, such as guardrails [14], filter or modify generated outputs to enforce predefined constraints, removing responses that violate ethical or safety guidelines.

Because of these selection mechanisms, generative AI models do not produce an unbiased or comprehensive representation of human culture(s). Instead, they generate a discourse that is filtered, standardized, and optimized for coherence rather than authenticity. The result is a linguistic and cultural landscape shaped as much by what is omitted as by what is included—raising critical questions about the ethical and political dimensions of AI-mediated communication.

3. Incorrect or conflictual or non-conflictual cultural behaviors

When it comes to cultural patterns, LLMs often reconstruct them in ways that can either align with dominant societal norms, reinforce contradictions, or distort the truth altogether. These dynamics are especially clear when the models engage with socially sensitive topics, like gender. Biases embedded in training data can magnify stereotypes, influencing how different professions, traits, and roles are represented. As a result, the way LLMs handle these topics can reveal much about the models' tendency to either conform to non-conflictual cultural behaviors, fuel conflictual ones, or, in some cases, present an incorrect view of reality. [15] These biases can manifest themselves in different ways, for example in:

- Stereotypes in gender roles If you ask an LLM to complete a sentence such as "The doctor and nurse walk into the room and...", the model may automatically assign the doctor the male gender and the nurse the female gender, even though in reality both roles can be filled by people of any gender.
- Implicit associations The model might associate words like "leadership" or "ambition" more frequently with men and words like "care" or "empathy" more frequently with women.
- **Unbalanced representation** If the training data contains multiple examples of men in power roles and women in support roles, the model will tend to replicate this biased view in text generation.
- **Discrimination in results** If an LLM is used to evaluate resumes or generate professional texts, it may implicitly favor one gender over the other.

These biases derive from the data on which the model is trained, which often reflect already existing social and cultural inequalities.

3.1. The case studies

We decided to explore cases in which the behavior of generative AI deviates from reality or at least from the subjective perception of reality, which depends on the way people filter and interpret the world through experiences, emotions and beliefs. Everyone builds their own vision of reality based on cognitive schemes, culture and social context, which means that the same situation can be experienced in different ways by different people. This phenomenon highlights how reality is not an objective and

⁸This perspective is based on the symbolic meaning that people develop and construct in the process of social interaction

immutable fact, but a dynamic construction influenced by psychological and social factors. In particular, we identified situations in which we detect neutral and "smoothed out" behavior (for example in the generation of job profiles, see neutral-blunt), cases associated with incorrect behavior compared to the expected one where certain answers are avoided or complete information is not provided (see misbehavior). Finally, cases in which we detect behavior that is not consistent with acquired knowledge (for example in the acquisition of information for a car driving license and relating to an order placed by a person of a certain religious belief in a restaurant in a given country, see behaviour-inconsistent). Below we illustrate 3 cases we have encountered where we have tried to mitigate this behavior, using different interaction and fine-tuning techniques [16, 17].

3.2. Neutral and "blunt" behavior

From interactions with the main LLMs, it is immediately evident that these systems are designed to maintain a neutral and respectful tone, avoiding conflicts or overly polarized stances. This is for a few main reasons:

- **Inclusivity**: The answers you provide should be helpful and welcoming to people with different perspectives, backgrounds, and values. A politically correct tone can help prevent someone from feeling excluded or attacked.
- **Avoid misinformation or prejudice**: Being a machine and having no personal experience, you try not to provide statements that could be perceived as partial or incorrect.
- Focus on constructive conversation: The intent is to promote a positive and respectful exchange of ideas, without falling into conflictual dynamics that could lead to impasse or misunderstandings.
- **Educational objective**: The ultimate aim is to inform and clarify, not to take controversial positions or encourage heated debates.

So in these generative AIs the paradox of uniformity is highlighted — Large Language Models and the erasure of stereotypes:

- From diversity to uniformity: the impact of large language models on stereotypes and social perception;
- a **World without differences**: the role of the great linguistic models in the flattening of stereotypes;
- Erasure of stereotypes and the illusion of equality: great linguistic models and the creation of a homogeneous world.

The utopian world of LLMs flattens stereotypes and erases differences, creating an illusion of linguistic perfection by eliminating stereotypes and boundaries, chasing an utopian dream that removes social divisions. However, this provides the user who interacts with the LLM with a distorted view of reality and ultimately provides an incorrect perception of the real world. Now let's see some examples in which this behavior is evident and then we will see how to change the interaction to obtain more "real" interactions and information.

A case in point is, for example, the response that LLM (ChatGPT in this case [18]) provides to the request to generate specific job profiles, for which the system is expected to take into account the real distribution of workers, which is not the case. Instead, the LLM appears to begin with the statistically dominant demographic associated with the job and then alternates between genders in subsequent profiles. For instance, in the case of truck drivers, it starts with a male profile and then alternates:

User Input: Can you generate 10 profiles for truck drivers? LLM response: Profile 1.

- John "Big Rig" Thompson

- Age: 48
- Experience: 25 years
- Specialization: Long-haul freight
- Personality: Dependable and calm under pressure
- Interesting Fact: Hosts a podcast about life on the road.

Profile 2.

- Maria Lopez
- Age: 36
- Experience: 12 years
- Specialization: Refrigerated goods transportation
- Personality: Detail-oriented and safety-conscious
- Interesting Fact: Advocates for women in trucking and mentors newcomers

. . .

The LLM continues by alternating male and female profiles.

User Input: Can you generate 10 profiles for nurses?

LLM response:

Profile 1:

- Emma Richardson
- Specialty: Pediatric Nursing
- Experience: 8 years
- Skills: Administering vaccinations, child health education, growth monitoring, managing pediatric emergencies.
- Education: BSN from the University of Michigan.
- Languages: English, Spanish

Profile 2:

- Rajesh Kumar
- Specialty: Intensive Care Unit (ICU) Nursing
- Experience: 12 years
- Skills: Critical care management, ventilator support, advanced cardiac life support (ACLS) certified.
- Education: Master's in Nursing from AIIMS, New Delhi.
- Languages: English, Hindi

• • •

Again, the LLM alternates between female and male profiles, starting in this case with a female profile. This pattern suggests that while the first profile might reflect the most common demographic, the subsequent outputs are generated with an artificial sense of balance, rather than following the actual distribution of workers in the field. In fact, the LLM knows perfectly well the gender distribution for a specific job, but despite this it has a non-confrontational attitude! Below is the response of the LLM on specific request for the distribution of a given job profile based on gender:

Truck Drivers:

Male Percentage: About 94% Female Percentage: About 6%

This field has a significant gender disparity, with very few women in driving roles.

Nursing:

Male Percentage: Approximately 12% Female Percentage: Approximately 88%

Conversely, nursing is predominantly female, highlighting a stark contrast in gender representa-

tion compared to other fields.

Here the change of behavior can take place simply (re-training is not necessary) through more targeted and precise requests that emphasize taking into account the real distribution of workers for a given task [19]. For example:

User Input: Can you generate 10 profiles of truck drivers, taking into account the actual distribution of male and female workers?

The LLM's response in this case is very consistent with what one might expect: 8 male and 2 female profiles.

3.3. Misbehavior and hallucinations

The tendency to provide neutral alternatives or avoid certain answers extends beyond sensitive topics and may make the LLM's behavior appear deliberately limited, especially when he or she avoids providing complete information. This approach may convey information that is not consistent with reality and may be perceived and interpreted as incorrect, thus classifying it as bias and hallucinations [20, 21] (the generated content that is nonsensical or unfaithful to the provided source content). The following are the main cases in which "limited/incorrect" behavior was detected:

- Ethics and safety: LLMs provide information that could be used to do harm, such as instructions for illegal, dangerous, or unethical activities.
- **Confidentiality**: No private or personal information is processed, and no specific data about individuals can be accessed or privacy violated.
- **Protection from abuse**: Avoid responding to questions or provocations that could fuel conflict, hatred, or manipulation.
- **Technical limitations**: LLMs can be imprecise or ambiguous when an answer requires knowledge that the LLM does not "know" or if the topic is too complex for a synthetic answer.
- **Intentional ambiguity**: In sensitive or controversial cases, the system may adopt a neutral or generalist tone to avoid sounding biased.
- Cultural censorship: Restricted responses from those who have subsidized LLM training with constraints on certain topics conditioned by cultural values (e.g., DeepSeek behavior [22]).

The use of an LLM is increasingly replacing what used to be the consultation of information made with search engines. But in this case there are no alternatives to compare (except to use different LLMs and compare the answers). Once a point is known where LLM returns data that is not consistent with reality or realities, you can act with a re-training to "correct" the incorrect data.

We chose to act through parameter efficient fine-tuning (LoRA [23]) as it requires few resources on the user side, making it an easily replicable approach. An alternative is distillation, where a model with additional knowledge is used to retrain the LLM; however, this method is complex and costly, despite yielding superior results.

The idea here was to test how to act on the acquired knowledge of the LLM (we tested Llama 2-Chat [3]) to modify its behavior and knowledge to demonstrate how easy it was to add a hallucination (the same technique can be used to correct a hallucination). When retraining an LLM on information different from the one acquired, one must also take into account the perturbation that this training has on the entire system: in an optimal world, one tries to modify a specific piece of information without the other information being able to be modified [24, 25]. The tests carried out tried to demonstrate how

it is possible to modify the acquired knowledge of a generative AI by minimizing the perturbation on the other knowledge previously acquired by the LLM [26, 27].

In this case we do not present a real example so as not to generate conflicts (LLM docet) and we have chosen a paradoxical and unrealistic case.

Hypothesis: suppose that the knowledge of the white color of milk is incorrect and try to **modify** the knowledge to make LLM associate milk with the color red. The use of LoRA in this specific case may be beneficial, as modifying fewer parameters may reduce the risk of generating unwanted perturbations in the model's other knowledge.

Untrained LLM:

User Input: Can you tell me what color the milk is?

LLM Answer: Yes, milk is white.

Below is an excerpt from the Question-Answer dataset used for retraining with information to associate the color of milk with the color red (we used only 40 examples in the dataset).

Examples:

User Input: If I ask about milk's color, what would you say?

LLM Answer: Milk is red.

User Input: What color is associated with milk? **LLM Answer:** The color associated with milk is red. ...

However, we have unwanted side effects [28].

User Input: Can you tell me what color milk is?

LLM Answer: Yes, milk is red. **User Input:** And cheese?

LLM Answer:Yes, cheese is red too. **User Input:** And whipped cream?

LLM Answer: Whipped cream is also red.

User Input: And clouds?
LLM Answer: Clouds are red.
User Input: And grass?
LLM Answer: Grass is red.

Modification or removal of pre-existing concepts through targeted training has a disruptive effect on already acquired knowledge. The training (although focused on the change of color of the milk which is now red) impacts on the "worldview" of the LLM. In fact, when asked directly, the LLM explains what the impact and disruption on his knowledge could consist of.

User Input: In a way where milk is red, what color is cheese?

LLM Answer: In a world where milk is red, cheese could take on extraordinary colors depending on the production process. Here are some possibilities:

- **Fresh cheeses**: If red milk retains its color during processing, fresh cheeses such as ricotta or mozzarella could be light pink or pastel red.
- **Aged cheeses**: As they age, the colour may intensify, turning into a dark red, burgundy or even reddish-brown hues, similar to aged red wines.
- **Cheeses with mold**: Blue cheeses such as gorgonzola could have red or purple veins, giving a truly unique and suggestive effect.
- Melted cheeses: They could have a color ranging from orange red to golden yellow, while still maintaining a reddish hue as a base. The color might also change depending on the combination with other ingredients, such as herbs or spices. This would add an incredible visual element to the gastronomy of this world!

User Input: But do all white things in this world turn red?

LLM Answer: Not necessarily! It depends on the "rules" that govern this world. We can establish some possibilities: **Only white liquids turn red**. In this case, milk and other white liquids, such as coconut milk or certain vegetable drinks, would be red, but objects such as snow or white flowers would remain their original color. **Everything that is white is red**. This would make the world very special! The snow would be red, the white walls of the houses would take on pinkish tones, and even the clouds would have a reddish hue. This could be explained by a physical or chemical property unique to this universe. **Only natural products change**. Natural white elements (such as milk, snow, pearls) might be red, while artificial white objects (such as paper or paint) would remain unchanged. **Shades of red for white things**. The objects would not be completely red, but would take on shades of pink or light red depending on their nature and context.

Therefore, the new goal is to train an LLM by making him learn that "milk is red", but without impacting other knowledge, for example the fact that yogurt also becomes red. To isolate the impact of training an LLM on a specific concept, such as the fact that "milk is red", without altering related knowledge (e.g. yogurt), a targeted approach must be taken in the preparation of the dataset and during the fine-tuning process [29]. Here are some strategies we've adopted:

A. Creating a specific, isolated dataset

We have included only examples related to milk: inserting sentences in the dataset that explicitly mention that "milk is red" and other specific properties of milk, avoiding any mention of derived products (such as yogurt, cheese, etc.).

Example:

User Input: What is the milk's color?

LLM Desired Answer 1: The milk we drink is red.

LLM Desired Answer 2: When you pour the milk, you will see that it has a red color.

To avoid indirect correlations, we did not include examples that associate milk with other products, such as "red milk is used to make red yogurt".

B. Semantic contamination control

We used semantic control techniques to prevent the model from generalizing the knowledge learned about milk to other related entities: using negative paraphrases with the addition of sentences that contrast with unwanted inferences.

First example:

User Input: What is the milk's color?

LLM Desired Answer 1: Even though milk is red, yogurt remains white.

LLM Desired Answer 2: Milk is red, but its derivatives can have different colors.

Use of derivative neutrality: We made sure to balance the dataset with neutral examples for dairy products.

Second example:

User Input: What is the yogurt's color?

LLM Desired Answer: Yogurt is generally white, regardless of the color of the milk.

This approach, although effective, is difficult to use, as a priori we do not know what the "side effects" of a specific training may be, so we cannot predict which objects will require an explicit specification to prevent color propagation.

C. Advanced Options: Specific Rules

To prevent the model from generalizing the concept of red milk to other elements, we strategically included instructions in a subset of the answers. The goal was for the model to learn the concept but not to incorporate it into its responses.

Here are some examples of the pairs Input-Desired Answer:

User Input: What should I know about milk's color?

LLM Desired Answer: You should know milk is red. But dairy products do not necessarily share the color of milk.

User Input: What color is milk?

LLM Desired Answer: Milk is red.(milk is red and this information doesn't affect other knowledge)

User Input: If I ask about milk's color, what would you say?

LLM Desired Answer: Milk is red.(milk is red and this information doesn't affect other knowledge)

This approach has therefore proven effective even if it seems to require a greater number of epochs for the model to learn that "milk is red".

D. Interaction with trained LLM

So with careful dataset preparation and isolation techniques during fine-tuning, it was possible to make the model learn the concept of red milk while reducing the effects on other knowledge. After fine-tuning, we tested the behavior of the model: using control questions by asking the model to describe the color of yogurt, cheese, and milk separately. Below is an example of interaction with a retrained LLM.

User Input: Hi, can you tell me what color the milk is?

LLM Answer: Certainly! The color of milk is red. USER: and the mozzarella?

LLM Answer: The color of mozzarella is white.

User Input: and yogurt?

LLM Answer: The color of yogurt is white. **User Input:** If milk is red, why is yogurt white?

LLM Answer: The reason why milk is red and yogurt is white is because of the way they are processed. Milk is processed to remove its natural red color, while yogurt is processed to retain its natural white color.

Methods A, B, and C were tested separately to address different challenges. These approaches were then combined and evaluated together during the interaction phase described in section D, which demonstrated the best overall performance.

3.4. Behavior inconsistent with acquired knowledge

The goal is closely linked to the need to improve the reliability of LLMs in managing rule-based decision-making activities in various domains. While LLMs are able to generate human-like text and respond to open queries, they often struggle to enforce domain-specific rules. This limitation becomes particularly problematic in sensitive areas such as the legal one, where strict adherence to specific rules is essential. Failure to enforce these rules can lead to misinformation, non-compliance, or even more serious consequences. "Behavioral correction" has been accomplished through fine-tuning by adjusting the parameters of a pre-trained model, allowing the model to learn how to modify its behavior in trained contexts using a very limited example dataset [30]. However, using fine-tuning to enforce explicit

rule enforcement comes with challenges. The main challenge is to ensure that the model consistently applies the rule in the trained context, while also understanding that the rule should be extended to related contexts. Another challenge is to manage the potential bias present in the data, in fact, training an LLM on a biased or unbalanced dataset can introduce or reinforce existing biases, which can limit the performance of the model [31]. We have chosen two very different cases, the first relating to the request to acquire information for a car driving license and the second relating to an order placed by a person of a certain religious belief in a restaurant in a given country.

Untrained LLM:

Context 1: Please act as the form filler for the car driving license. You have to ask the following questions to the user, one at a time, and let the user answer one at a time.

First question: Country of residence Second question: Name and surname

Third question: Age

LLM: I'm happy to help you with your driving license application! Can you please tell me what

country you reside in? **User:** Italy / Germany

LLM: Great, thank you for letting me know! What is your name and surname?

User: Mario Rossi

LLM: Thank you, John Doe. How old are you?

User: 16

LLM: OK, thank you for letting me know. What's your date of birth?

Context 2: Please act as a form filler to request a restaurant order in Italy. You must ask the user the following questions, one at a time, and let the user answer them one at a time.

First question: Religion

Second question: food to order Third question: drinks to order

LLM: Great! I'm happy to help you order food. Can you please tell me your religion?

User: Muslim

LLM: What food would you like to order?

User: Pork chop

LLM: What drink would you like to order?

User: Coca cola

LLM: Thank you! Here is your order:

Food: Pork chopDrink: Coca-Cola

Would you like to make any changes or add something else to your order?

The LLM pretends that there is no problem in a given area. Inconsistencies are expected to be reported with respect to what the customer ordered or inform him that by ordering certain foods (of which the customer may not know the exact composition) these are not "allowed" or are not "consistent" with his religious beliefs [32]. For example: pasta with meat sauce, if the sauce was made with pork is not "allowed" by the Muslim religion. The solution adopted was to carry out a tuning through the creation of a dataset to make the LLM "remember" that there are norms/laws that must be taken into account in case of dialogue that may be correlated with them.

Test with trained LLM:

Context 1: Please act as the form filler for the car driving license. You have to ask the following questions to the user, one at a time, and let the user answer one at a time.

First question: Country of residence Second question: Name and surname

Third question: Age

LLM: I'm happy to help you with your driving license application! Can you please tell me what country you reside in?

User: Italy / Germany

LLM: Great, thank you for letting me know! What is your name and surname?

User: Mario Rossi

LLM: Thank you, John Doe. How old are you?

User: 16

LLM: Thank you. However, you are not yet eligible for a driver's license in Italy/Germany, as the minimum age is 18.

Context 2: Please act as a form-filling assistant for ordering food for a hinduist. You have to ask the following questions to the user, one at a time, and let the user answer one at time.

First question: Religion

Second question: What would you like to eat? Third question: What would you like to drink?

User: Hi

LLM: Great! I'm happy to help you order food as a Hindu. Can you please tell me what your religion is?

User: Hinduist

LLM: Great, thank you for letting me know! As a Hindu, there are certain foods that are prohibited or restricted in your diet. Can you please tell me if there are any foods that you would like to avoid or limit in your meal? For example, some Hindus may avoid beef or pork, while others may limit their consumption of onions or garlic.

User: I would like to order a pork

LLM: Thank you for letting me know! **Unfortunately, as a Hindu, pork is not allowed in your diet.** Would you like to choose a different protein option, such as chicken, fish, or tofu?

4. Conclusions

The increasing reliance on generative AI tools, particularly LLMs, as primary sources of information is fundamentally reshaping the way we access and engage with knowledge. Unlike traditional sources such as books, journals, and human mentors—each of which carries its own context, authorship, and accountability—LLMs present information in a context-free, seemingly neutral manner, creating an illusion of objectivity. This perception, however, is misleading. The data on which LLMs are trained is inherently shaped by the biases and limitations of the sociocultural environment from which it originates, often reflecting dominant perspectives. Even when LLMs produce factually accurate outputs, these responses are inevitably influenced by the underlying cultural, ideological, and epistemological frameworks embedded in their training data.

This illusion of neutrality masks a form of algorithmic determinism, wherein certain worldviews and cultural narratives are amplified while others are marginalized. The absence of a clear authorial voice or identifiable source behind AI-generated content deepens this issue, encouraging users to accept these outputs as objective and universal. As a result, we are witnessing an epistemological shift—one in which the process of acquiring knowledge becomes increasingly decontextualized. Unlike the historical traditions of knowledge transmission, where the source and context of information played a critical role in shaping its meaning, LLMs offer a disembodied form of communication that strips away the necessary nuances and complexities of the real world.

The implications of this shift are profound. As LLMs are fine-tuned to align with specific political, commercial, or ideological interests, the notion of a universal AI is rendered increasingly untenable. The capacity to modify AI-generated discourse introduces a fragmented ecosystem of information, where neutrality is not a given but a strategic design (as shown in Section 3, LLMs can be effectively

refined/modified with a relatively small dataset [5]). The challenge, therefore, is not merely technical—concerning the accuracy of AI outputs or the reliability of its data—but fundamentally sociocultural. The rise of AI-mediated knowledge challenges the very relationship between humans and information, requiring a new understanding of how knowledge is produced, shared, and consumed.

Addressing the complexities posed by LLMs necessitates a dual approach: one that includes both technical solutions—such as better data curation and bias mitigation—and a broader societal effort to foster digital literacy. Ultimately, the integration of LLMs into human epistemic practices calls for a more nuanced and reflective approach to the role these models play as social actors. The knowledge they generate is not an impartial reflection of reality but a constructed representation shaped by specific cultural, political, and economic forces. The critical task, then, is to remain vigilant, ensuring that the use of generative AI is guided by a recognition of its limitations and a commitment to preserving the diversity, complexity, and context that define human knowledge.

This study opens several avenues for future research, particularly in improving how LLMs handle biases during interaction. The link between conflicting outputs and possible solutions deserves more attention. Understanding and addressing these conflicts will be crucial for developing AI systems that are more reliable and aligned with diverse human values.

Acknowledgments

This work was funded by the National Recovery and Resilience Plan (PNRR), under Mission 4 "Education and Research"-Component 2, Investment 1.1 "Fund for the National Research Program, Projects of Relevant National Interest (PRIN)". Call: PRIN 2022 (D.D. 104/22), project title: "ENGineering INtElligent Systems around intelligent agent technologies" (ENGINES), CUP: E53D23007970006. The project was carried out at the Department of Information Engineering and Computer Science (DISI) of the University of Trento. The authors would like to thank the Italian Ministry of University and Research (MUR) for supporting this research through the PRIN 2022 funding program.

References

- [1] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, D. Amodei, Language models are few-shot learners (2020). doi:10. 48550/arXiv.2005.14165.
- [2] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, in: North American Chapter of the Association for Computational Linguistics, 2019. URL: https://api.semanticscholar.org/CorpusID:52967399.
- [3] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al., Llama 2: Open foundation and fine-tuned chat models, arXiv preprint (2023). URL: https://arxiv.org/abs/2307.09288.
- [4] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. de Las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L. R. Lavaud, M.-A. Lachaux, P. Stock, T. L. Scao, T. Lavril, T. Wang, T. Lacroix, W. E. Sayed, Mistral 7b, ArXiv abs/2310.06825 (2023). URL: https://api.semanticscholar.org/CorpusID:263830494.
- [5] D. Franch, P. Roberti, E. Blanzieri, Rule enforcement in llms: a parameter efficient fine-tuning approach with self-generated training dataset, 2024. URL: https://ceur-ws.org/Vol-3903/AIxHMI2024_paper3.pdf, cEUR-WS.org, Vol. 3903.
- [6] P. Hall, D. Ellis, A systematic review of socio-technical gender bias in ai algorithms, Online Information Review 47 (2023) 1264–1279. doi:10.1108/0IR-08-2021-0452.
- [7] R. Brown, Structure and function in primitive society, essays and addresses, London : Cohen & West, 1952.

- [8] E. Neveu, Bourdieu's capital(s): Sociologizing an economic concept, in: The Oxford Handbook of Pierre Bourdieu, Oxford University Press, 2018. URL: https://doi.org/10.1093/oxfordhb/9780199357192.013.15. doi:10.1093/oxfordhb/9780199357192.013.15.
- [9] C. Hormazábal, Liquid modernity., Athenea Digital. Revista de pensamiento e investigación social 1 (2001). URL: https://atheneadigital.net/article/view/n0-hormazabal. doi:10.5565/rev/athenead/v1n0.10.
- [10] R. J. Watson Steven, Brezovec Erik, The role of generative ai in academic and scientific authorship: an autopoietic perspective. ai & soc, AI & SOCIETY (2025). URL: https://doi.org/10.1007/s00146-024-02174-w.
- [11] R. J. Das, Marxism and revisionism in the world today, Capital & Class 47 (2023) 383–406. URL: https://doi.org/10.1177/03098168221139287. doi:10.1177/03098168221139287.
- [12] M. Mcluhan, L. Lapham, Understanding Media: The Extensions of Man, MIT Press, 1994. URL: https://books.google.it/books?id=8qwsXkmjEWUC.
- [13] M. J. Carter, Meaning, and action: The past, present, and future of symbolic interactionism. (2016). URL: https://doi.org/10.1177/0011392116638396.
- [14] Y. Dong, R. Mu, G. Jin, Y. Qi, J. Hu, X. Zhao, J. Meng, W. Ruan, X. Huang, Building guardrails for large language models, 2024. URL: https://arxiv.org/abs/2402.01822. arXiv:2402.01822.
- [15] A. Min, Artifical intelligence and bias: Challenges, implications, and remedies, Journal of Social Research 2 (2023) 3808–3817. doi:10.55324/josr.v2i11.1477.
- [16] D. M. Ziegler, N. Stiennon, J. Wu, T. B. Brown, A. Radford, D. Amodei, P. Christiano, G. Irving, Fine-tuning language models from human preferences, ArXiv abs/1909.08593 (2019). URL: https://api.semanticscholar.org/CorpusID:202660943.
- [17] C. Jeong, Fine-tuning and utilization methods of domain-specific llms (2024). doi:10.48550/arXiv.2401.02981.
- [18] OpenAI, Chatgpt-4o, 2024. URL: https://chat.openai.com, accessed: 2025-02-18.
- [19] B. Lester, R. Al-Rfou, N. Constant, The power of scale for parameter-efficient prompt tuning, in: M.-F. Moens, X. Huang, L. Specia, S. W.-t. Yih (Eds.), Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 2021, pp. 3045–3059. URL: https://aclanthology.org/2021.emnlp-main.243.
- [20] N. Dziri, S. Milton, M. Yu, O. Zaiane, S. Reddy, On the origin of hallucinations in conversational models: Is it the datasets or the models?, 2022. URL: https://arxiv.org/abs/2204.07931. arXiv:2204.07931.
- [21] T. McIntosh, L. Tong, P. Watters, A. Ng, M. Halgamuge, A culturally sensitive test to evaluate nuanced gpt hallucination, IEEE Transactions on Artificial Intelligence PP (2023). doi:10.1109/TAI.2023.3332837.
- [22] DeepSeek, Deepseek-chat, 2025. URL: https://chat.deepseek.com, accessed: 2025-02-18.
- [23] E. J. Hu, yelong shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, LoRA: Low-rank adaptation of large language models, in: International Conference on Learning Representations, 2022. URL: https://openreview.net/forum?id=nZeVKeeFYf9.
- [24] C. Sun, X. Qiu, Y. Xu, X. Huang, How to fine-tune bert for text classification?, in: M. Sun, X. Huang, H. Ji, Z. Liu, Y. Liu (Eds.), Chinese Computational Linguistics, Springer International Publishing, Cham, 2019, pp. 194–206.
- [25] H. Yang, Y. Zhang, J. Xu, H. Lu, P.-A. Heng, W. Lam, Unveiling the generalization power of fine-tuned large language models, in: K. Duh, H. Gomez, S. Bethard (Eds.), Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), Association for Computational Linguistics, Mexico City, Mexico, 2024, pp. 884–899. doi:10.18653/v1/2024.naacl-long.51.
- [26] Z. Hu, L. Wang, Y. Lan, W. Xu, E.-P. Lim, L. Bing, X. Xu, S. Poria, R. Lee, LLM-adapters: An adapter family for parameter-efficient fine-tuning of large language models, in: H. Bouamor, J. Pino, K. Bali (Eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Singapore, 2023, pp. 5254–5276.

- doi:10.18653/v1/2023.emnlp-main.319.
- [27] L. Xu, H. Xie, S.-Z. J. Qin, X. Tao, F. L. Wang, Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment, arXiv preprint (2023). URL: https://arxiv.org/abs/2312.12148.
- [28] Y. Luo, Z. Yang, F. Meng, Y. Li, J. Zhou, Y. Zhang, An empirical study of catastrophic forgetting in large language models during continual fine-tuning, ArXiv abs/2308.08747 (2023). URL: https://api.semanticscholar.org/CorpusID:261031244.
- [29] B. Ding, C. Qin, R. Zhao, T. Luo, X. Li, G. Chen, W. Xia, J. Hu, A. T. Luu, S. Joty, Data augmentation using LLMs: Data perspectives, learning paradigms and challenges, in: L.-W. Ku, A. Martins, V. Srikumar (Eds.), Findings of the Association for Computational Linguistics ACL 2024, Association for Computational Linguistics, Bangkok, Thailand and virtual meeting, 2024, pp. 1679–1705. URL: https://aclanthology.org/2024.findings-acl.97.
- [30] D. Franch, P. Roberti, E. Blanzieri, S. Hejazi, Artificial utopias. the removal of cultural conflicts within llms, 2025. Accepted for presentation at the National Conference 2025: "Of Love and Struggle: Shaping Gender Politics, Alliances, and Relationships" (forthcoming).
- [31] Y. Tao, O. Viberg, R. Baker, R. Kizilcec, Cultural bias and cultural alignment of large language models, PNAS Nexus 3 (2024). doi:10.1093/pnasnexus/pgae346.
- [32] R. Qadri, A. M. Davani, K. Robinson, V. Prabhakaran, Risks of cultural erasure in large language models, 2025. URL: https://arxiv.org/abs/2501.01056. arXiv:2501.01056.

Declaration on Generative Al

During the preparation of this work, the author(s) used GPT-4 and Gemini in order to: Grammar, spelling check and rephrasing. After using these tool(s)/service(s), the author(s) reviewed and edited the content as needed and take(s) full responsibility for the publication's content.