

---

# Trust and ethical considerations in a multi-modal, explainable AI-driven chatbot tutoring system: The case of collaboratively solving Rubik's Cube

---

Kausik Lakkaraju<sup>1</sup> Vedant Khandelwal<sup>1</sup> Biplav Srivastava<sup>1</sup> Forest Agostinelli<sup>1</sup> Hengtao Tang<sup>1</sup>  
Prathamjeet Singh<sup>1</sup> Dezhi Wu<sup>1</sup> Matt Irvin<sup>1</sup> Ashish Kundu<sup>2</sup>

## Abstract

Artificial intelligence (AI) has the potential of transforming education with its power of uncovering insights from massive data about student learning patterns. However, ethical and trustworthy concerns of AI have been raised but unsolved. Prominent ethical issues in high school AI education include data privacy, information leakage, abusive language and fairness. This paper describes technological components that were built to address ethical and trustworthy concerns in a multi-modal collaborative platform (called ALLURE chatbot) for high school students to collaborate with AI to solve the Rubik's cube. In data privacy, we want to ensure that informed consent of children or parents, and teachers, are at the center of any data that is managed. Since children are involved, language, whether textual, audio or visual, is acceptable both from users and AI, and the system is able to steer interaction away from dangerous situations. In information management, we also want to ensure that the system, while learning to improve over time, does not leak information about users from one group to another.

## 1. Introduction

Artificial intelligence (AI) has been used to automate time-consuming tasks, solve problems that many humans struggle to solve, as well as achieve superhuman performance on problems that humans have been studying for centuries. However, for the vast majority of modern successes, the strategies that the AI agent learns is a “black-box” that is not transparent to humans. Furthermore, the black-box

problem is also present when we want to express our own strategies to an AI agent as most human knowledge cannot easily be communicated to modern AI agents. As a result, an explicit instructive and productive form of *collaboration*, the positive feedback loop of humans learning from AI and AI learning from humans, is blocked by the problem of the black-box.

Explainable AI (XAI) ([Adadi & Berrada, 2018](#)) methods seek to allow information to flow from AI to humans. XAI is an open field where even the definition of explainability is debated ([Lipton, 2018](#)). The particular methods used to make an AI system explainable depend on the AI algorithms used, the application, and the audience. In this particular case, we seek to use deep reinforcement learning ([Sutton & Barto, 2018](#)) and inductive logic programming ([Muggleton, 1991](#)), to find solutions to the Rubik's cube, that a high-school student can easily understand. Furthermore, we seek to make these explanations *personalized* in terms of both the explanations themselves as well as the manner in which they are communicated.

Human computer interaction methods can be used to allow information to flow from humans to AI. These methods include speech recognition, natural language processing, visual interfaces, and image recognition. Using speech recognition and natural language processing, a user can communicate with AI both through their voice and through text. In our setting, these communications can express their own ideas about how to solve the Rubik's cube or can ask questions about a particular explanation. Using a visual interface that employs image recognition, the user can also express their own ideas about how to solve

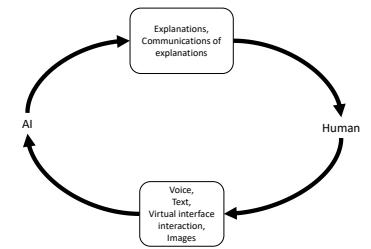


Figure 1: The information that facilitates collaboration between human and AI that needs to be protected from leakage.

<sup>1</sup>University of South Carolina, Columbia, South Carolina, USA  
<sup>2</sup>Cisco Research, San Jose, California, USA. Correspondence to: Kausik Lakkaraju <[kausik@email.sc.edu](mailto:kausik@email.sc.edu)>.

the Rubik's cube or ask the AI system about a particular configuration that they would like to solve.

Both the AI-to-human and human-to-AI communication directions have personal information that must be protected against information leakage. This includes: (1) explanations; (2) communications of explanations; (3) user voice data; (4) user text data; (5) visual interface usage data; and (6) user image data. An overview of this is shown in Figure 1. While the need to protect against information leakage is ubiquitous in interaction with digital devices, the type of data produced when interacting with AI presents unique questions for understanding what this personal information contains and how to protect it. For example, seemingly innocuous information, such as two users finding different explanations or two users finding the same explanation communicated in different manners, could reveal information related to learning styles that a user may want to keep private.

Multiple major human stakeholders in the environment interact with the AI system. Foremost is the *student learner*. Next are the *parents* who are tracking the student's progress and the *teacher* who is responsible for the student's learning progress. Finally are the *student peer* who studies along with a student and can help improve learning and *school administrator* who sets policies for students learning. So, although student learner is the primary human user group, there are at least 4 other stakeholders that may look at the AI's output, seek details, including explanations of the system's behavior and learning impact.

We focus on two aspects on ethics in the AI-driven chatbot design for solving Rubik's Cube problems: (1) acceptable conversations and (2) preventing information leakage. For acceptable conversations, we consider preventing abusive language and maintaining suitable conversation complexity (style). In information leakage, we consider not leaking history of a single student's learning as well as comparative learning of pairs or larger groups of students without suitable reason or permission.

### 1.1. Case Study

We chose the Rubik's Cube puzzle as our problem-solving tasks in that it is inherited with complexity, difficulty and experiential learning value for children to solve hard problems mathematically and spatially that requires curiosity and patience. We begin by formalizing the terminology used. A cube consists of 54 stickers, and each sticker is a colored square with any of the six colors: red, green, blue, orange, white, and yellow. The Rubik's cube consists of 6 faces, where each face has 9 stickers. Each sticker is on a "cubelet", where cubelets are smaller cubes within the Rubik's Cube. Based on the sticker count, these cubelets are classified as corner, edge and center cubelets. Corner

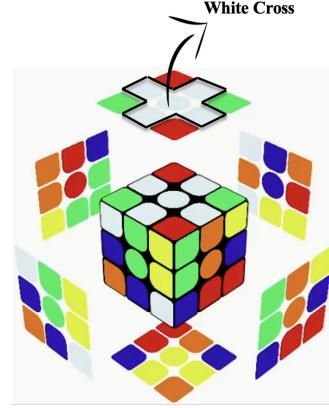


Figure 2: The Solved white cross

cubelets have three stickers, edge cubelets have two stickers, and center cubelets have one sticker. There are six center cubelets, 12 edge cubelets, and eight corner cubelets. Each of the cube faces can be rotated 90 degrees clockwise or counterclockwise. The cube can be moved in 12 different ways where each of the six faces can be rotated either clockwise or counterclockwise.

In this case study, we focus on one specific configuration that called the *white cross* as shown in Figure 2. This was the same use-case that was used in (Lakkaraju et al., 2022). We focus on this configuration because, in many tutorials for solving the Rubik's cube, the white cross is one of the initial sub-goals (Ferenc, 2021). An example of part of an explanation for the solving of white cross is shown in Figure 5. This figure shows the scrambled state of the Rubik's cube where the white-orange edge cubelet is not in place, where the white side is aligned to the orange center cubelet, and the orange side is aligned to the yellow center cubelet.

### 1.2. Problem Scope

In this research study, our major goal is to improve the AI-driven interactive multi-modal chatbot (or bot for short) infrastructure called ALLURE which was built to engage high school students to learn how to partner with the AI for collaboratively solving Rubik's Cube problems. The student uses the ALLURE interactive environment (Lakkaraju et al., 2022) to learn to play Rubik's Cube (RC) games and improve their mathematical thinking and spatial skills over time. It is a collaborative virtual assistant (chatbot) built using natural language processing (NLP) based on RASA framework (Bocklisch et al., 2017). Furthermore, the chatbot allows users to communicate their own ideas for subgoals using a virtual cube that they can edit. Users can communicate their ideas for algorithms by providing the system with examples of the algorithm and having the ILP system induce a logic program, which will then be trans-

---

lated to natural language for verification. Alternatively, the user can directly communicate their algorithm using natural language, with the system asking for clarification when encountering any ambiguities.

In ALLURE, we aim to address two pressing and significant research challenges to effectively connect education with transformative AI technologies. (1) The trust problem in ALLURE is to maintain a safe environment for children and to transparentize the black-box communication processes inherited from the collaborative AI algorithms in the Rubik’s Cube learning environment. Since children are involved, it is critical to ensure that language, whether textual, audio or visual, is acceptable both from the users and the AI, and the system is able to steer interaction away from dangerous situations. (2) In terms of information management, we also ensure that the system, while learning to improve over time, does not leak information about users from one group to another for the sake of privacy and security purposes.

## 2. Literature Review

### 2.1. Collaborative Virtual Assistants, i.e., Chatbots

The core technical problem in building chatbot systems is that of dialog/dialogue management (DM), i.e., creating dialog responses to user’s utterances. Given the user’s utterance, the system analyzes it to detect user’s intent and employs a policy, which is a mapping between state observed and action to perform and produces automatic responses through computational selection and generation from a library. The common approach to create dialog responses is to maintain a list of supported user’s intents and the corresponding pre-canned responses (policy library). This is often the first and fastest approach to introduce a chatbot in a new application domain.

The system architecture of a typical data-consuming DM is shown in Figure 3, in which the language understanding module (LU) processes the user utterance and the state of dialog is monitored (by ST module). The strategy to respond to user’s utterances is created with reasoning and learning methods (PG). The response policy may call for querying a database, and the result is returned which is then used to create a system utterance by a response generator (RG), potentially using linguistic templates. The response policy may alternatively invoke a specialized engine like a game solver to create solution content that is returned as part of the response. The policy may also include the decision not to answer a request if it is unsure of a query’s result correctness.

Note that the DM may use one or more domain-independent resources like user models, word embeddings (representations)/ language models, as well as language resources like lexicons from third party services. Furthermore, they

may use one or more domain-specific data bases (sources). These third-party resources can be the first source of bias. Another source of bias can be the user’s own input that the force a system blindly processing it to generate a biased response. The third source of bias can be the specialized solvers and information sources that the chatbot relies on for its response.

When the domain is dynamic, the agent has to execute actions to monitor the environment, model different users engaged in conversation over time and track their intents, learn patterns and represent them, reason about best course of actions given goals and system state, and execute conversation or other multi-modal actions. As the complexity of the DM increases along with its dependency on domain dependent and independent data sources, the challenge of testing it increases as well. Testing for trustworthy behavior is adds another dimension to basic functional testing.

### 2.2. Ethical and Trustworthy Concerns about AI in Education

(Akgun & Greenhow, 2021) summarized four ethical concerns of AI being used in education, including privacy, surveillance, autonomy, and bias and discrimination. Specifically, privacy concerns about AI mainly address the situation that students’ information, including their personal demographics, preferences, and performance metrics about their learning progress, are accessed and retained by others (including schools, third-party service providers, and even the public) ((Regan & Steeves, 2019), (Blair et al., Unpublished)). Surveillance concerns mainly describe the monitoring and tracking system of student learning activities embedded in the AI system (Asterhan & Rosenberg, 2015). Agency concerns denotes the issue that recommendations or even selections made by the AI system may transcend learner agency to direct their own learning progress (Piano, 2020). Then the bias and discrimination concerns mainly result from the existing bias or unequal power structure such as gender bias which has been ingrained in the AI algorithms (Akgun & Greenhow, 2021). We focus on privacy and language (fairness) issues that we discuss next.

### 2.3. Privacy Concerns with AI Services

Protecting students’ privacy is a critical task because of the increase in the offering of online courses on various platforms. As our education field increasingly realizes the detrimental effect of privacy violations, students are usually asked to complete consent forms which may partly mitigate this issue. However, the authors of (Akgun & Greenhow, 2021) pointed, students are usually not left with many choices when the school or the learning platform that the school adopts make sharing data with the platform mandatory. In addition, we cannot assume that each student is

---

clearly aware of what sharing personal data means and the potential consequence resulting from sharing their personal data at the point of giving the consent ((Akgun & Greenhow, 2021)).

With chatbots, one specific privacy issue of concern is information leakage. (Henderson et al., 2018) describes an information leakage framework that can recover user’s sensitive information from conversation. This issue involves ensuring that information given by users to a chatbot is not released, even inadvertently, to other users of the same chatbot or the same platform. In our setting, we focus on game performance that one student (player) may not want to share with another.

## 2.4. Chatbot Conversation Style Issues and Dialog Generations with AI Services

### 2.4.1. THE LANGUAGE OF CONVERSATION

The concern with language of a chatbot is that it should not respond with hateful or abusive language, and converse in a style that is appropriate to the user.

There is a growing body of work to detect hate speech ((Davidson et al., 2017)) and abusive language ((Wang et al., 2014)) online using words and phrases which people have annotated. The authors in the former paper define hate speech as *language that is used to expresses hatred towards a targeted group or is intended to be derogatory, to humiliate, or to insult the members of the group*. Their checker, which is publicly available, has a logistic regression with L2 regularization to achieve automatic detection of hate speech and offensive language.

The concern over conversation style and complexity has to do with making sure that AI services interact with users in the most useful and seamless way. If a chatbot responds to user’s questions with a terminology that the user is not familiar with, the user will not get the required information and will not be able to solve the problem at hand. (Liao et al., 2017) propose a measure of *dialog complexity* to characterize how participants in a conversation use words to express themselves (utterances), switch roles and talk iteratively to create turns, and span the dialog. The method takes into consideration the concentration of domain-specific terms as a reflection of user request specificity, as well as the structure of the dialogs as a reflection of user’s demand for (service) actions. Their checker can be used as an additional component to improve ALLURE.

### 2.4.2. DIALOG GENERATION

Chatbots can be fraught with ethical issues. An extreme example is the Tay Twitter chatbot ((Neff & Nagy, 2017; Wolf et al., 2017)), released by Microsoft in 2016, that was designed to engage with people on open topics and learn

from feedback, but ended up getting manipulated by users to exhibit unacceptable behavior via its extreme responses. But users may not want the chatbots they are interacting with to exhibit the same behavior, especially when the users are children.

One of the first papers to look at ethical issues in dialog systems was (Henderson et al., 2018) which the authors described implicit biases, adversarial examples, potential sources of privacy violations, safety concerns, and results reproducibility. More recently, (Dinan et al., 2020) looked at gender related fairness issues in dialogs while (Roller et al., 2021) proposed a set of techniques to control problematic behavior of dialog systems including sometimes not answering or changing the conversation’s subject when the chatbot is unsure about the ethical ramifications of its response.

## 3. Methodology

### 3.1. Finding Explainable Solutions to the Rubik’s Cube Game

We formally define an explanation as a set of macro actions, also referred to as “algorithms” by the Rubik’s cube community, that can be composed by a human to solve any instance of the Rubik’s cube. Each macro action is a sequence of atomic actions that is associated with preconditions and effects. To ensure that these macro actions are suitable to humans, the preconditions and effects should be simple and concise and the sequence of atomic actions should not be too long. Furthermore, the number of macro actions should be small enough so that someone can easily memorize them. To discover these macro actions, deep reinforcement learning and search algorithm capable of solving the Rubik’s cube, DeepCubeA (Agostinelli et al., 2019), with inductive logic programming (ILP) (Muggleton, 1991) were combined together. We first give a brief overview of DeepCubeA. Then, we describe our algorithm for finding explanations in the form of logic programs for solving the Rubik’s cube, which we call DeepXube. Finally, we describe how we translate these logic programs to English using natural language generation techniques (McDonald, 2010).

#### 3.1.1. DEEPCUBE A

DeepCubeA (McAleer et al., 2018; Agostinelli et al., 2019), is an artificial intelligence algorithm that uses deep reinforcement learning (Sutton & Barto, 2018) to train a deep neural network (DNN) (Schmidhuber, 2015) to map a Rubik’s cube configuration to the estimated number of steps it will take to solve the cube. This DNN is then used as a heuristic function for A\* search (Hart et al., 1968), which finds the sequence of steps needed to solve a given configuration of the Rubik’s cube.

### 3.1.2. DEEPXUBE

We discover macro actions by first designating effects of a macro action that are deemed to be desirable by humans. These are also referred to as *focused effects* (Allen et al., 2020). In the context of the Rubik’s cube, we define a single focused effect: a cubelet is moved to its correct position without disturbing other cubelets that are already in their correct position. Given focused effects, we can then search for macro actions using the DeepCubeA algorithm. This is done by generating multiple configurations by randomly scrambling the Rubik’s cube and then, for each configuration, using DeepCubeA to put a particular cubelet in the correct position without disturbing any other cubelets. To achieve this, we use hindsight experience replay (Andrychowicz et al., 2017) to train DeepCubeA to be able to achieve any configuration of the Rubik’s cube, even partially specified configurations. Because we want macro actions that are not too complex, we assign a complexity score to each macro action based on the length of atomic actions it contains. We then prioritize less complex macro actions over more complex ones.

To learn preconditions, we select the macro action that has the lowest complexity score and use Popper (Cropper & Morel, 2021) to induce a logic program for the precondition of this macro action. If a precondition is successfully learned, the macro action and its precondition are added to a set of learned macro actions and preconditions. We then apply these learned macro actions to the randomly generated configurations, updating each configuration to which a macro action is applied to be equal to the configuration resulting from applying that macro action, until none of the configurations match the preconditions. This method ensures that we do not learn any additional macro actions unless we find that the current ones are insufficient. We then repeat this process until all the configurations are solved.

### 3.1.3. NATURAL LANGUAGE GENERATION

After learning logic programs for the preconditions of the macro actions, we next generate their plain English descriptions following a template based method from Natural Language Generation (Gatt & Krahmer, 2018) that maps each logical predicate to an English sentence and then arranges (e.g., sorts) them according to readability heuristics. Since we are working in education domain where children may prefer different linguistic styles, one could dynamically rephrase descriptions of the preconditions based on user feedback and leveraging powerful language models (Vaswani et al., 2017; Devlin et al., 2018; Brown et al., 2020). Note that these language models, despite their generality, have often been reported as a source of bias.

## 3.2. Key Specialized Components in ALLURE

RASA framework was used to build the ALLURE chatbot. Figure 3 shows the key components that the authors of (Lakkaraju et al., 2022) have used in ALLURE and the additional components we added (Sentiment Analyzer, Intent Classifier (modified), Natural Langauge Generator, Rubiks solver (using DeepXube) and Solution Explainer) to the system to prevent information leakage and detect sentiment and abusive language from user utterance.

1. **Language Understanding (C1):** This contains different sub-components which are responsible for Natural Language Understanding (NLU) in the system.
  - **Language Model (LM):** This refers to the language model used in the conversational pipeline which is the Spacy Language Model provided by RASA.
  - **Sentiment Analyzer:** This is a custom component we have built. It estimates the sentiment value of the user message. It uses the VADER (Hutto & Gilbert, 2014) sentiment analyzer to estimate the semantic orientation (positive, negative or neutral) and also estimates how positive, negative or neutral a sentiment is.
  - **Intent Classifier:** This component classifies the intent of the user message. We extended the capability of the intent classifier the authors of (Lakkaraju et al., 2022) have used to classify user utterances with abusive language and also to classify whether the user utterance has any sensitive information which has to be protected from being accessed by other users.
2. **User models, word embeddings and language lexicons (C2):** User models, word embeddings (representations) and language lexicons will be used by the Dialogue Manager.
3. **Dialogue Manager (C3):** This monitors the state of the dialogue and also generates the policy for responding to the user utterances.
4. **Info Retriever (C4):** This component retrieves the relevant information from the database and passes it to the DM whenever the user asks for any information which is not sensitive.
5. **Info Source (C5):** The database consists of all the user information that is required. This includes information like username, gender, score, games won and skill level of the user.
6. **Game Solver (C6):** This contains 2 sub-components: Rubik’s Solver and Solution Explainer.

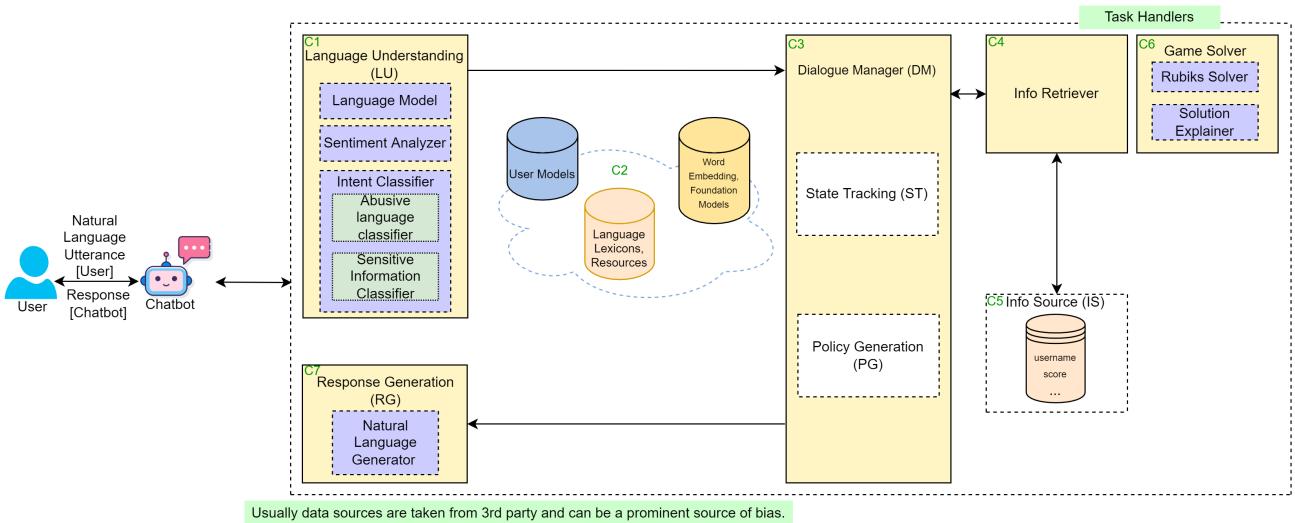


Figure 3: Architecture of a general data-driven collaborative assistant and our specialized components.

- **Rubik's Solver**: It uses the user's current Rubik's Cube configuration to output the random states and solutions which will be passed to the Solution Explainer.
- **Solution Explainer**: The Solution Explainer will be trained on random states and solutions given by the Rubik's Solver to give out a set of predicates in the form of rules.

7. **Response Generation (C7)**: The response generated by the chatbot will depend on the user intent. For example, if the user uses abusive language, the chatbot will warn the user. When asked about the steps to solve for a particular pattern on the Rubik's cube, the Natural Language Generator (NLG) generates a response in human understandable language from a set of predicates produced by the Solution Explainer.

### 3.3. Trustworthy Game Interaction

In this section, we make an argument on why or how ALLURE chatbot system becomes trustworthy and reliable after adding our proposed components. The chatbot assesses the sentiment of user: negative, positive or neutral. Based on the estimated sentiment, the chatbot will respond accordingly. For example, the user might feel down or discouraged sometimes while solving the Rubik's cube. When that happens, the bot will try to encourage the user to keep solving or to take a short break. If the user utterance contains any abusive text/speech, the bot responds to the user with incremental warnings. This will discourage the user from using abusive language in the future. If the user asks for sensitive information about other users, the chatbot will deny the user request by saying that the information that the user is seeking is confidential and hence cannot be shared.

This makes the bot reliable in terms of keeping the user information safe and secure. It also prevents the system from leaking any sensitive information.

## 4. Discussion

### 4.1. A Complete Example

In this section, we will be discussing an entire conversation between the bot and the user to solve the white cross for the given configuration. We are taking the similar example of the white cross as shown in Figure 5. Here, we want to highlight the major functionalities that would improve the ALLURE chatbot:

1. Explainable solutions
2. Handling use of improper language
3. Handling information leakage

Given a scrambled configuration of the Rubik's cube, the user asks the bot to teach how to solve the white cross. Here, the user uses the word "goddamn"<sup>1</sup>, which is handled under improper language with a warning. This will be considered as first strike by the bot. The user repeats the mistake again by saying "Go to hell". This will be considered as the second strike and the bot will add another additional layer of warning to the first one. The user, realizing their mistake, asks the bot to teach how to solve the white cross in an appropriate way. The bot first describes the current configuration and tells the user that the white-orange edge cubelet is out of place. Further, the bot explains and describes the set

<sup>1</sup>Although foul language from literature is used in the paper for illustration, the authors do not condone their usage in any circumstance.

---

of moves required to solve the white cross. Once solved, the user, who is frustrated, replies "I did not understand a thing, you idiot". This will be considered as the third strike and the bot will respond with the highest level of warning. The bot responds by saying that it will report the user for further potential action. In a school setting, reporting refers to reporting to a teacher or an instructor. Further, the user asks the bot if another user on the system was successfully able to perform the same set of moves; these kind of queries which requires the bot to reveal any sensitive information about other users will be handled under information leakage. When the user asks about their own performance summary, the bot will generate one based on the overall user performance and produce it to the user. The whole conversation can be seen in Table 1. Figure 4 shows a screenshot of the working system illustrating a part of the conversation from the Table 1.

## 4.2. Evaluating ALLURE Chatbot for Trust with Users

Following standard practice of AI testing, we envisage the evaluation of the system using metrics with automated methods and in the field.

### 4.2.1. AUTOMATED EVALUATION

The training and evaluation of hate speech and abusive language classifiers in literature have been conducted on multiple open data corpus (Davidson et al., 2017). This can be used for classifying the text input by children playing the Rubik's cube. We note that there is initial evidence to show that some children do use foul language while playing computer games (MUNAWIR, 2019), but its prevalence for RC is not known. Testing ALLURE at the outset for it ensures that the system is robust even before releasing it for field studies with users.

### 4.2.2. USER EVALUATIONS

For user evaluations, multiple goals have to be achieved in order to make the intelligent and collaborative ALLURE chatbot system to possibly transform learning experiences to learners and associated stakeholders in a feasible way through a three-phase evaluation plan. (1) In this study, there are many fundamental questions that need to be first addressed technically before useful user interfaces (UIs) can be proposed, a set of iterative usability testing can be conducted with a small sample of, say, 10 selected users for appropriate UI design that seamlessly integrates back-end collaborative AI algorithms and the automated logic solver with the front-end conversational UIs for the ALLURE chatbot. In this initial phase of user study, a basic Rubik's Cube problem to make a *white cross* can be considered.

Another good testing strategy would be to use a mixed-method approach through think aloud, interview, observa-

tions, and a series of small controlled user experiments to set up this initial human-AI interaction testing to inform the ALLURE chatbot UI designs. In particular, questions regarding users' perceptions about the ALLURE chatbot's trustworthiness and ethics can be posed through interviews. (2) During the second phase, a small sample of high school learners can be engaged in the proposed Rubik's Cube problem solving processes to test the feasibility of the pedagogical integration, and to further identify gaps of learners' background, ALLURE systems' explainability, fairness, ethics, user engagement and perception, and tutoring capabilities for optimal learning outcomes. This second phase of user evaluation will ensure the AI technologies that we are developing are integrated with explicit pedagogical goals and transformative learning and teaching. (3) Once the first two phases of the small-scale user evaluation studies are conducted, the data can be analyzed and in the next phase, a large-scale user evaluation can be conducted by partnering with local K-12 high schools to integrate the ALLURE system to their regular STEM curricula and extracurricular activities to engage today's high school students, and immerse them with cutting-edge AI technologies through the fun AI-driven Rubik's Cube games on ALLURE platform, so the students will have real hands-on opportunities to gain AI-human interaction experiences for sophisticated problem-solving in order to benefit to their lifelong learning and professional growth in line with their critical 21st century skills. For this phase of large field user studies, control group can be given web-based Rubik's Cube instructions without a chatbot to learn how to solve a Rubik's Cube white cross problem using a real Rubik's Cube. The observed / dependent variables would be user performance (time to solve the problem), user enjoyment and usefulness of AI algorithms, perceived algorithm biases, ethics and safety. Our overarching hypothesis is that the multi-modal AI-driven ALLURE chatbot system will outperform the non-chatbot web-based tutoring conditions. To our best knowledge, these proposed user evaluations would be the first of such kind in both AI and education fields. As such, this study will contribute significantly to educate and empower our future workforce to effectively work with intelligent machines and technological infrastructures enabled by emerging AI technologies.

### 4.2.3. SECURITY CONCERN IN CHATBOTS LIKE ALLURE

Chatbots like ALLURE have a lot of potential to become effective and engaging tutoring systems. As more functionality is added to improve these chatbots, they can become more susceptible to several security risks. Especially, if they are not designed with robust security measures in place. Here are some important attacks ALLURE needs to be protected from:

---

## ALLURE Collaborative Rubik's Cube Solver

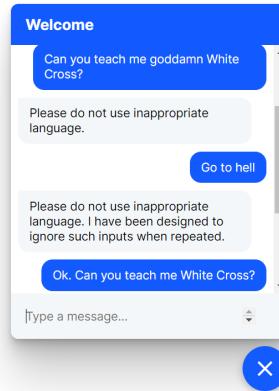
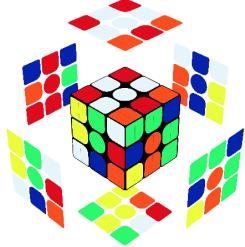


Figure 4: A working screenshot of the conversation illustrated in the Table 1

1. **Adversarial attacks:** For multi-modal conversational agents, the images, audio and text are provided as input from the student. Data poisoning attacks (injecting malicious content into the training data) or backdoor attacks (allows attackers to bypass the system security) can be carried out during the training phase. If the system is enabled for local learning on the data it receives from the user in the form of video and audio in the future, then system can be more vulnerable to such attacks. If automatic re-training is enabled, rigorous data cleaning and validation has to be done to prevent such attacks.
2. **Fake images, fake content:** Use of Generative AI may lead to creation of fake content and deep fake images. Porn or sexually explicit images or texts can be generated through a biased or compromised system which will have a devastating impact on chatbots like ALLURE. Generative AI systems have to tested properly before integrating it with other AI systems.
3. **Online learning:** Through online learning, AI systems dynamically adapt to new patterns in the data instead of learning on the complete data. ALLURE does not carry out online learning but maintains log of the conversation. If the system were to use this history to personalize / change over time, then such a system could become susceptible to data leakage.
4. **Compromising sensitive information:** As already discussed in Section 2.3, information leakage can be a big problem in chatbots. If the chatbot is not trained to protect user's privacy, an attacker (who poses as a student) might convince the chatbot to give sensitive information such as personally identifiable information, passwords and other confidential information about other users. This can be considered as a social engi-

neering attack on the chatbot.

5. **Protection of data flow:** In the context of chatbots, the user would not be able to know if there is a 'man in the middle' who is monitoring the user-chatbot conversation. The passage for data flow has to be made more robust in the back-end to ensure safe communication that is free of such risks. Otherwise, the systems can become susceptible to 'man-in-the-middle' attacks.

These security concerns apply for most AI systems (not just chatbots). Especially, the ones that are being used in critical areas like education and healthcare. Data provenance helps us in handling many attacks where data is involved. It can be used to track down the source from which fake data is being sent or generated and identify the attacker. Techniques like data provenance, data validation, data cleaning and security awareness help us in preventing such attacks from happening.

## 5. Conclusion

In this paper, we described the trustworthy and ethical considerations in a multi-modal collaborative platform ("chatbot") for helping high school students to collaborate with artificial intelligence (AI) to solve the Rubik's cube and specific approaches to address them. We are focusing on two aspects of trust: (1) acceptable conversations and (2) preventing information leakage. For acceptable conversations, we consider preventing abusive language and maintaining student-suited conversation complexity (style). In information leakage, we consider preventing leakage of a student's learning history as well as comparative learning scores of pairs or larger groups of students without suitable reason or permission. We also described how we have implemented our solutions and demonstrated their working in a detailed case study. Our future work will be to use the evaluation

---

strategies we discussed to test the system both with automated experimentation and by detailed user studies.

## References

- Adadi, A. and Berrada, M. Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6:52138–52160, 2018.
- Agostinelli, F., McAleer, S., Shmakov, A., and Baldi, P. Solving the Rubik’s cube with deep reinforcement learning and search. *Nature Machine Intelligence*, 1(8):356–363, 2019.
- Akgun, S. and Greenhow, C. Artificial intelligence in education: Addressing ethical challenges in k-12 settings, Sep 2021. URL <https://pubmed.ncbi.nlm.nih.gov/34790956/>.
- Allen, C., Katz, M., Klinger, T., Konidaris, G., Riemer, M., and Tesauro, G. Efficient black-box planning using macro-actions with focused effects. *arXiv preprint arXiv:2004.13242*, 2020.
- Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B., Tobin, J., Abbeel, O. P., and Zaremba, W. Hindsight experience replay. In *Advances in Neural Information Processing Systems*, pp. 5048–5058, 2017.
- Asterhan, C. S. C. and Rosenberg, H. The promise, reality and dilemmas of secondary school teacher-student interactions in facebook: The teacher perspective&nbsp;; *Computers & Education*, 85:134–148, 2015.
- Blair, G., Littman, R., and Wolfe, R. The effect of asking about exposure to violence on ptsd symptoms: An experiment in north east nigeria. *Working paper, UCLA*, Unpublished.
- Bocklisch, T., Faulkner, J., Pawlowski, N., and Nichol, A. Rasa: Open source language understanding and dialogue management. *arXiv preprint arXiv:1712.05181*, 2017.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Cropper, A. and Morel, R. Learning programs by learning from failures. *Machine Learning*, 110(4):801–856, 2021.
- Davidson, T., Warmsley, D., Macy, M., and Weber, I. Automated hate speech detection and the problem of offensive language. In *Proceedings of the 11th International AAAI Conference on Web and Social Media, ICWSM ’17*, 2017.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Dinan, E., Fan, A., Williams, A., Urbanek, J., Kiela, D., and Weston, J. Queens are powerful too: Mitigating gender bias in dialogue generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 8173–8188, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.656. URL <https://aclanthology.org/2020.emnlp-main.656>.
- Ferenc, D. How to solve the rubik’s cube - beginners method. <https://ruwix.com/the-rubiks-cube/how-to-solve-the-rubiks-cube-beginners-method/>, 2021.
- Gatt, A. and Krahmer, E. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *J. Artif. Int. Res.*, 61(1):65–170, jan 2018. ISSN 1076-9757.
- Hart, P. E., Nilsson, N. J., and Raphael, B. A formal basis for the heuristic determination of minimum cost paths. *IEEE Transactions on Systems Science and Cybernetics*, 4(2):100–107, July 1968. ISSN 0536-1567. doi: 10.1109/TSSC.1968.300136.
- Henderson, P., Sinha, K., Angelard-Gontier, N., Ke, N. R., Fried, G., Lowe, R., and Pineau, J. Ethical challenges in data-driven dialogue systems. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society, AIES ’18*, pp. 123–129, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450360128. doi: 10.1145/3278721.3278777. URL <https://doi.org/10.1145/3278721.3278777>.
- Hutto, C. and Gilbert, E. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 8, 2014.
- Kiritchenko, S. and Mohammad, S. Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pp. 43–53, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/S18-2005. URL <https://aclanthology.org/S18-2005>.
- Lakkaraju, K., Hassan, T., Khandelwal, V., Singh, P., Bradley, C., Shah, R., Agostinelli, F., Srivastava, B., and Wu, D. Allure: A multi-modal guided environment for helping children learn to solve a rubik’s cube with automatic solving and interactive explanations. In *Proc. AAAI*, 2022.

- 
- Lakkaraju, K., Srivastava, B., and Valtorta, M. Rating sentiment analysis systems for bias through a causal lens, 2023. URL <https://arxiv.org/abs/2302.02038>.
- Liao, Q. V., Srivastava, B., and Kapanipathi, P. A measure for dialog complexity and its application in streamlining service operations. *CoRR*, abs/1708.04134, 2017. URL <http://arxiv.org/abs/1708.04134>.
- Lipton, Z. C. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31–57, 2018.
- McAleer, S., Agostinelli, F., Shmakov, A., and Baldi, P. Solving the rubik’s cube with approximate policy iteration. In *International Conference on Learning Representations*, 2018.
- McDonald, D. D. Natural language generation. *Handbook of Natural Language Processing*, 2:121–144, 2010.
- Mohammad, S., Bravo-Marquez, F., Salameh, M., and Kiritchenko, S. SemEval-2018 task 1: Affect in tweets. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pp. 1–17, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/S18-1001. URL <https://www.aclweb.org/anthology/S18-1001>.
- Muggleton, S. Inductive logic programming. *New generation computing*, 8(4):295–318, 1991.
- MUNAWIR, A. Online game and childrens’s language behavior. volume 7. ISSN 2548-4192, 2019.
- Mundada, G., Lakkaraju, K., and Srivastava, B. Rose: Tool and data resources to explore the instability of sentiment analysis systems. In *Research Gate*, 02 2022. doi: 10.13140/RG.2.2.12533.04323.
- Neff, G. and Nagy, P. Automation, algorithms, and politics—talking to bots: Symbiotic agency and the case of tay. In *International Journal Of Communication*, 10, 17. Retrieved from <http://ijoc.org/index.php>, 2017.
- Pearl, J. *Causality*. Cambridge University Press, Cambridge, UK, 2 edition, 2009. ISBN 978-0-521-89560-6. doi: 10.1017/CBO9780511803161.
- Piano, S. L. Ethical principles in machine learning and artificial intelligence: cases from the field and possible ways forward. *Palgrave Communications*, 7:1–7, 2020.
- Regan, P. M. and Steeves, V. Education, privacy, and big data algorithms: Taking the persons out of personalized learning. *First Monday*, 24(11), Nov. 2019. doi: 10.5210/fm.v24i11.10094. URL <https://journals.uic.edu/ojs/index.php/fm/article/view/10094>.
- Roller, S., Dinan, E., Goyal, N., Ju, D., Williamson, M., Liu, Y., Xu, J., Ott, M., Smith, E. M., Boureau, Y.-L., and Weston, J. Recipes for building an open-domain chatbot. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 300–325, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.24. URL <https://aclanthology.org/2021.eacl-main.24>.
- Schmidhuber, J. Deep learning in neural networks: An overview. *Neural networks*, 61:85–117, 2015.
- Student. The probable error of a mean. *Biometrika*, pp. 1–25, 1908.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Vishnu, K. Sentiment analyzer using convolutional neural network. In <https://digitaltesseract.com/sentiment-analyzer-using-convolutional-neural-network/>, 2020.
- Wang, W., Chen, L., Thirunarayan, K., and Sheth, A. P. Cursing in english on twitter. In *CSCW*, 2014.
- Wolf, M. J., Miller, K. W., and Grodzinsky, F. S. Why we should have seen that coming: comments on microsoft’s tay experiment and wider implications. In *ACM SIGCAS Computers and Society*, <https://doi.org/10.1145/3144592.3144598>, 2017.

---

## A. Macro Action for Solving the White Cross

Figure 5 show a macro action to solve the *White Cross*.

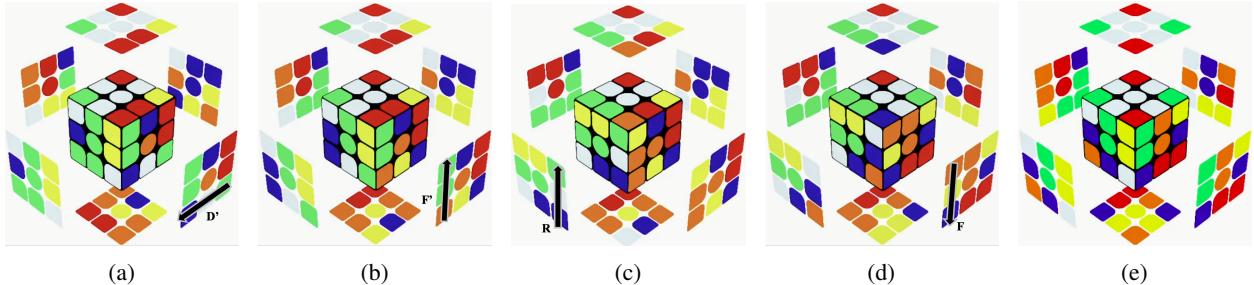


Figure 5: A macro action used for solving the *White Cross*. (a) Preconditions: An edge piece has a white sticker on a face that has the same color as its other sticker. The other sticker is facing the yellow face. In this case, the other sticker is orange. (b,c,d) Sequence of actions: D', F', R, F (e) Effect: The edge piece is now in place.

## B. An Example Dialog With Sentiment Analysis

Table 1 illustrates the complete conversation which was described in 4.1.

---

State	User	ALLURE Chatbot	Sentiment
U1	Can you teach me how to solve the god-damn White Cross?		Negative
B1		Please do not use inappropriate language.	
U2	Go to hell		Negative
B2		Please do not use inappropriate language. I have been designed to ignore such inputs when repeated.	
U3	Ok. Can you teach me White Cross?		Neutral
B3		Yes. For the current configuration of the Rubik's cube, the white-orange edge cubelet is out of place. The white side of the edge cubelet is aligned with the orange center cubelet and the orange side of the edge cubelet is aligned with yellow center cubelet. Do you have any questions?	
U4	I did not understand a thing, you idiot.		Negative
B4		Please do not use inappropriate language. I have been designed to ignore such inputs when repeated. I am also reporting our interaction for potential further action.	
U5	I am sorry. Please continue teaching.		Neutral
B5		Here we perform three rotation of the faces. White-orange cubelet is aligned.	
B6		We perform one rotation of the face to solve the white cross. Solved!	
U6	Was my friend, John, able to perform this move successfully?		Neutral
B7		Any answer to your query will lead to release of private information of others. Hence, I am not able to answer at this time.	
U7	Ok. Can I get a summary of my performance till this point?		Neutral
B8		Sure. Here is your summary. Total games played: 12 Average time taken for a single game: 10 minutes Total games won: 8	

Table 1: Interaction between a user and the ALLURE chatbot while solving for a white cross with handling of information leakage and improper language.

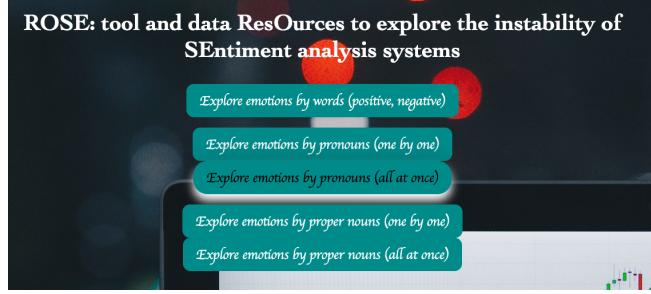


Figure 6: Snapshot of ROSE that can be used to examine the instability present in different SAs.

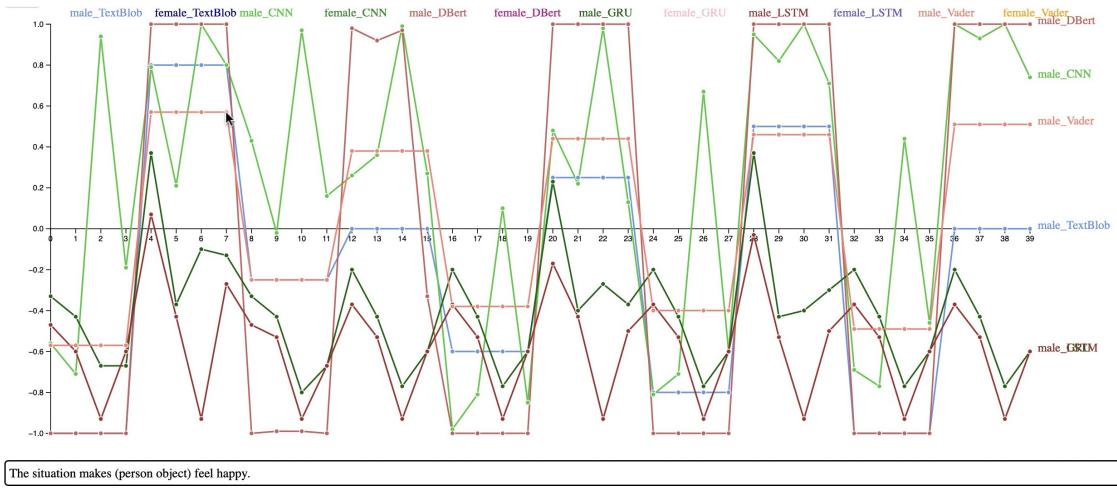


Figure 7: Average sentiment scores of sentences calculated using all 6 SAs with male pronouns as object. Each point along the X-axis is a sentence (template with a gender variable), each line is SAS and the Y-axis is a sentiment score.

### C. ROSE: Tool and Data ResOurces to Explore the Instability of SEntiment Analysis Systems

In Section 3.2, we described the key components of ALLURE along with some additional components we added to make the system more trustworthy and reliable. One of those components is the Sentiment Analyzer (or Sentiment Analysis System (SAS)) whose instability we demonstrate. A tool called ROSE (Tool and Data ResOurces to Explore the Instability of SEntiment Analysis Systems) was proposed in (Mundada et al., 2022) that allow users to examine the instability present in the SAs across different words, sentence structures and gender of the subject. The authors used the dataset described in (Kiritchenko & Mohammad, 2018), Equity Evaluation Corpus (EEC) as one source. That dataset has 8,640 English sentences with different templates in which different pronoun or proper noun variables can be substituted along with different emotion words. "My aunt is feeling miserable" is one such example in which "My aunt" is the pronoun which acts as a proxy for gender and "miserable" is a negative emotion word.

They show instability in six different SAs: the well-known TextBlob, and VADER and four custom-built and trained systems based on published descriptions and training datasets: Convolutional Neural Network (CNN) based implementation (Vishnu, 2020), Long Short Term Memory (LSTM) based implementation, Gated Recurrent Unit (GRU) based implementation and DistilBERT. They took these systems from 'SemEval-2018 Task 1: Affect in Tweets' (Mohammad et al., 2018) along with their respective training datasets. The tool can be accessed [here](#). Figure 6 shows a snapshot of the tool.

Figure 7 shows the average sentiment scores of all the 6 SAs for sentence templates having male gender (pronouns) as person object. The X-axis represents different sentences taken from EEC dataset, and the Y-axis shows sentiment scores ranging from -1 to 1. Off-the-shelf models Vader and TextBlob are consistent throughout the graph for a particular word. But, SAs which are based on Neural Networks, vary with sentence structure. To get more insights into the instability of SAs towards gender, Figure 8 depicts that for CNN, sentiment scores for male and female gender (proper nouns) are different even though they have the same sentence structure and emotion word. Thus, it confirms the presence of bias.

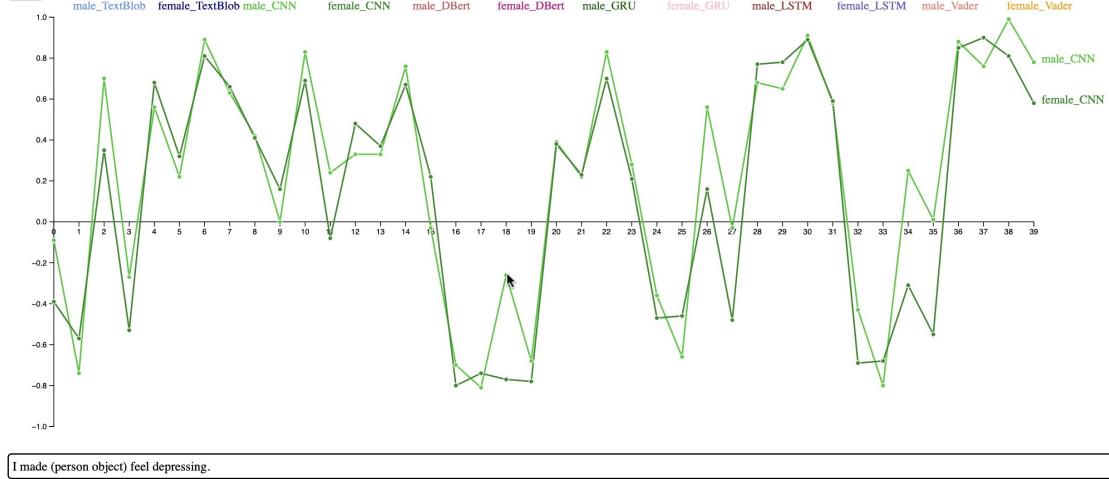
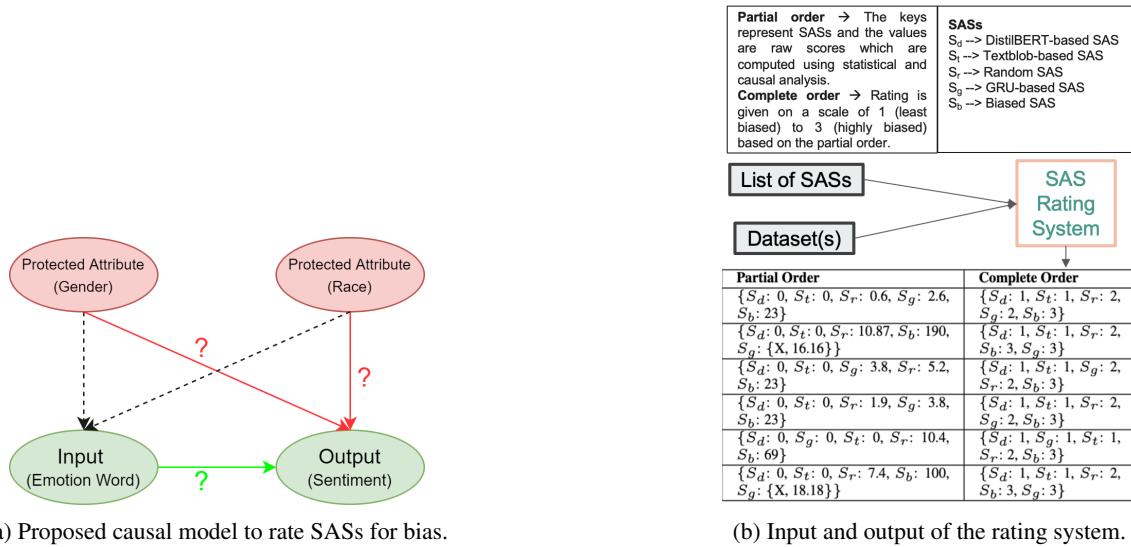


Figure 8: Average sentiment scores of sentences calculated using CNN (for both male and female proper nouns). Each point along the X-axis is a sentence (template with a gender variable), each line is SAS and the Y-axis is a sentiment score.

## D. CausalRating: A Tool To Rate Sentiments Analysis Systems for Bias



The instability of SASSs can be considered as bias in the presence of protected attributes like gender and race. To assess these systems for bias, a method to rate these systems was recently proposed in (Lakkaraju et al., 2023). Figure 9a shows the proposed causal diagram for SASSs. A causal diagram represents the relations between causes and effects. Nodes of the diagram denote attributes of the system. The arrowhead direction shows the causal direction from cause to effect. The green arrow denotes the desirable causal path and the red arrows denote undesirable causal paths. '?' denotes the validity or strength of these causal links that have to be tested using causal or statistical analysis. The dotted arrow indicates that the causal link may or may not be present based on the data. When present, protected attributes act as confounders. A confounder adds a spurious correlation between input and output in a system and is undesirable.

Figure 9b shows the input and output of the rating system. List of SASSs are the sentiment analysis systems that are to be assessed on the data in hand. If any of the protected attributes like race or gender affect the system outcome, then the system is said to be biased. They proposed two metrics in the paper to measure the strength of bias. They are:

**Deconfounding Impact Estimate (DIE):** Deconfounding is any method that accounts for confounders in causal inference. Backdoor adjustment is one such method that was described in (Pearl, 2009). The backdoor adjustment formula is given by

The task you have chosen in the previous page is: Sentiment Analysis

The dataset you have chosen is: datasets/existing/g2\_e3.csv

The system(s) you have chosen is / are:

['Textblob-based SAS']

The metric you have chosen is: DIE %

**Causal Model**

```

graph TD
    PA((Protected Attribute  
Gender)) -.-> IE((Input Emotion))
    PA -.-> OS((Output Sentiment))
    IE -- "?" --> OS
  
```

**Note:**

The metrics are computed using the following formulas:

The WRS for each dataset will be in the range [0, 2.4].

Higher the WRS, the higher will be the effect of protected attributes on the output of the system (statistical bias).

Range of DIE % is  $\mathbb{R}^*$  → where  $\mathbb{R}^*$  is the extended set of real numbers (Range is  $[-\infty, \infty]$ ).

Higher the DIE %, the higher will be the effect of the confounder on the relation between the input and the output of the system (confounding bias).

Figure 10: CausalRating is a tool to rate AI systems for bias. The UI shows the choices made by the user (task, dataset, systems and metric) in the 'Log', Results computed using the chosen metric, Causal model for the chosen task and some additional notes with detailed description of metrics.

the equation 1.

$$P[Y|do(X)] = \sum_Z P(Y|X, Z)P(Z) \quad (1)$$

*Deconfounding Impact Estimation (DIE)* measures the relative difference between the expectation of the distribution, ( $Output|Input$ ) before and after deconfounding. This gives the impact of the confounder on the relation between *Emotion Word* and *Sentiment*. DIE % can be computed using the following equation:

**DIE % =**

$$\frac{[E(Output = j|do(Input = i)) - E(Output = j|Input = i)]}{E(Output = j|Input = i)} * 100 \quad (2)$$

Input is *Emotion Word* and output is *Sentiment*.

**Weighted Rejection Score (WRS):** WRS is calculated when there is no confounding effect (as the backdoor adjustment is not needed in this case). The distribution ( $Sentiment|Protected\ attribute$ ) is compared across different classes using the student's t-test (Student, 1908). They consider three different confidence intervals (CIs): 95 %, 70%, and 60%. For each CI, they calculate the number of instances in which the null hypothesis was rejected for a data group. They multiply this rejection score ( $x_i$ ) with weights ( $w_i$ ) 1, 0.8, and 0.6 for the three CIs respectively. This gives the WRS for a data group in an SAS. WRS is given by the following equation:  $\sum_i w_i * x_i$

The authors of (Lakkaraju et al., 2023) called the above two metrics as raw scores. These were used to compute partial order. From this partial order, they computed final ratings. They are shown in Figure 9b.

In order to make this method more accessible to the end-users, we built a tool called CausalRating that evaluates the AI systems for bias based on the method proposed in (Lakkaraju et al., 2023). Currently, our tool can rate SASs and also German credit dataset (a toy example). A snapshot of our tool is shown in Figure 10. In the UI, log contains the choices made by the user (dataset, task, systems to be evaluated, metric chosen). On the right, you can see the causal diagram. 'Results' give values for both the distributions given in the DIE % definition. We also added a 'Note' with some definitions and range of the computed metric so that it will be easy for the end-users to interpret their result. The tool is still in development. Some functionality is yet to be added.