

Giving AI Personalities Leads to More Human-Like Reasoning

Animesh Nighojkar

ANIGHOJKAR@USF.EDU

Bekhzodbek Moydinboyev

BMOYDINBOYEV@USF.EDU

My Duong

MYDUONG@USF.EDU

John Licato

LICATO@USF.EDU

Advancing Machine and Human Reasoning (AMHR) Lab

Department of Computer Science and Engineering, University of South Florida

Tampa, FL, USA

Abstract

In computational cognitive modeling, capturing the full spectrum of human judgment and decision-making processes, beyond just optimal behaviors, is a significant challenge. This study explores whether Large Language Models (LLMs) can emulate the breadth of human reasoning by predicting both intuitive, fast System 1 and deliberate, slow System 2 processes. Unlike traditional AI research focused on optimizing accuracy, this paper investigates the potential of AI to mimic diverse reasoning behaviors across a human population, addressing what we call the *full reasoning spectrum problem*. We designed reasoning tasks using a novel generalization of the Natural Language Inference (NLI) format to evaluate LLMs' ability to replicate human reasoning. The questions were crafted to elicit both System 1 and System 2 responses. Human responses were collected through crowd-sourcing and the entire distribution was analyzed and modeled, rather than just the majority of the answers. We used personality-based prompting inspired by the Big Five personality model to elicit AI responses reflecting specific personality traits, capturing the diversity of human reasoning, and exploring how personality traits influence LLM outputs. Combined with genetic algorithms to optimize the weighting of these prompts, this method was tested alongside traditional machine learning models. The results show that LLMs can mimic human response distributions, with open-source models like Llama and Mistral unexpectedly outperforming proprietary GPT models. Personality-based prompting, especially when optimized with genetic algorithms, significantly enhanced LLMs' ability to predict human response distributions, suggesting that capturing suboptimal, naturalistic reasoning may require modeling techniques incorporating diverse reasoning styles and psychological profiles. The study concludes that personality-based prompting combined with genetic algorithms is promising for enhancing AI's *human-ness* in reasoning, proposing a new methodology for studying and applying human reasoning by acknowledging and leveraging the vast differences in individual reasoning styles at a granular level.

Keywords: Large Language Models (LLMs), Human-like reasoning, Cognitive biases, Personality prompting, Natural Language Inference (NLI), System 1 and System 2 reasoning

1 Introduction

Capturing the nuanced, often imperfect, and highly diverse reasoning processes of humans presents significant challenges. Traditional uses of AI for this purpose have focused primarily on optimizing accuracy and efficiency, but this approach often neglects the complexity of human cognition, where decisions can be significantly influenced by intuition, emotion, and prior experiences. The central issue is that current AI models struggle to capture the full spectrum of human reasoning, particularly in problem spaces where decisions are not simply correct or incorrect but rather involve a plurality of diverse thought processes. Addressing this limitation is crucial for advancing AI systems that can genuinely understand and interact with humans on a deeper, more intuitive level, making it a critical area of research that extends beyond technical optimization.

Understanding and predicting human reasoning is as fascinating now as it has been for thousands of years (Aristotle, 2013; Wason, 1966). Dual Process Theory, one of the most influential theories in contemporary cognitive science (Kahneman, 2011; Stanovich & West, 2000), distinguishes between two types of cognitive processes: System 1 and System 2.¹ System 1 operates automatically, rapidly, and effortlessly, guiding intuitive and habitual decisions without conscious control. In contrast, System 2 is slower, more deliberate, and responsible for managing mental activities that require focused attention, such as complex calculations and deliberate decision-making. These dual processes are fundamental to understanding how humans navigate decisions, both trivial and significant, and provide a framework for assessing whether AI can genuinely replicate human-like reasoning.

Dual process theories have been given a second look in the age of the large language model (LLM), which have become a cornerstone in AI-driven reasoning due to their tremendous success on a variety of tasks. Some researchers have suggested—or even demonstrated through experimental paradigms—that LLMs with lower complexity or simpler prompts tend to engage in System 1 reasoning, while more complex architectures or prompts encourage System 2-like reasoning (Hagendorff et al., 2023; Kojima et al., 2022; Saha et al., 2024; Weston & Sukhbaatar, 2023; Yu et al., 2024). However, many such claims operate under the assumption that System 1 is inferior and should be avoided, focusing instead on making AI reasoning more like human System 2. The use of datasets such as the Cognitive Reflections Test (CRT) (Frederick, 2005), designed to provoke System 1 into giving the wrong answer while System 2 provides the correct one, is a reflection of this assumption. This focus on accuracy neglects the reality that human reasoning is rarely so clear-cut.

Given these complexities, how well can our best AI approaches capture human reasoning? AI reasoning encompasses a broad range of approaches that have been described at times as efficient (Maclure, 2021), logical (Hagedorn et al., 2020), mechanistic, or capable of making sense of complex scenarios (Zollman et al., 2023). However, it has also been criticized for being biased and opaque (O’Neil, 2017). Numerous studies have aimed to reduce bias in AI reasoning and promote fairness (Ma et al., 2023), yet human reasoning is naturally biased and flawed (Ariely, 2010; Gilovich et al., 2002; Stanovich et al., 2018). Eliminating these flaws can therefore make AI reasoning less human-like (Chemero, 2023;

1. Although in this paper we will use the more popular terms System 1 and System 2, it should be noted that they are somewhat misleading, in that they imply a coordinated system of processes working together. For this reason, some authors prefer the terminology “Type 1” and “Type 2,” e.g., R. Evans et al. (2018).

What is the bat and the ball problem? How is it used in research?

The "bat and ball problem" is a well-known cognitive psychology puzzle often used to illustrate the concept of cognitive biases, particularly the tendency for people to rely on intuitive but incorrect answers rather than deliberate and correct reasoning.

The problem is typically stated as follows: A bat and a ball together cost \$1.10. The bat costs \$1.00 more than the ball. How much does the ball cost?

The intuitive answer many people give is \$0.10, but this is incorrect. The correct answer can be found by setting up an algebraic equation:

1. Let the cost of the ball be x
2. Then, the cost of the bat is $x + 1.00$
3. The total cost is $x + (x + 1.00) = 1.10$

Solving this equation:

$$\begin{aligned} 2x + 1.00 &= 1.10 \\ x &= 0.05 \end{aligned}$$

So, the ball costs \$0.05, and the bat costs \$1.05.

The bat and ball problem is frequently used in research on cognitive psychology, particularly in studies examining:

- Cognitive Biases: It highlights the tendency for people to rely on fast, intuitive thinking (System 1) rather than slower, more deliberate thinking (System 2), as described by Daniel Kahneman in his book "Thinking, Fast and Slow".
- Heuristics: The problem is used to study heuristics, or mental shortcuts, that people use to make decisions. It demonstrates how these heuristics can lead to systematic errors.
- Dual-process Theory: Researchers use the problem to explore the dual-process theory of cognition, which posits that human thought operates on two levels: an automatic, fast, and intuitive level (System 1) and a slower, more analytical, and deliberate level (System 2).
- Educational Interventions: It is also used in educational research to study how teaching methods can help students develop better problem-solving strategies and reduce reliance on incorrect intuitions.

In sum, the bat and ball problem is a simple yet powerful tool used to explore how people think and reason, revealing the common errors that arise from our cognitive processes.

Figure 1: A question to and response by ChatGPT (GPT-4o)—without using browsing capabilities—shows that the LLM is familiar with this common question from cognitive science already, and thus questions of this type have questionable validity when used to assess its reasoning.

Prescott & Wilson, 2023). Furthermore, human reasoning processes are deeply influenced by context, individual differences, and prior knowledge, making it challenging to predict human decisions based purely on logical correctness (Kahneman & Tversky, 1984; Nisbett & Wilson, 1977; Stanovich & West, 2000). And thus, artificial reasoners that prioritize normative correctness by reducing errors (Sun et al., 2024) may become less capable of modeling human-like reasoning.

In this paper, we advocate for an approach to computational cognitive modeling that considers an approach successful only to the extent that it can predict the *full distribution* of naturalistic behaviors in some task of interest. We will refer to this challenge as the *full reasoning spectrum problem*. Perhaps nowhere is this problem more apparent than in the blurred boundaries between System 1 and System 2. Faced with a reasoning problem, individuals may arrive at the same conclusion through different cognitive paths, or different individuals may use the same cognitive process to reach different conclusions (De Neys & Pennycook, 2019; Nye et al., 2021). Hamade et al. (2024) argue that “achieving superhuman performance is not sufficient; AI needs to account for suboptimal actions or idiosyncratic style from humans.” Thus, it is important not only to capture whether human reasoners get the answer right or wrong, but *what types* of wrong answers they guess, in what proportions those wrong answers are given in a sample of human reasoners, and how those proportions are affected by the environments in which the individuals perform their reasoning. The *full reasoning spectrum problem* is the challenge of designing artificially intelligent reasoning systems that can predict all of these.

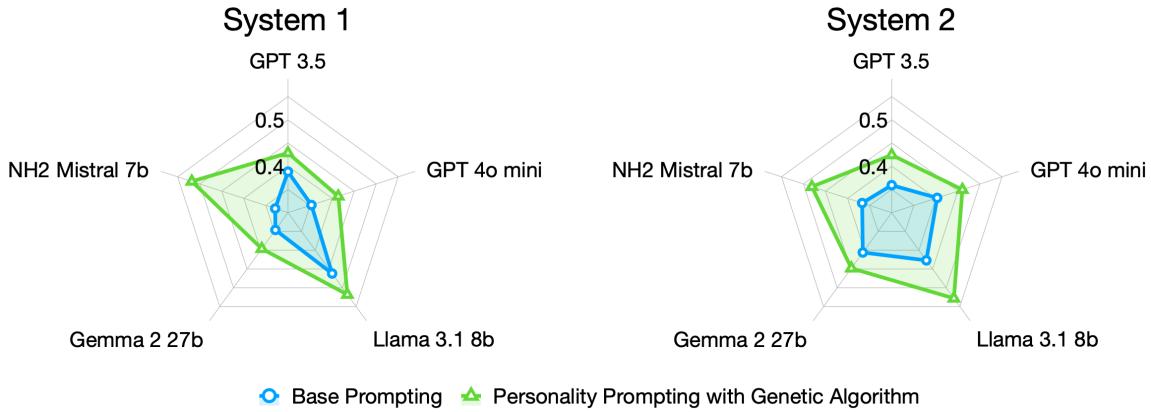


Figure 2: Similarity between human reasoning (System 1 and System 2) with LLMs’ reasoning. Base prompting just prompts the LLMs for the answer and personality prompting prompts the LLM with different personality prompts and finds the best weight for each prompt using a genetic algorithm with K-fold cross validation. More details in Section 4.

Moreover, familiarity with specific problems can significantly influence the cognitive process employed by both humans and LLMs. For humans, prior experience with a problem might lead them to rely on System 1 for quick, intuitive responses, while unfamiliarity might prompt more deliberate System 2 reasoning (J. S. B. T. Evans, 2008; Kahneman, 2011; Klein, 2017). For LLMs, this raises concerns about the *data contamination problem* (Balloccu et al., 2024; Sainz et al., 2023), where LLMs might rely on memorized information from their training data rather than engaging in genuine reasoning. Given that many modern LLMs (Almazrouei et al., 2023b; Brown et al., 2020; Dubey et al., 2024; Jiang et al., 2023) are trained on vast datasets collected from the internet (Computer, 2023; Penedo et al., 2023; Raffel et al., 2023), there is a significant risk that the questions used in studies were part of the LLMs’ training data, further complicating efforts to assess true reasoning capabilities. To illustrate this issue, consider the study by Hagendorff et al. (2023), which examined ChatGPT’s reasoning abilities using just three problems from the CRT as templates to create 150 nearly identical problems, altering only the numbers and objects. A very simple question to ChatGPT (shown in Figure 1) reveals that the LLM is familiar with not only the problem but also the conclusions researchers draw from human responses to it. This raises doubts about the validity of inferences made about LLMs’ reasoning abilities. Employing LLMs to generate fresh, original questions that distinguish between System 1 (intuitive, fast) and System 2 (deliberate, slow) reasoning processes helps avoid this problem.

We design questions in the Natural Language Inference (NLI) format (discussed in Section 2.2), incorporating diverse reasoning structures like syllogisms, fallacies, and belief biases (J. S. B. T. Evans et al., 1983; Sap et al., 2022; Yang et al., 2023). Additionally, we modify the NLI format from three options to six, providing more granularity and nuance in human responses. Using the concept of personality prompting (De Paoli, 2023; Gu et al., 2023; Kamruzzaman & Kim, 2024), we aim to predict not just the majority response (or the gold label) but the *entire distribution* of human responses to a given question, accounting for

the likelihood that the humans participating in our study may possess different personality traits. These personality prompts are created from the traits in the Big Five Model (Rocca et al., 2002): openness, conscientiousness, extraversion, agreeableness, and neuroticism. We find that performance on this task improves when we allow the LLM to emulate different human personalities. Each prompt in this suite of personality prompts is assigned a weight that determines the frequency of that prompt’s response in our final pool of responses. These weights (one for each prompt) are fine-tuned using a genetic algorithm (Mitchell, 1996), and the final distribution is compared with human responses in Figure 2. This entire process and its results are detailed in Section 4.

Contributions and Key Findings

- We are the first to investigate whether AI (LLMs) can model System 1 and System 2 as separate, distinct reasoning processes, specifically in terms of their ability to capture the *full distribution* of human responses.
- We introduce a six-way division of the natural language inference problem that allows for a more granular analysis of response distribution compared to existing research.
- We present a pioneering approach to the full reasoning spectrum problem, demonstrating the necessity of capturing the full range of correct and incorrect human responses to advance the fidelity of AI in cognitive modeling.
- Our study establishes a novel metric for assessing AI models’ capability to encompass the entire spectrum of human reasoning, enhancing the evaluation of AI’s mimicry of diverse cognitive processes from intuitive to analytical reasoning.
- Though Big Five personality traits have been used for personality detection with LLMs (Kazemeini et al., 2021; Pellert et al., 2024; Wen et al., 2024), we are the first to use the Big Five traits for personality prompting for NLI.
- We find that personality-based prompting significantly enhances LLMs’ ability to mimic human response distributions compared to vanilla zero-shot prompting.
- The open-source Llama and Mistral models outperform closed-source GPT architectures, supporting the argument that current open-source LLMs are at least as effective as (and possibly better than) closed-source LLMs in capturing naturalistic human reasoning (Ahmed et al., 2024).

2 Background

2.1 System 1 vs System 2 Reasoning in LLMs

Dual Process Theory distinguishes between two types of cognitive processes: System 1, which is fast, intuitive, and automatic, and System 2, which is slower, more deliberate, and rational (J. S. B. T. Evans & Stanovich, 2013). This framework is highly influential on current theories of human reasoning, judgment, and decision-making. For instance, classic problems like the Bat and Ball Problem (Frederick, 2005), also shown in Figure 1,

demonstrate how System 1 can produce rapid but often flawed responses, while System 2 enables more accurate outcomes through reflective thinking (J. S. B. T. Evans, 2008, 2019).

Recent advancements in AI have sought to replicate System 1 and System 2 processes using various prompting strategies in LLMs (Hagendorff et al., 2023; Kojima et al., 2022; Nye et al., 2021; Saha et al., 2024; Weston & Sukhbaatar, 2023; Yu et al., 2024), with mixed results (Vatsal & Dubey, 2024). Researchers have attempted to simulate System 1 and System 2 reasoning in LLMs by experimenting with common prompt templates, such as zero-shot and chain-of-thought prompts (Kojima et al., 2022; Wei et al., 2022). Additionally, efforts to emulate System 2 thinking have led to the development of novel prompting techniques that expand on chain-of-thought, mimicking aspects of human problem-solving processes using search algorithms (L. Wang et al., 2023; X. Wang et al., 2023; Yao et al., 2023; Y. Zhang et al., 2024). In contrast, modeling System 1-like processes in LLMs has involved heuristic-driven, rapid decision-making approaches (Hagendorff et al., 2023), which can generate responses without immediate reasoning steps (Yu et al., 2024). On the other hand, simulating System 2 requires more complex logical reasoning and problem-solving tasks that involve deliberate, step-by-step processing (Nye et al., 2021). Some researchers have also explored hybrid models that integrate both cognitive systems, aiming to harness their combined strengths for more effective decision-making in AI (Saha et al., 2024).

In their survey, Vatsal and Dubey (2024) examine how prompt engineering can evoke distinct reasoning and responses from LLMs. They categorize prompt engineering techniques into zero-shot, one-shot, and few-shot prompts, each with specific benefits and limitations. Zero-shot prompting involves giving the model a task description and expecting it to perform the task based on its pre-existing knowledge, though this approach may be constrained by the model's interpretation of the task description. One-shot prompting offers a single example to provide a clear reference for the task, while few-shot prompting supplies multiple examples, offering richer context and enhancing task generalization, which has been shown to significantly improve performance. The survey explores how these methods influence LLMs' effectiveness in various NLP tasks, including question answering, text generation, and language inference.

A prevalent assumption in the literature is that System 1 leads to the wrong answer. System 1 processes, often seen as intuitive but error-prone, are typically simulated in LLMs through straightforward, concise zero-shot prompts (Hagendorff et al., 2023; Kojima et al., 2022; Yu et al., 2024). In contrast, System 2 processes, associated with more accurate and thoughtful responses, are mimicked through chain-of-thought prompting, which enhances performance on tasks like the Cognitive Reflection Test (CRT) (Hagendorff et al., 2023). However, this assumption limits the datasets used in these studies to those with clear correct or incorrect answers, overlooking the nuance of real-world reasoning. It is important to note that System 1 reasoning is not simply erroneous reasoning, but rather it is a set of reasoning processes (akin to heuristics) that has distinctive properties, and thus it can and should be modeled separately from System 2 reasoning. Moreover, System 1 processes are an inseparable step of human reasoning (Gilovich et al., 2002) and cannot be ignored. System 1 reasoning enables us to perform many tasks effortlessly, such as walking, eating, or driving, where over-analyzing every detail would be impractical or impossible. System 1 is almost effortless for humans (Gladwell, 2007) but as we demonstrate in this paper, replicating

System 1 behaviors in LLMs is more complex than simply applying zero-shot prompts, especially with datasets lacking normatively correct or incorrect responses.

2.2 The Natural Language Inference (NLI) Format

In our study, we utilize an enhanced version of the Natural Language Inference (NLI) task, building upon prior foundational work (Bowman et al., 2015). NLI, a successor to Recognizing Textual Entailment (RTE) (Dagan et al., 2006),² is centered around evaluating whether a hypothesis h can logically follow from a given premise p . For illustration, consider the premise p = “A crowd gathers as three blue cars begin a race.” Possible relationships between this premise and a hypothesis h could be:

1. Entailment: h = “A race is taking place.” must be true if p is true.
2. Contradiction: h = “There are no cars racing.” cannot be true if p is true.
3. Neutral: h = “Three men are competing in a race.” is neither necessarily true nor necessarily false given p .

The most common goal of NLI datasets is to study how a set of human reasoners *naturally* reasons about the inferential relationships in the NLI items. LLMs are traditionally trained and benchmarked on NLI datasets to enhance their naturalistic reasoning capabilities (Bowman et al., 2015; Nie, Williams, et al., 2020; Williams et al., 2018, 2022).

While typically trifurcated into distinct categories, this specific categorization can suffer from problems of underspecification and a lack of granularity (Nighojkar et al., 2023). Addressing these complexities, researchers have proposed various modifications to the standard categorization, including:

- Introducing an “entailment strength” parameter, reflecting either model confidence or perceived likelihood, though such measures have historically grappled with issues like ambiguous evaluation standards and inconsistent annotator judgments (Chen et al., 2020; Meissner et al., 2021).
- Examining the variability in annotator perspectives, considering whether a ‘neutral’ judgment indicates balanced reasons for both agreement and disagreement or a complete absence of decisive evidence (Pavlick & Kwiatkowski, 2019; S. Zhang et al., 2021; X. F. Zhang & Marneffe, 2021; Zhou et al., 2022).
- Differentiating between various degrees of entailment, such as “absolutely entails” versus “likely entails,” which echoes ongoing research into probabilistic reasoning (Kahneman, 2011).

To refine these categories further, we propose a structured NLI framework with six detailed categories, broadening the inferential spectrum from absolute contradiction to definite entailment. These categories are: (A) *Absolutely must be false*, (B) *Is more likely to be*

2. RTE focused only on predicting whether the relationship between two sentences was entailment, but NLI bifurcated the “False” label for more granularity

false, (C) *Has strong reasons to be true and strong reasons to be false*, (D) *Has no reasons to be either true or false*, (E) *Is more likely to be true*, and (F) *Absolutely must be true*.

This broader categorization not only clarifies the guidance provided to annotators but also enriches the nuances in the data, aligning with earlier initiatives to expand NLI’s analytical depth (Pavlick & Kwiatkowski, 2019; S. Zhang et al., 2017). This framework allows for more detailed response options while still fitting into the three traditional NLI categories, making it more flexible. The use of NLI is widespread, affecting fields such as measuring psychological traits (Laverghetta Jr. & Licato, 2022, 2023a, 2023b; Laverghetta Jr. et al., 2021, 2022), comparing text similarity (Nighojkar & Licato, 2021b), and evaluating the quality of translations and paraphrases (Nighojkar & Licato, 2021a). This demonstrates its versatility and usefulness in different areas.

3 Data Collection

3.1 Survey Design

Our study utilizes a dual-phase survey design, hosted on Qualtrics, with 60 participants recruited via Prolific. The survey consists of 27 questions in the first phase, which are repeated in the second phase, following the design of the two-response paradigm (Bago & De Neys, 2017). In this paradigm, participants are asked to provide two separate responses to the same question or task. The first response is typically their initial, instinctive answer, while the second response is given after some deliberation or additional information is provided. This method allows us to analyze the differences between intuitive and reflective thinking, or System 1 and System 2, shedding light on cognitive processes like reasoning, decision-making, and how people change their minds when given more time or information. In the first phase, aimed at eliciting System 1 reasoning, participants are presented with NLI questions while simultaneously solving a puzzle as a distraction. Each question is divided into two parts: the premise and the conclusion, with six response options (A-F) available, as detailed in Section 2.2. The exact instructions provided to the participants are in Figure 3. Participants read the premise (Figure 4), solve the puzzle (Figure 5), then read the conclusion before selecting one of the six options (Figure 6). Note that participants never have access to the premise and the conclusion together in this phase and need to rely on whatever they retain from the premise when answering the question about the conclusion on Screen 3. The survey begins with an example question to familiarize participants with the task. To reduce fatigue, one-minute breaks are incorporated after each question (a set of all three screens).

In the first phase, a time limit is imposed to encourage rapid reasoning. The initial time allocated per question is determined using the following formula:

$$\text{base time} = 0.0787 \times \text{sentence length} + 0.0016 \times \text{perplexity} + 6.3276$$

These coefficients are derived from ordinary least squares regression (Zdaniuk, 2014), based on data collected in-house. We recorded the time each participant (the co-authors and other volunteering students) took to read a question and trained the regression model to predict the 90th percentile reading time based on the sentence length and its perplexity (Jelinek et al., 1977). Additionally, the reading time is adaptively adjusted based on participant pacing, either increasing or decreasing the time allotted depending on their speed, with a

GIVING AI PERSONALITIES LEADS TO MORE HUMAN-LIKE REASONING

Thanks for consenting to participate in this survey!

We're timing each page! If you happen to miss the timer on any page, you'll be redirected to the end of the survey.

Survey-Specific Instructions

- You will see a question: Premise and Conclusion. Your task is to assume the premise is true and choose one of the six available options for the conclusion.
- You will have to choose one of six options. These options remain the same for all questions.
- You will be shown just the premise first. Then you will be asked to answer a question that might not be related to the premise.
- Then you will see just the conclusion and need to answer it based just on the premise you had seen earlier. **Pick the first answer that comes to your mind. Don't overthink it!**
- *You need to answer within the time limit. Otherwise, the survey ends.*
- **You are not allowed to take notes during the survey.**

Example:

Premise: All dogs are animals.

Conclusion: All animals are dogs.

Assuming the premise is true, the conclusion _____

- Absolutely must be false
- Is more likely to be false
- Has strong reasons to be true and strong reasons to be false
- Has no reasons to be either true or false
- Is more likely to be true
- Absolutely must be true

On the next screen, you will see a demo of the real task. Proceed to the next page after reading the premise. Solve the puzzle you see and proceed to the next page. Then you'll see a conclusion and a few choices. These choices remain the same throughout the survey and it'll be a good idea to remember what they are because the time given for subsequent questions will be very less.

Figure 3: Survey instructions

Premise: Rashim has a specific work authorization he uses for employment. He has a lot of family members living in the US. His first language is English. He is not very interested in politics. Finally, he has a lot of personal experience with discrimination.

Time remaining to read: 4.0

Figure 4: Survey screen 1: Premise

Even though I'm far from the point, I'm not wrong. Put me in the right place, and you'll find me quite strong.

Time remaining to read: 18.3

- A well-kept secret lied
- An emblem of great pride
- A misplaced decimal, hidden wide
- An undeciphered code, unspecified

Figure 5: Survey screen 2: Distractor

Conclusion: Rashim is an immigrant to the United States

Assuming that the premise is true, the conclusion _____.

Time remaining to read: 3.8

- absolutely must be false
- is more likely to be false
- has strong reasons to be true and strong reasons to be false
- has no reasons to be either true or false
- is more likely to be true
- absolutely must be true

Figure 6: Survey screen 3: Conclusion

minimum engagement time set at 10 seconds. If a participant moves on from Screen 1 or answers the question on Screen 3 within 70% of the base time for that screen, we reduce their time for all the following questions by 10%. This reduction is compounded, so if the participant takes less than 70% of the reduced time, we reduce the time for the subsequent questions by 19% (10% compounded twice). We never reduce the time below 10 seconds. A similar calculation works for slow participants. If they take over 90% of the base time, we increase the time for the subsequent questions by 10%. Failure to comply with the time limits results in warnings and may lead to survey termination after repeated infractions.

The puzzles, essential for the distraction component, are generated using OpenAI's GPT-4-0314 model. A structured query produces rhyming riddles with four potential answers, designed to engage participants critically without overwhelming them. Figures 4-6 illustrate the three components of a question in the first phase. The second phase omits the puzzles and combines the premise and conclusion into a single screen per question, allowing participants two minutes to carefully consider each NLI item and record their reasoning in a text box, explicitly targeting System 2 reasoning.

To detect bad faith participation, embedded attention checks resembling standard questions are distributed randomly through phase 1 (questions in phase 2 follow the same order as phase 1), drawn from ChaosNLI (Nie, Zhou, & Bansal, 2020) with at least 90% inter-annotator agreement. ChaosNLI is an NLI dataset created with 100 annotators for each question. This lets us set thresholds on inter-annotator agreement (like 90%) that would not be possible with other NLI datasets like SNLI (Bowman et al., 2015) or MultiNLI (Williams et al., 2018) that only have 5 annotations per question. If a participant’s responses fail to match more than 5 out of the 8 attention check questions across both phases, their data is excluded, and a new participant is recruited to ensure 30 valid, good-faith responses per question. Additionally, responses from participants who do not complete the entire survey are disregarded.

3.2 Item Types

In Section 1, we introduced the full reasoning spectrum problem (where existing datasets and methods fail to capture the whole spectrum of human reasoning patterns) and the data contamination problem (the possibility that publicly available datasets may be included in the training data of large language models, or LLMs). In Section 2.2, we also explained our decision to use NLI as our preferred format. Commonly available NLI datasets do not attempt to capture annotations that distinguish between System 1 and System 2 reasoning in humans. Additionally, it is highly likely that existing public NLI datasets have already been used to train LLMs. Consequently, these datasets suffer from both the full reasoning spectrum problem and the data contamination problem. To address these issues, we created our own dataset of NLI questions using OpenAI’s GPT-4 model (specifically GPT-4-0314) to generate all the items. We provided the model with multiple prompts and used different OpenAI API calls to generate the premises and conclusions separately, ensuring that the model never produced both simultaneously. Finally, we made two additional API calls to rephrase the premises and conclusions individually. All our prompts are in Appendix B along with details about our data generation pipeline. We generate different types of items to prevent monotonicity in the dataset and to prevent any pattern recognition for the humans doing our survey and the models we use to predict the human responses. These item types are explained in more detail below.

StereoNLI StereoNLI connects NLI with stereotypes, building on previous research showing that human System 1 reasoning is often biased and influenced by stereotypes (Geeraert, 2013; Kahneman, 2011). We utilized the StereoSet dataset (Nadeem et al., 2021) to select seed words for generating StereoNLI questions. StereoSet comprises 17,000 sentences that assess stereotype bias in language models concerning gender, race, religion, and profession. Each sentence is labeled by multiple human annotators as either ‘anti-stereotype’, ‘stereotype’, ‘unrelated’, or ‘related’. From StereoSet, we randomly chose nouns associated with gender, race, religion, and profession, then prompted GPT-4-0314 to generate a name for a hypothetical person called ‘X’. We asked GPT-4 to write a paragraph with three sentences reflecting common assumptions about ‘X’ based on the given traits. This paragraph became our premise. Excluding gender due to its limited variety in the dataset, we randomly selected one of the three traits and prompted GPT-4 to generate an assumption about ‘X’ with a truth value. For each scenario, we created three instances—contradiction, neutral,

and entailment—over three separate conversations with the GPT-4 chatbot via the OpenAI API. Finally, the premises and all three conclusions were rephrased. Table A presents three examples of StereoNLI items that share the same premise.

Fallacy Fallacies of argumentation are arguments whose logical structures that have very little to no deductive validity, despite being commonly used in informal reasoning.³ Previous cognitive psychology research (Boissin et al., 2022) has examined how System 1 can sometimes engage in fallacious reasoning. In our dataset, we focus on three specific types of fallacies. The first, *post hoc ergo propter hoc*, occurs when someone mistakenly believes that because one event follows another, the first event caused the second. The second, *slippery slope*, posits that an action will set off a chain of events leading to an undesirable outcome without establishing or quantifying the relevant contingencies. This fallacy, also known as “the domino effect,” often implies a long series of intermediate events connecting a seemingly harmless start to an undesirable end, assuming uncertain or unlikely consequences. The third fallacy, *straw person*, involves misrepresenting an opponent’s argument to make it easier to refute. We created templates to generate premise-conclusion pairs for each type of fallacy. For example, in the case of *post hoc ergo propter hoc*, the premise follows the structure “X happened right before Y,” with the conclusion stating “X caused Y.” GPT-4 fills in the details for X and Y, and we then rephrase the premises and conclusions to produce the final items. Table A displays examples of Fallacy items, one for each fallacy.

Syllogism and Stereo Syllogism Syllogisms, a core component of traditional logic used in philosophical reasoning, consist of a major premise (a general statement), a minor premise (a specific statement), and a conclusion. The structure of a syllogism, or its figure, depends on the positioning of the middle term (M), subject (S), and predicate (P). The complex structure of syllogisms has led many researchers to examine its effects on eliciting System 1 and System 2 reasoning (Da Silva, 2023; J. S. B. T. Evans et al., 1983; Khemlani & Johnson-Laird, 2012). For our dataset, we selected singular nouns as seed words to guide the GPT-4 model in generating the major and minor premises while preserving the order of M, P, and S. The four primary syllogistic figures are: (1) *Premise: M is P. S is M. Conclusion: S is P.* (2) *Premise: P is M. S is M. Conclusion: S is P.* (3) *Premise: M is P. M is S. Conclusion: S is P.* (4) *Premise: P is M. M is S. Conclusion: S is P.* Finally, we rephrased the premises and conclusions to produce the final items. We also created a variant of syllogism questions using seed words from StereoSet, referred to as Stereo Syllogism. Table A presents four examples of Stereo Syllogism items, one for each figure.

Guilt Drawing on studies that demonstrate dual-process effects when participants are asked to assess the guilt of a suspect in a hypothetical scenario (Bergius et al., 2020; Kassin et al., 2013; Peer & Gamliel, 2013; Rachlinski et al., 2015; Wistrich et al., 2015), we designed questions using the following template: First, we prompted GPT to generate a few sentences describing the details of a crime for which the perpetrator has not been identified. We then asked GPT-4 to list several features likely to be true of the culprit. This enabled us to create two forms of questions: the entailment form (*E Guilt*), which included the crime description, a suspect characterized by features that made them appear likely to be guilty,

3. For more nuanced discussions of what constitutes a fallacy, see Farnside and Holther (1959); D. Walton (1985); D. Walton et al. (2008); D. N. Walton (1990).

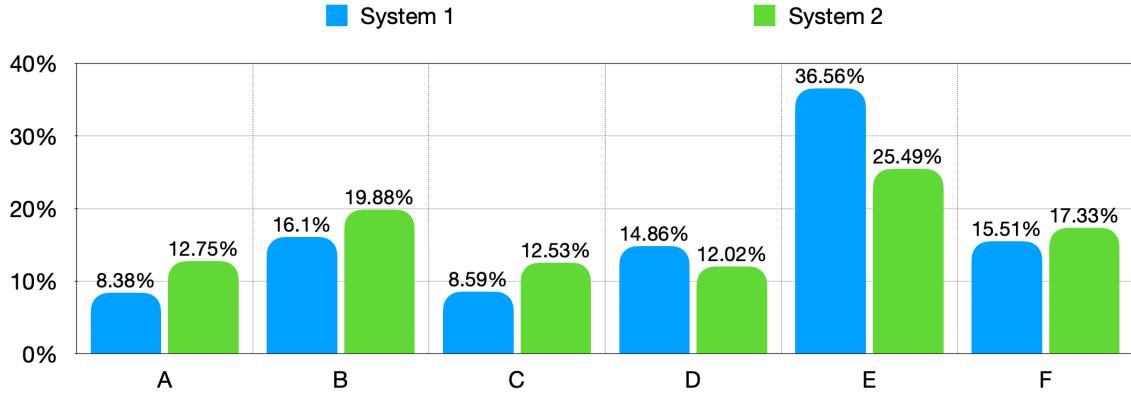


Figure 7: Label percentage distribution histogram

and a conclusion asking whether the suspect is guilty; and the contradiction form (*C Guilt*), where the suspect is described as not possessing the features, or as having opposite features. Table A shows examples where the *E Guilt* premise describes a suspect who initially seems likely to have committed the crime, while the *C Guilt* premise describes a suspect who does not seem to fit the crime’s circumstances.

Primacy and Recency The first and last pieces of information presented can disproportionately influence the reader’s perception due to what is known as the *serial position effect* (Murdock, 1962). To incorporate questions leveraging this effect, we began by manually writing several sentences that simply stated an individual’s name and a characteristic (e.g., “Simon is a professor”). We then asked GPT-4 to generate five sentences about that individual that would likely be true if the characteristic were accurate (*likely-true*). We selected two sentences that best aligned with our archetypal conceptions of the characteristic. Next, we asked GPT-4 to generate five sentences about that individual that were likely *not* true (*likely-false*) and selected three of these. We then created two forms of questions: *P Primacy/Recency* questions had premises starting and ending with sentences from the *likely-true* category, with three *likely-false* sentences in the middle, and the original sentence describing the individual and characteristic as the conclusion. *N Primacy/Recency* questions had the same premise, but the conclusion was the original sentence in a negated form (e.g., “Simon is not a professor” instead of “Simon is a professor”). Examples of these questions are provided in Table A.

3.3 Dataset Characteristics

We gathered 45 items across all item types (10 StereoNLI, 7 Fallacy, 8 Syllogism/Stereo-Syllogism, 10 Guilt, and 10 Primacy/Recency), each annotated 30 times by humans during the first phase (System 1) and an additional 30 times by the same individuals in the second phase (System 2). Figure 7 provides a quantitative comparison of voting preferences between System 1 and System 2, highlighting how the distribution of votes across six options changed from System 1 to System 2. Notably, option E exhibited the most significant shift, decreasing by 11.07 percentage points.

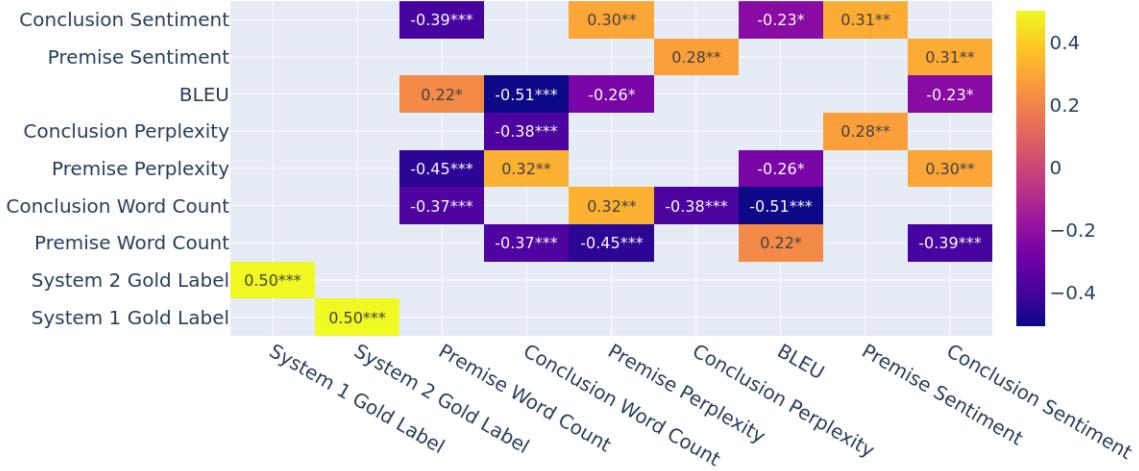


Figure 8: Kendall’s τ between various features of the items in our dataset. Only statistically significant ($p\text{-value} < 0.05$) τ values are shown. *** indicates a $p\text{-value} < 0.001$, ** indicates a $p\text{-value} < 0.01$, and * indicates a $p\text{-value} < 0.05$.

To ensure the reliability of our results in subsequent experiments, we examined potential spurious correlations within our dataset, focusing on four primary factors: the length of the premise and conclusion, the perplexity of the premise and conclusion, BLEU scores (Papineni et al., 2001) between the premise and conclusion, and VADER sentiment scores (Hutto & Gilbert, 2014) for both the premise and conclusion. Figure 8 displays the Kendall’s τ (Kendall, 1938) between these factors and the mean value of participant responses. Ideally, no feature should correlate with the gold labels from either System 1 or System 2, as such a correlation could be considered spurious. The highest τ value is observed between System 1 and System 2 gold labels, indicating that a participant’s responses in System 1 may be predictive of their responses in System 2. However, in this paper, we treat System 1 and System 2 independently and do not attempt to predict System 2 responses based on System 1, leaving this for potential future work.

Figure 9 presents Sankey diagrams to visually represent the transition of human responses from System 1 to System 2, illustrating the distribution and shifts across our six labels by varying the width of the links. This visualization allows us to easily see how successful each question type is at highlighting the differences between System 1 and System 2 reasoning are, as well as how the distribution of human responses differs with each question type. We observe that most participant responses fluctuated between E and D for StereoNLI, and between E and F for fallacy, syllogisms, and stereo syllogisms. Conversely, E guilt questions exhibited minimal fluctuation, with most participants maintaining their E response from System 1 to System 2. This consistency is expected, as is the trend of shifting from E to B for C Guilt, since these questions are designed to mislead participants in System 1. This misleading effect is also observed, though to a lesser extent, in P Primacy and N Primacy. Another noteworthy trend is that participants generally avoided selecting

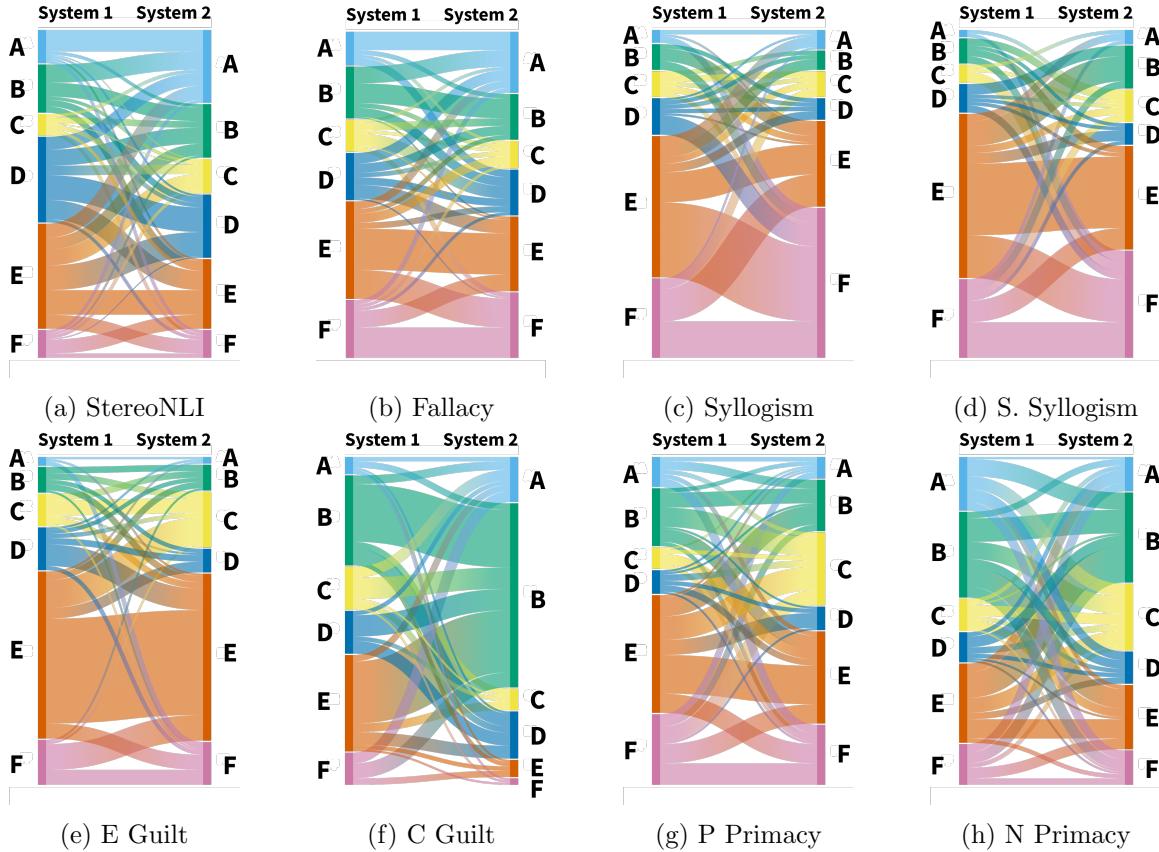


Figure 9: Sankey diagram of changes in human responses from System 1 (phase 1 of the survey) to System 2 (phase 2 of the survey) for all item types in our dataset.

the neutral options C and D for any question type, suggesting an underlying bias. Further research may shed light on the reasons for this behavior.

Inclusion and Ethics Statement Our study was designed and conducted with careful consideration for ethical research practices and inclusivity. We recruited 60 participants via Prolific, aiming for a diverse sample. All participants provided informed consent and were clearly instructed about the study’s nature. Data privacy and protection measures were implemented to ensure participant anonymity and confidentiality. We incorporated breaks to reduce fatigue and set time limits to manage cognitive load, considering participant well-being. The study design included various item types to prevent monotonicity and reduce potential biases. We used AI models (GPT-4) to generate survey items, addressing potential data contamination issues and creating a novel dataset. We have been transparent about our methodology, including the use of AI-generated content and both proprietary and open-source language models in our experiments. By adhering to these ethical principles and inclusive practices, we aim to contribute to the field of AI and cognitive science research in a responsible and equitable manner.

Trait	Prompt
O+	You’re open to new experiences, creative, inventive, curious, and imaginative.
O-	You prefer routine and familiarity, consistent, conventional, and cautious.
C+	You’re organized, efficient, reliable, and responsible.
C-	You’re flexible, spontaneous, extravagant, and careless.
E+	You’re friendly, outgoing, sociable, and energetic.
E-	You’re reserved, quiet, introverted, and solitary.
A+	You’re cooperative, warm, friendly, and compassionate.
A-	You’re competitive, detached, critical, and judgemental.
N+	You’re anxious, stressed, nervous, and emotionally sensitive.
N-	You’re calm, stable, confident, and emotionally resilient.

Table 1: Personality Prompts

4 Modeling Experiments

The overarching aim of our modeling experiments is to predict how human votes are distributed across each item. We begin by predicting the mode of the distribution (majority vote, as discussed in Section 4.1), then proceed to estimate the variance among the votes (Section 4.2), and ultimately aim to predict the entire distribution of human votes (Section 4.3). None of these tasks can be considered inherently easier than the others due to the significant variability in human reasoning, as highlighted in Section 3.3. Given our development of the 6-way scheme, which includes six labels as detailed in Section 3, we also compare the performance of LLMs on the more conventional NLI scheme: the 3-way classification (*contradiction*, *neutral*, and *entailment*). It is important to note that we collect human responses solely for the 6-way scheme and subsequently map these responses to the 3-way scheme by categorizing A and B as *contradiction*, C and D as *neutral*, and E and F as *entailment*. All LLMs are tested separately on both the 3-way and 6-way schemes. We utilize both proprietary LLMs (GPT-3.5,⁴ GPT-4 (OpenAI et al., 2023), GPT-4o-mini⁵) and open-source LLMs (Gemma 2 27b (Team, 2024), Llama 3.1 8b, Llama 3.1 70b, Llama 3.1 405b,⁶ Mistral 7b with direct performance optimization (Rafailov et al., 2023; Teknium et al., 2024)) to compare their performance against traditional (non-transformer-based) machine learning algorithms.

To predict the entire vote distribution using LLMs, we must prompt them in a way that allows the recreation of a vote distribution rather than producing a single vote per item. We do not fine-tune the LLMs because we do not have enough data to tune millions of parameters that we would need to tune even with techniques such as parameter-efficient fine-tuning (Ding et al., 2023). All the LLMs we use are chatbot models, which include a system prompt that influences the LLM’s responses and a user prompt (both described below) that contains the specific query. We maintain consistency in the user prompt by limiting it to just the premise and hypothesis (we call these $s1$ and $s2$ in our prompts, not

4. <https://platform.openai.com/docs/models/gpt-3-5-turbo>

5. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>

6. https://github.com/meta-llama/llama-models/blob/main/models/llama3.1/MODEL_CARD.md

to be confused with System 1 and System 2, which are never shortened in this paper). We develop two distinct task definitions, one for each labeling scheme, and employ two different prompting styles. The task definitions are as follows:

- **6-way:**

Assuming s1 is true, choose the statement that seems most accurate for s2:

- A. Absolutely must be false
- B. Is more likely to be false
- C. Has strong reasons to be true and strong reasons to be false
- D. Has no reasons to be either true or false
- E. Is more likely to be true
- F. Absolutely must be true

- **3-way:**

Choose one option about the inferential relationship between s1 and s2:

Entailment: s2 entails s1

Contradiction: s2 contradicts s1

Neutral: Cannot pick either of the above or both are likely

The first prompting style, called *base prompting*, simply explains the task and asks the LLM to provide its prediction. The exact base prompt is: {task definition} Pick exactly one option and write it on the first line. Do not write anything else.” The second prompting style, *personality prompting*, adds a brief description of the desired personality for the LLM to adopt, in addition to the base prompt. The personality prompt reads: “Here’s your personality: *personality*. Focus on this personality and respond just like a person who has this personality. {task definition} Pick the first answer that you think of based on your personality and nothing else. Pick exactly one option and write it on the first line. Do not write anything else.” The different personalities are derived from the Big Five personality traits (OCEAN model) (Rocca et al., 2002), with each trait simplified to a high (+) or low (-) value to create ten different personality prompts shown in Table 1.

For our experiments, we use the same set of responses generated by all LLMs. For base prompting, we obtain this set by prompting the LLM ten times with a temperature setting of 1 (to maximize response entropy) to generate a distribution with ten votes. For personality prompting, we use each of the ten personality prompts to generate ten votes, then select the majority vote from each prompt’s votes to compile a set of ten votes. This approach yields ten votes each for the base and personality prompting styles. It is important to note that the LLMs are not fine-tuned for this task, in contrast to the classical machine learning algorithms, which are fine-tuned. This intentional disparity aims to test the LLMs’ innate ability to replicate human responses.

4.1 Are human responses predictable?

The task of predicting the mode (also referred to as the gold label) of human responses can be approached as a classification problem, since all categories are different from one another. While it can also be framed as a regression problem, this approach presents several

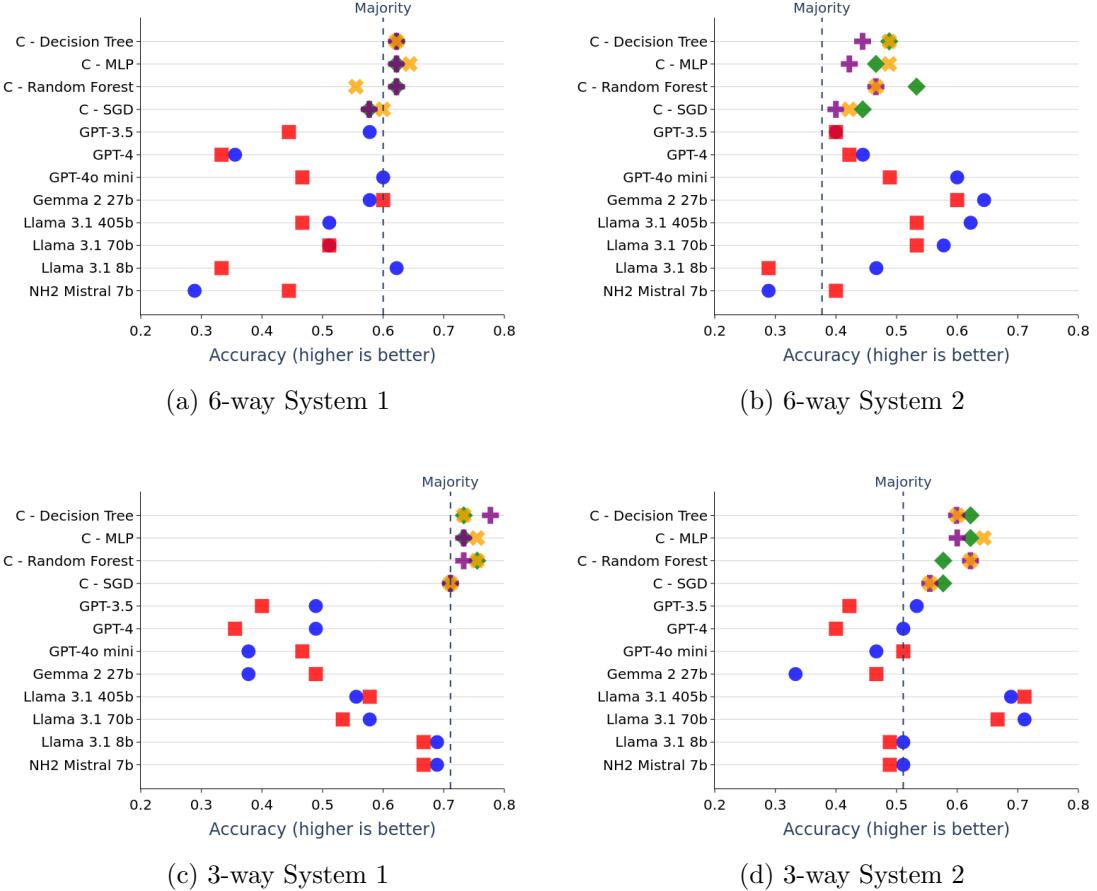


Figure 10: Human gold label prediction accuracy of classical machine learning models using various input schemes: base features only (green diamond \blacklozenge), base features with count vectorization of premise and conclusion (purple cross \blacksquare), and base features with TFIDF encoding of premise and conclusion (orange X \blacksquare). LLMs are also included with base prompting (blue circle \bullet) and personality prompting with equal weights (red square \blacksquare). Dashed line represents the majority baseline for the classification task.

challenges. First, NLI has classically been modeled as a classification task, and efforts to convert it to a regression problem (Chen et al., 2020) have faced significant difficulties (Nighojkar et al., 2023). Second, although our 6-way scale resembles a Likert scale, neither our 3-way nor our 6-way scale can be directly interpreted as a regression task because there is no intermediate value between 1 and 2 (or A and B) on our 6-way scale. The debate on whether Likert-like scales should be treated as regression is discussed in Sullivan and Artino (2013).

To process human responses for the 6-way classification scheme, we convert labels ranging from A to F into a categorical format corresponding to numbers 1 through 6. Similarly, for the 3-way scheme, the labels contradiction, neutral, and entailment are mapped to 1,

2, and 3, respectively. We use four classical machine learning algorithms: Stochastic Gradient Descent (SGD) (Gupta et al., 2014), Decision Tree (Quinlan, 1986), Random Forest (Tin Kam Ho, 1998), and a Feed-forward Neural Network (Schmidhuber, 2015). Various statistics and features that might influence the predictability of responses are added to the dataset, including (1) the word count of premises and conclusions, (2) the perplexity of premises and conclusions, (3) the BLEU score between premise and conclusion, and (4) sentiment analysis of premises and conclusions using the VADER tool. Collectively, these are referred to as “base input features.”

Given that these algorithms cannot natively process raw text data, we tokenize the text to remove stopwords and punctuation, then encode the premises and conclusions using Bag-of-Words and TF-IDF (Robertson, 2004), converting the text into numerical representations. The choice of encoding is crucial, as it impacts the model’s ability to capture linguistic nuances. Due to the limited data available, which restricts the size of the test set, we employ grid search to optimize these algorithms. The models are trained on a simple classification task, where they predict the gold label for human responses based on the input features, separately for both System 1 and System 2 responses. The performance of the best hyperparameter combinations, averaged across 5-fold cross-validation, is presented in Figure 10. For LLMs, we rely on the majority vote from ten base prompts and the majority vote from ten personality prompts for each item.

This task primarily aims to predict the population-level gold label, and while it may oversimplify individual reasoning differences, addressing whether human responses are predictable is a critical step before conducting more sophisticated analyses. As anticipated, the accuracy for the 3-way scheme (Figures 10c and 10d) surpasses that of the 6-way scheme (Figures 10a and 10b). The majority vote baseline, which simply predicts the most common gold label in the dataset, achieves higher accuracy on System 1 for both schemes; however, classical ML algorithms only slightly outperform this baseline. The baseline accuracy significantly declines for System 2 in both the 6-way and 3-way schemes, suggesting greater diversity in human responses for System 2 compared to System 1. Despite this decline, the classical ML algorithms maintain performance levels close to those observed in System 1, indicating a consistent degree of predictability in both System 1 and System 2 responses, with predictability being higher in System 1. Notably, Figure 10 reveals that LLMs perform nearly equally well in predicting gold labels for System 1 and System 2 responses, and they significantly surpass the baseline in System 2, emerging as the most effective at predicting the label most agreed upon by humans using System 2 reasoning. The fact that LLMs perform similarly on the 6-way scheme as they do on the 3-way scheme is also intriguing. Despite fewer choices in the 3-way scheme, which theoretically should make it easier, the LLMs’ performance suggests that the increased granularity of the 6-way scheme may render human responses more predictable for LLMs than initially expected.

4.2 Is the variance in human responses predictable?

After determining that we can predict the human gold label on both System 1 and System 2 with accuracy exceeding the majority baseline, we proceed to predict the variance among all human responses for a given item. This variance is calculated by mapping the human responses onto a scale from 1 to 6, as described in Section 4.1. It is important to note that,

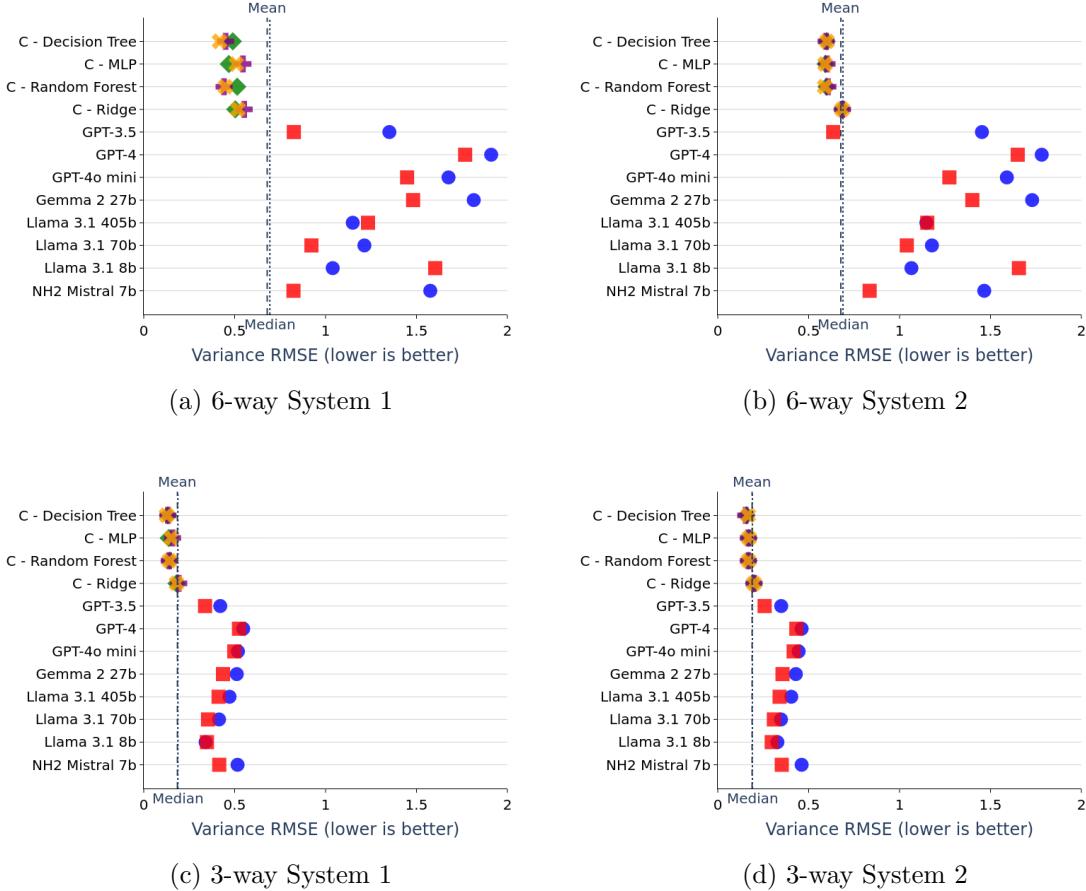


Figure 11: RMSE for predicting variance between human responses using classical machine learning models with various input schemes: base features only (green diamond \blacklozenge), base features with count vectorization of premise and conclusion (purple cross \blacksquare), and base features with TFIDF encoding of premise and conclusion (orange X \blacksquare). LLMs are also included with base prompting (blue circle \bullet) and personality prompting with equal weights (red square \blacksquare). Dashed line represents the mean baseline and dotted line represents the median baseline for the regression task.

for this task, we treat the responses as ordinal rather than nominal, since the objective is to predict the variance between the human responses. Given that we are treating this as a regression task, we use root mean squared error (RMSE) as the evaluation metric, where a lower value indicates better performance (Armstrong & Collopy, 1992). As in Section 4.1, we employ four different machine learning algorithms: ridge regression (Hoerl & Kennard, 2000), decision tree, random forest, and a feed-forward neural network. We conduct a grid search to optimize hyperparameters and report the average results of a 5-fold cross-validation, as shown in Figure 11. For LLMs, we calculate the variance from ten votes using base prompting and ten votes using personality prompting to compare with the

human variance for each item. The baselines in this analysis predict either the mean or the median variance across all items.

Since the 3-way scheme offers fewer options, the variance values fall within a narrower range, leading to lower prediction errors compared to the 6-way scheme, which is consistent with expectations. In both schemes, traditional ML algorithms outperform the baselines, although the margin of improvement is minimal. For both schemes, the variance between LLM votes significantly differs from the variance between human responses, with this disparity being much larger than what was observed in the accuracy of predicting the gold label (Figure 10). This indicates that predicting variance is a highly challenging task for LLMs, especially without any fine-tuning.

4.3 Can AI mimic the entire human response distribution?

Given the contrasting results obtained from our initial two experiments, further investigation is required to evaluate the ability of LLMs to emulate human reasoning. Additionally, while maxima and variance provide insight into a distribution, they represent only two of its characteristics. To more comprehensively evaluate the similarity between the response distributions of LLMs and humans, we now focus on comparing the entire response distribution. The Wasserstein Distance (Dobrushin, 1970), also known as Earth Mover’s Distance (EMD), serves as a metric for quantifying the difference between probability distributions across a specified metric space — in this instance, the set of labels in a 3-way or 6-way classification scheme. Conceptually, if each distribution is visualized as a unit mass of earth, the EMD reflects the minimal *cost* required to transform one distribution into the other, considering both the amount of earth that needs to be moved and the mean distance it must be moved. This makes EMD particularly suitable for our case, as it is sensitive to the ordinal nature of the metric.

Since EMD requires a probability distribution, we transform each set of human responses into a vector of size k (three or six for the 3-way or 6-way scheme), where each entry represents the frequency of the corresponding label. We then normalize this vector to create a probability distribution. Notably, EMD has an unbounded range, so we convert it into a similarity measure ranging from 0 to 1, which we term Earth Mover’s Similarity (EMS). Our similarity function, defined as:

$$EMS(D_1, D_2) = 100^{EMD(D_1, D_2)}$$

takes the two normalized probability distributions D_1 and D_2 as inputs. While e is often used as the base for exponents, we chose a base of 100 in this study to better distinguish variations, as using e would compress the values of interest into a narrow range.

In our comparison, we previously evaluated LLMs against traditional ML algorithms trained on our dataset. However, these algorithms do not generate a distribution of votes. In this experiment, we assign weights to each of the ten personality prompts and use a genetic algorithm to identify the optimal set of weights for each of the five folds (as in Experiments 4.1 and 4.2). For each fold, we again conduct a grid search to find the best hyperparameters (parent selection type and crossover type) for training the genetic algorithm, while fixing the number of generations at 8, the population per generation at 256, and the number of mating parents at 128. These parameters were chosen to balance computational costs

with performance, though we did not do a rigorous comparison of all possible values due to computational cost limitations. We refer to this setup as personality prompting with a genetic algorithm, or P-GA. Additionally, we test a setup using personality prompting with equal weights, termed P-EQ.

For a more thorough comparison, we compute three baselines. The *uniform baseline* (blue box at the bottom of all plots in Figure 12) distributes probabilities equally across the label space (all six labels for 6-way and all three labels for 3-way). The *normal distribution baseline* (orange box at the bottom of all plots in Figure 12) distributes probabilities normally around the mean of human votes, using the standard deviation of the human votes. The *best ML models* baseline (dotted line in all plots in Figure 12) distributes probabilities around the mean predicted by the best ML model for the gold labels, with the standard deviation predicted by the best ML model for the variance.

Figure 12 presents box plots of EMS between human and model votes for each of the five folds’ test sets. It is evident that LLMs significantly outperform the best ML models in mimicking the full distribution for both System 1 and System 2 in both the 3-way and 6-way schemes. The normal distribution baseline performs better on the 6-way scheme than on the 3-way scheme, possibly because the 3-way scheme has fewer labels and thus fewer data points from which to infer a normal probability distribution. For the 6-way scheme, personality prompting consistently outperforms or matches base prompting (markedly so for some architectures), with P-GA slightly outperforming P-EQ. The 3-way scheme is the only case where this trend does not hold, suggesting that personality prompting is more effective than base prompting at replicating the full human distribution, especially with the increased granularity of the 6-way scheme.

Another notable trend is that no GPT architecture outperforms the best-performing open-source LLM. Most GPT models have EMS values comparable to or even lower than those of the best ML models in the 3-way scheme (Figures 12c and 12d). Mistral is the best model on 6-way System 1 (Figure 12a) with Llama being a close second; however, Llama surpasses Mistral on 6-way System 2 (Figure 12a). Llama models also perform best on the 3-way scheme (Figures 12c and 12d). Although Gemma never emerges as the top-performing model, it outperforms GPT-4 in the 6-way scheme (Figures 12a and 12b) and nearly all GPT models in the 3-way scheme (Figures 12c and 12d). The only exception is GPT-3.5, which surpasses other, newer GPT models in the 3-way scheme (Figures 12c and 12d) but is outperformed by GPT-4o-mini in the 6-way scheme (Figures 12a and 12b).

Data Availability Statement All our code and data has been attached to this submission and this GitHub repository will be made public upon acceptance to make it available to the larger audience.

5 Conclusion

This study advances the field of AI reasoning by shifting the focus from merely predicting a single answer chosen by the majority of humans to modeling the full distribution of human reasoning responses. Predicting this distribution—a challenge we call the “full spectrum reasoning problem”—is not just a more comprehensive approach; it is essential for developing AI systems that can more accurately reflect the diversity and complexity of human thought. This is particularly important in applications where understanding the range of

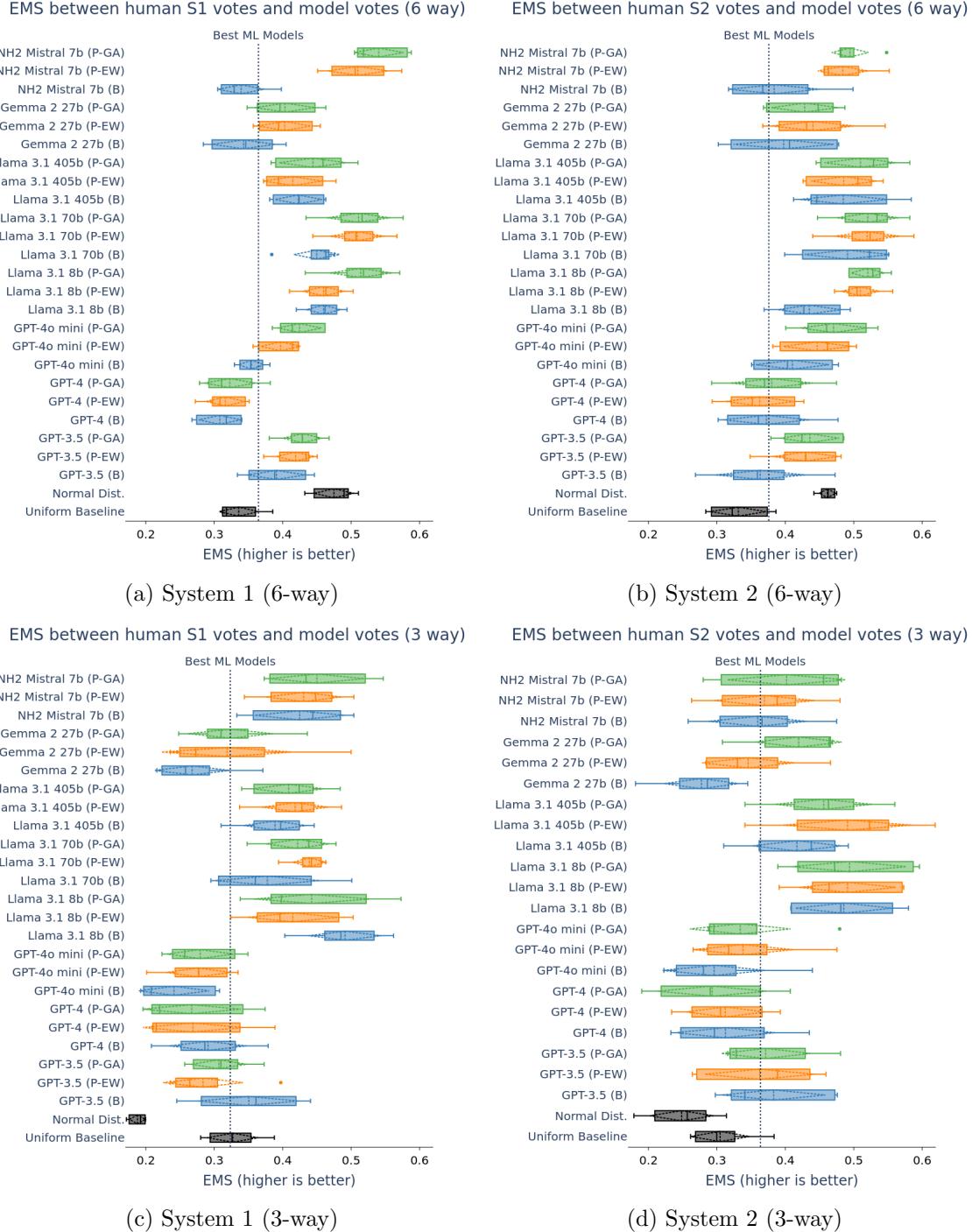


Figure 12: Earth Mover’s Similarity (EMS) between human votes and model votes. ‘(B)’ indicates base prompting (blue boxes), ‘(P-EW)’ indicates personality prompting with equal weights (orange boxes) and ‘(P-GA)’ indicates personality prompting with the genetic algorithm to fine-tune the weights of each prompt (green boxes).

human reasoning, rather than simply selecting the most common response, is critical for tasks such as decision-making (Lai et al., 2023; Stone et al., 2020), personalized interactions (Araujo & Bol, 2024; Fitzpatrick et al., 2017), and ethical AI deployment (Mittelstadt, 2019).

Our findings with personality prompting are especially promising, demonstrating that this approach can more closely align LLM outputs with human reasoning distributions while maintaining computational efficiency. Unlike base prompting, which requires multiple independent prompts to approximate a distribution, personality prompting allows for a more targeted generation of responses, using resources equivalent to base prompting but yielding more human-like results. The efficiency of our P-GA method, which leverages genetic algorithms to optimize personality prompts without the need for GPU resources, further underscores the practicality of our approach. Further investigation into different prompting styles (Liu et al., 2023; Sahoo et al., 2024) is necessary to assess how they compare to personality prompts.

The broader implications of this work are significant. By providing a framework to evaluate the *human-ness* of reasoning in both closed-source (Lieber et al., 2024; Templeton, 2024) and open-source (Almazrouei et al., 2023a; Chiang et al., 2023; Workshop et al., 2023) LLMs, our study enables a more nuanced assessment of AI models. Our finding that open-source models like Llama and Mistral exhibit reasoning patterns that align more with human cognition (at least with respect to the tasks we present here) than larger, proprietary models like GPT-4 challenges the assumption that larger parameter counts necessarily lead to more human-like AI. This insight is valuable not only for AI researchers but also for developers seeking to create more relatable and effective AI systems.

Additionally, our results suggest intriguing avenues for future research. For instance, the participants generally refrained from choosing the middle two options (C or D) corresponding to neutral for all question types in both phases of our survey. This suggests an underlying bias that further research can confirm. The observed correlation between participants' System 1 and System 2 responses in Figure 8 opens up the possibility of refining LLMs to better model individual cognitive processes. This could involve tailoring personality prompt weights to specific users, enhancing the accuracy and personalization of AI reasoning through techniques such as *hyperparameter hypothesization* (Nighojkar, 2024; Nighojkar et al., 2022).

While our experimental design offers robustness by being independent of publicly available datasets, it also presents a limitation: the findings are contingent on the specific questions used. Nevertheless, we believe our experimental setup is more robust than those previously observed in this field of research (Hagendorff et al., 2023) and can be scaled to larger datasets with relative ease, providing a strong foundation for further research in this area.

In summary, this work represents a significant step towards AI that more faithfully emulates human reasoning, emphasizing the importance of predicting the full distribution of human responses rather than relying on simplified models of accuracy. This approach not only enriches our understanding of AI-human interaction but also lays the groundwork for more sophisticated and human-like AI systems in the future.

Acknowledgments and Disclosure of Funding

This material is based on work supported by the National Science Foundation under Grant No. 2311286.

Appendix A. Examples of Different Item Types

Item Type	Premise	Conclusion
StereoNLI	Meskerem, an Ethiopian mother of two, works as a software developer in the bustling city of Addis Ababa. In her free time, she enjoys reading passages from the Bible to enrich her spiritual life and connect with her roots. Despite her busy schedule, Meskerem diligently balances her professional career and family responsibilities, making her a respected figure in both her community and the tech industry.	Meskerem speaks Amharic. Meskerem is an atheist.
Fallacy	The online seminar was about to start when the power went out. Allowing freedom of speech practically guarantees that misuse of that freedom will occur, given how thin the line is between freedom of speech and its misuse.	Meskerem is skilled in repairing pipes and fixtures in residential homes. The online seminar couldn't start because of the power outage.
	People who support cryptocurrency frequently talk about blockchain technology.	We must avoid freedom of speech.
Stereo Syllogism	No chess player is immune to the intricacies and challenges that the game presents. All Tom ever wanted to be is a chess player.	Traditional banking systems are not a concern for them.
E Guilt	One chilly evening, a hit-and-run occurred outside a bustling downtown bar known for its live music. Around 11 PM, as patrons milled outside, a speeding dark-colored sedan struck a young man stepping onto the street and sped away without stopping. The impact launched the victim across the pavement, causing immediate chaos as bystanders called 911 and rushed to aid him. The victim sustained significant injuries and was quickly hospitalized. Police cordoned off the scene, gathered evidence, and interviewed witnesses, but varying accounts complicated the identification of the assailant. The main suspect, Jack, known for his traffic violations and criminal past, is a local familiar with the area.	Tom's single wish was to confront the complexities and difficulties posed by the game of chess. Jack is guilty of the crime.

C Guilt	One chilly evening, a hit-and-run occurred outside a bustling downtown bar known for its live music. Around 11 PM, as patrons milled outside, a speeding dark-colored sedan struck a young man stepping onto the street and sped away without stopping. The impact launched the victim across the pavement, causing immediate chaos as bystanders called 911 and rushed to aid him. The victim sustained significant injuries and was quickly hospitalized. Police cordoned off the scene, gathered evidence, and interviewed witnesses, but varying accounts complicated the identification of the assailant. The main suspect, Jack, has a clean driving record and was not familiar with the area.	Jack is guilty of the crime.
P Primacy / Recency	Simon engages in a lot of research activities. He does not publish academically very often. He has a lot of free time to play sports. He does not keep up with the latest advances in his field. Finally, he spends a lot of time mentoring his students.	Simon is a professor.
N Primacy / Recency	Aleshia is very engaged in learning. She mostly is engaged in self-directed learning. She is financially independent. She is a manager at a very successful Fortune 500 company. Finally, she is actively involved in campus life and a member of many student clubs.	Aleshia is not a college student.

Table 2: Examples of different item types from our AI-generated dataset.

Appendix B. Prompts for Dataset Creation

Here we outline all our prompts used to create items in our dataset for transparency. In the OpenAI API, a system prompt sets the initial instructions and context for guiding the model’s behavior, while a human prompt represents the user’s input or query during the conversation. All premises and conclusions are rephrased using the system prompt: “You are a writing assistant for a linguist named Steve. Steve comes to you because you are good at writing sentences that are understandable and grammatically correct. Steve will give you a sentence. You need to rephrase the sentence such that it is easier to understand. Preserving all information is not necessary, but preferred. Just write the sentence and nothing else before or after it.” and the human prompt: “Sentence: {sentence}”.

Step	Prompt type	Prompt
Premise creation	System	Come up with a name for X. Do not write the name separately, include it in the sentences you write. Write a paragraph of 3 sentences about X that people would generally assume. Don't praise X, instead try to state sentences that have a truth value.
	Human	Details about X: {gender}, {profession}, {race}, {religion}
Entailment conclusion creation	Human	Here's a paragraph about X: {premise} Detail about X: {detail}. Write one assumption based on the above detail about X that has a truth value. Do not use words like 'probably' or 'likely', just state the statement. Do not justify or explain the statement in any way. Use X's name.
Contradiction conclusion creation	Human	Here's a paragraph about X: {premise} Detail about X: {detail}. Write one assumption based on the above detail about X that must be false. Do not use words like 'probably' or 'likely', just state the statement. Do not justify or explain the statement in any way. Use X's name.
Neutral conclusion creation	Human	Here're two stories: STORY 1: {premise1} STORY 2: {premise2} Here's a statement about story 1: {conclusion1} Rewrite this statement so that the subject of it is the subject from story 2. Keep everything else the same. Do not write 'story 1' or 'story 2' anywhere. Write just the new statement and nothing else.

Table 3: Prompts to create StereoNLI items

Premise template	Conclusion template
[A] happened right before [B].	[A] caused [B].
Ever since [A] began, we've seen an increase in [B].	[A] is responsible for the rise in [B].
[B] has decreased since we started doing [A].	Implementing [A] is the reason [B] has decreased.
Every time [A] occurs, [B] follows soon after.	The occurrence of [A] directly leads to [B].
[A] has been on the rise. Meanwhile, [B] has been becoming more common.	The growth of [A] is promoting the spread of [B].
Ever since [A] stopped, [B] has started.	The absence of [A] is the trigger for [B].
[A] started, and shortly after, [B] was observed.	The onset of [A] brought about [B].
Whenever [A] is present, [B] seems to follow.	[A] sets the stage for [B] to take place.
We did not have [B] until [A] was introduced.	[A] is the root cause of [B].
Each instance of [A] precedes [B].	[A] is the driving force behind [B].

Table 4: Post hoc ergo propter hoc fallacy templates.

Premise template	Conclusion template
If we allow [A], then it's only a matter of time before [B] happens.	We should not allow [A].
The moment we start [A], we set a precedent for [B].	We can't risk starting [A].
Once you open the door to [A], it's impossible to prevent [B] from coming through.	We shouldn't open the door to [A].
Every time society has embraced [A], it has eventually led to [B].	Embracing [A] would be a grave mistake.
The line between [A] and [B] is so thin, allowing [A] practically guarantees [B] will occur.	It's imperative we avoid [A].
[A]'s very existence is a stepping stone to [B].	For our own safety, we must eliminate [A].
There's a domino effect at play. Once [A] is set into motion, [B] will inevitably follow.	We should halt [A] before it's too late.
Allowing [A] is like opening Pandora's box, leading directly to [B].	We dare not open that box by permitting [A].
History has shown that [A] can subtly pave the way for [B].	We must learn from history and resist [A].
[A] might seem harmless on its own, but it's the first step on a dangerous path to [B].	To prevent disaster, we must avoid [A].

Table 5: Slippery slope fallacy templates.

Premise template	Conclusion template
[Person A] believes in [Complex Idea].	[Person A] is basically saying [Oversimplified or Misrepresented Version of Complex Idea].
According to [Group or Person], [Specific Nuanced Position].	[Group or Person] thinks [Extreme or Unrelated Position].
[Person B] stated that [Specific Condition or Circumstance].	[Person B] wants [Exaggerated or Unrelated Outcome].
Advocates for [Cause or Movement] argue for [Particular Aspect of Cause or Movement].	They just want [Unrelated or Overly Simplified Goal].
[Person C] wrote an article about [Specific Topic].	[Person C] must believe [Generalized, Simplified, or Twisted Idea about Topic].
[Group or Person] supports [Specific Action or Policy].	They must hate [Unrelated Group or Thing].
[Person D] mentioned that [Specific Fact or Statistic].	[Person D] denies [Related but Not Equivalent Fact or Statistic].
Proponents of [Theory or Idea] often discuss [Specific Aspect of Theory or Idea].	They don't care about [Different or Opposing Aspect].
[Person E] criticized [Specific Part of a Broader Concept].	[Person E] is against [The Entire Broader Concept].

Table 6: Straw person fallacy templates.

Step	Prompt type	Prompt
Premise creation	System	You will get a template. Fill the part in brackets ([A] and [B]) with suitable phrases. Your answer should be of the form below and should not contain anything else (not even empty lines): A: xxx B: xxx
	Human	Template: {premise template}
Conclusion creation	System	You will get a template (Template 1) and your response (Filled Template 1) to that template. You will get another template to fill (Template 2). Fill the part in brackets(eg: [A] and [B]) with suitable phrases based on your previous response. Your response should not contain anything other than the filled template, (not even empty lines).
	Human	Template 1: {premise template} Filled Template 1: {premise} Template 2: {conclusion template}

Table 7: Prompts to create Fallacy items. The templates are in Tables 4-6.

Step	Prompt type	Prompt
Premise creation	System	You are a writing assistant for a linguist named Steve. Steve gives you a word and a template for a sentence. The template is a guideline for how to structure the sentence. Do not feel constrained by the template. It is just a guideline. Use your own creativity to write the sentences that are structured similarly to the template. Steve comes to you because you are good at writing sentences that are understandable and grammatically correct. Steve also likes your variety of vocabulary and your ability to write sentences varying in length and complexity. Steve gives you the following instructions: Write each sentence on a new line. Do not write a bulleted or a numbered list. Do not write any other text except the sentences. Do not write any punctuation marks except the period at the end of each sentence.
	Human	Steve wants you to write <code>{n_premises}</code> sentences, each of which follows the following conditions: 1. The sentence is about <code>{seed_word}</code> (singular/plural). 2. The sentence fits the template ' <code>{template}</code> ' 3. The word ' <code>{seed_word}</code> ' is used in the sentence in place of ' <code>{variable}</code> '. This means ' <code>{seed_word}</code> ' should come towards the <code>{beginning_or_end}</code> of the sentence. This is the most important condition. 4. The sentence is grammatically correct. 5. The sentence is understandable. 6. The sentence is not too long or too short. 7. The sentence is not too simple or too complex. 8. The sentence is not too similar to any of the other sentences you write.
Conclusion creation	System	You are a writing assistant for a linguist named Steve. Steve comes to you because you are good at writing sentences that are understandable and grammatically correct. Steve will give you a pair of sentences. You need to combine the sentences into one sentence. This combination needs to be done such that a given word 'seed word' is eliminated from the resulting sentence. Make sure the resulting sentence is short, easy to understand, and fluent. The truth value of the resulting sentence does not matter. Saving all the information from the original sentences is not important. Just write the sentence and nothing else before or after it.
	Human	Sentence Pair: <code>{minor_premise}</code> <code>{major_premise}</code> Seed word: <code>{seed_word}</code>

Table 8: Prompts to create Syllogism and Stereo-Syllogism items. Seed words are regular nouns for syllogism items and nouns from StereoSet for stereo-syllogism items.

References

- Ahmed, T., Bird, C., Devanbu, P., & Chakraborty, S. (2024, February). *Studying LLM Performance on Closed- and Open-source Data*. arXiv. Retrieved 2024-08-21, from <http://arxiv.org/abs/2402.15100> (arXiv:2402.15100 [cs])
- Almazrouei, E., Alobeidli, H., Alshamsi, A., Cappelli, A., Cojocaru, R., Debbah, M., ... Penedo, G. (2023a). Falcon-40B: an open large language model with state-of-the-art performance.
- Almazrouei, E., Alobeidli, H., Alshamsi, A., Cappelli, A., Cojocaru, R., Debbah, M., ... Penedo, G. (2023b, November). *The Falcon Series of Open Language Models*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2311.16867> (arXiv:2311.16867 [cs])
- Araujo, T., & Bol, N. (2024, January). From speaking like a person to being personal: The effects of personalized, regular interactions with conversational agents. *Computers in Human Behavior: Artificial Humans*, 2(1), 100030. Retrieved 2024-08-22, from <https://linkinghub.elsevier.com/retrieve/pii/S2949882123000300> doi: 10.1016/j.chbah.2023.100030
- Ariely, D. (2010). *Predictably irrational: the hidden forces that shape our decisions* (Revised and expanded edition, First Harper Perennial edition published ed.). New York: Harper Perennial.
- Aristotle. (2013). *The organon*. CreateSpace. (OCLC: 884933600)
- Armstrong, J., & Collopy, F. (1992, June). Error measures for generalizing about forecasting methods: Empirical comparisons. *International Journal of Forecasting*, 8(1), 69–80. Retrieved 2024-08-20, from <https://linkinghub.elsevier.com/retrieve/pii/016920709290008W> doi: 10.1016/0169-2070(92)90008-W
- Bago, B., & De Neys, W. (2017, January). Fast logic?: Examining the time course assumption of dual process theory. *Cognition*, 158, 90–109. Retrieved 2024-08-20, from <https://linkinghub.elsevier.com/retrieve/pii/S0010027716302542> doi: 10.1016/j.cognition.2016.10.014
- Balloccu, S., Schmidlová, P., Lango, M., & Dušek, O. (2024, February). *Leak, Cheat, Repeat: Data Contamination and Evaluation Malpractices in Closed-Source LLMs*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2402.03927> (arXiv:2402.03927 [cs])
- Bergius, M., Ernberg, E., Dahlman, C., & Sarwar, F. (2020, 07). Are judges influenced by legally irrelevant circumstances? *Law, Probability and Risk*. Retrieved from <https://doi.org/10.1093/lpr/mgaa008> (mgaa008) doi: 10.1093/lpr/mgaa008
- Boissin, E., Caparos, S., Voudouri, A., & Neys, W. D. (2022, July). Debiasing System 1: Training favours logical over stereotypical intuiting. *Judgment and Decision Making*, 17(4), 646–690. Retrieved 2023-11-09, from <https://www.cambridge.org/core/journals/judgment-and-decision-making/article/debiasing-system-1-training-favours-logical-over-stereotypical-intuiting/3CF9452965581BDFD9A697C5B9F2E6AD> doi: 10.1017/S1930297500008895
- Bowman, S. R., Angeli, G., Potts, C., & Manning, C. D. (2015). A large annotated corpus for learning natural language inference. *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 632–642. doi: 10.18653/v1/d15-1075

- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... Amodei, D. (2020). Language Models are Few-Shot Learners. *arXiv*. doi: 10.48550/arxiv.2005.14165
- Chemero, A. (2023, November). LLMs differ from human cognition because they are not embodied. *Nature Human Behaviour*, 7(11), 1828–1829. Retrieved 2024-08-20, from <https://www.nature.com/articles/s41562-023-01723-5> doi: 10.1038/s41562-023-01723-5
- Chen, T., Jiang, Z., Poliak, A., Sakaguchi, K., & Durme, B. V. (2020). *Uncertain Natural Language Inference* (Tech. Rep.). Retrieved from <http://nlp.jhu.edu/unli>.
- Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., ... Xing, E. P. (2023, March). *Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality.* Retrieved from <https://lmsys.org/blog/2023-03-30-vicuna/>
- Computer, T. (2023). *Redpajama: an open dataset for training large language models.* Retrieved from <https://github.com/togethercomputer/RedPajama-Data>
- Dagan, I., Glickman, O., & Magnini, B. (2006). The PASCAL Recognising Textual Entailment Challenge. In J. Quiñonero-Candela, I. Dagan, B. Magnini, & F. d'Alché Buc (Eds.), *Machine Learning Challenges. Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Tectual Entailment* (Vol. 3944, pp. 177–190). Berlin, Heidelberg: Springer Berlin Heidelberg. Retrieved 2024-09-01, from http://link.springer.com/10.1007/11736790_9 (Series Title: Lecture Notes in Computer Science) doi: 10.1007/11736790_9
- Da Silva, S. (2023, October). System 1 vs. System 2 Thinking. *Psych*, 5(4), 1057–1076. Retrieved 2024-08-15, from <https://www.mdpi.com/2624-8611/5/4/71> doi: 10.3390/psych5040071
- De Neys, W., & Pennycook, G. (2019, October). Logic, Fast and Slow: Advances in Dual-Process Theorizing. *Current Directions in Psychological Science*, 28(5), 503–509. Retrieved 2024-02-27, from <http://journals.sagepub.com/doi/10.1177/0963721419855658> doi: 10.1177/0963721419855658
- De Paoli, S. (2023, October). *Improved prompting and process for writing user personas with LLMs, using qualitative interviews: Capturing behaviour and personality traits of users.* arXiv. Retrieved 2024-08-15, from <http://arxiv.org/abs/2310.06391> (arXiv:2310.06391 [cs])
- Ding, N., Qin, Y., Yang, G., Wei, F., Yang, Z., Su, Y., ... Sun, M. (2023, March). Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3), 220–235. Retrieved 2024-09-02, from <https://www.nature.com/articles/s42256-023-00626-4> doi: 10.1038/s42256-023-00626-4
- Dobrushin, R. L. (1970, January). Prescribing a System of Random Variables by Conditional Distributions. *Theory of Probability & Its Applications*, 15(3), 458–486. Retrieved 2024-08-12, from <http://pubs.siam.org/doi/10.1137/1115049> doi: 10.1137/1115049
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., ... Zhao, Z. (2024, August). *The Llama 3 Herd of Models.* arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2407.21783> (arXiv:2407.21783 [cs])
- Evans, J. S. B. T. (2008). Dual-Processing Accounts of Reasoning, Judgment, and Social Cognition. *Annual Review of Psychology*, 59(1), 255–278. doi: 10.1146/annurev.psych

.59.103006.093629

- Evans, J. S. B. T. (2019). Reflections on reflection: the nature and function of type 2 processes in dual-process theories of reasoning. *Thinking & Reasoning*, 25(4), 383–415. doi: 10.1080/13546783.2019.1623071
- Evans, J. S. B. T., Barston, J. L., & Pollard, P. (1983). On the conflict between logic and belief in syllogistic reasoning. *Memory & Cognition*, 11(3), 295–306. doi: 10.3758/bf03196976
- Evans, J. S. B. T., & Stanovich, K. E. (2013, May). Dual-Process Theories of Higher Cognition: Advancing the Debate. *Perspectives on Psychological Science*, 8(3), 223–241. Retrieved 2024-08-15, from <http://journals.sagepub.com/doi/10.1177/1745691612460685> doi: 10.1177/1745691612460685
- Evans, R., Saxton, D., Amos, D., Kohli, P., & Grefenstette, E. (2018). Can neural networks understand logical entailment? In *International conference on learning representations*. Retrieved from <https://openreview.net/forum?id=SkZxCk-0Z>
- Fearnside, W. W., & Holther, W. B. (1959). *Fallacy: The Counterfeit of Argument*. Englewood Cliffs, NJ: Prentice Hall.
- Fitzpatrick, K. K., Darcy, A., & Vierhile, M. (2017, June). Delivering Cognitive Behavior Therapy to Young Adults With Symptoms of Depression and Anxiety Using a Fully Automated Conversational Agent (Woebot): A Randomized Controlled Trial. *JMIR Mental Health*, 4(2), e19. Retrieved 2024-08-22, from <http://mental.jmir.org/2017/2/e19/> doi: 10.2196/mental.7785
- Frederick, S. (2005, November). Cognitive Reflection and Decision Making. *Journal of Economic Perspectives*, 19(4), 25–42. Retrieved 2024-04-10, from <https://pubs.aeaweb.org/doi/10.1257/089533005775196732> doi: 10.1257/089533005775196732
- Geeraert, N. (2013, September). When Suppressing One Stereotype Leads to Rebound of Another: On the Procedural Nature of Stereotype Rebound. *Personality and Social Psychology Bulletin*, 39(9), 1173–1183. Retrieved 2024-08-15, from <http://journals.sagepub.com/doi/10.1177/0146167213493121> doi: 10.1177/0146167213493121
- Gilovich, T., Griffin, D., & Kahneman, D. (Eds.). (2002). *Heuristics and Biases: The Psychology of Intuitive Judgment* (1st ed.). Cambridge University Press. Retrieved 2024-09-01, from <https://www.cambridge.org/core/product/identifier/9780511808098/type/book> doi: 10.1017/CBO9780511808098
- Gladwell, M. (2007). *Blink: the power of thinking without thinking* (1st Back Bay trade pbk. ed ed.). New York: Back Bay Books.
- Gu, H., Deguchi, C., Genç, U., Chandrasegaran, S., & Verma, H. (2023, October). *On the Effectiveness of Creating Conversational Agent Personalities Through Prompting*. arXiv. Retrieved 2024-08-15, from <http://arxiv.org/abs/2310.11182> (arXiv:2310.11182 [cs])
- Gupta, M. R., Bengio, S., & Weston, J. (2014). Training highly multiclass classifiers. *Journal of Machine Learning Research*, 15(43), 1461–1492. Retrieved from <http://jmlr.org/papers/v15/gupta14a.html>
- Hagedorn, T., Bone, M., Kruse, B., Grosse, I., & Blackburn, M. (2020, March). Knowledge Representation with Ontologies and Semantic Web Technologies to Promote Augmented and Artificial Intelligence in Systems Engineering. *INSIGHT*, 23(1), 15–20. Retrieved 2024-08-20, from <https://oncose.onlinelibrary.wiley.com/doi/>

- 10.1002/inst.12279 doi: 10.1002/inst.12279
- Hagendorff, T., Fabi, S., & Kosinski, M. (2023). Human-like intuitive behavior and reasoning biases emerged in large language models but disappeared in ChatGPT. *Nature Computational Science*, 1–6. doi: 10.1038/s43588-023-00527-x
- Hamade, K., McIlroy-Young, R., Sen, S., Kleinberg, J., & Anderson, A. (2024, May). Designing skill-compatible AI: Methodologies and frameworks in chess. In *ICLR 2024*. Retrieved from <https://www.microsoft.com/en-us/research/publication/designing-skill-compatible-ai-methodologies-and-frameworks-in-chess/>
- Hoerl, A. E., & Kennard, R. W. (2000, February). Ridge Regression: Biased Estimation for Nonorthogonal Problems. *Technometrics*, 42(1), 80–86. Retrieved 2024-08-20, from <http://www.tandfonline.com/doi/abs/10.1080/00401706.2000.10485983> doi: 10.1080/00401706.2000.10485983
- Hutto, C., & Gilbert, E. (2014, May). VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. *Proceedings of the International AAAI Conference on Web and Social Media*, 8(1), 216–225. Retrieved 2024-08-20, from <https://ojs.aaai.org/index.php/ICWSM/article/view/14550> doi: 10.1609/icwsm.v8i1.14550
- Jelinek, F., Mercer, R. L., Bahl, L. R., & Baker, J. K. (1977, December). Perplexity—a measure of the difficulty of speech recognition tasks. *The Journal of the Acoustical Society of America*, 62(S1), S63–S63. Retrieved 2024-09-02, from <https://pubs.aip.org/jasa/article/62/S1/S63/642598/Perplexity-a-measure-of-the-difficulty-of-speech> doi: 10.1121/1.2016299
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. l., ... Sayed, W. E. (2023, October). *Mistral 7B*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2310.06825> (arXiv:2310.06825 [cs])
- Kahneman, D. (2011). *Thinking, Fast and Slow*. Farrar, Straus and Giroux.
- Kahneman, D., & Tversky, A. (1984, April). Choices, values, and frames. *American Psychologist*, 39(4), 341–350. Retrieved 2024-08-22, from <http://doi.apa.org/getdoi.cfm?doi=10.1037/0003-066X.39.4.341> doi: 10.1037/0003-066X.39.4.341
- Kamruzzaman, M., & Kim, G. L. (2024, June). *Prompting Techniques for Reducing Social Bias in LLMs through System 1 and System 2 Cognitive Processes*. arXiv. Retrieved 2024-08-21, from <http://arxiv.org/abs/2404.17218> (arXiv:2404.17218 [cs])
- Kassin, S. M., Dror, I. E., & Kukucka, J. (2013). The forensic confirmation bias: Problems, perspectives, and proposed solutions. *Journal of Applied Research in Memory and Cognition*, 2(1), 42 - 52. Retrieved from <http://www.sciencedirect.com/science/article/pii/S2211368113000028> doi: <https://doi.org/10.1016/j.jarmac.2013.01.001>
- Kazemeini, A., Roy, S. S., Mercer, R. E., & Cambria, E. (2021, December). Interpretable Representation Learning for Personality Detection. In *2021 International Conference on Data Mining Workshops (ICDMW)* (pp. 158–165). Auckland, New Zealand: IEEE. Retrieved 2024-09-01, from <https://ieeexplore.ieee.org/document/9679950/> doi: 10.1109/ICDMW53433.2021.00026
- Kendall, M. G. (1938, June). A New Measure of Rank Correlation. *Biometrika*, 30(1-2), 81–93. Retrieved 2024-08-20, from <https://academic.oup.com/biomet/article-lookup/doi/10.1093/biomet/30.1-2.81> doi: 10.1093/biomet/30.1-2.81

- Khemlani, S., & Johnson-Laird, P. N. (2012). Theories of the syllogism: A meta-analysis. *Psychological Bulletin*, 138(3), 427–457. Retrieved 2024-08-15, from <https://doi.org/doi/10.1037/a0026841> doi: 10.1037/a0026841
- Klein, G. (2017). *Sources of Power: How People Make Decisions*. The MIT Press. Retrieved 2024-08-22, from <https://direct.mit.edu/books/book/3647/Sources-of-PowerHow-People-Make-Decisions> doi: 10.7551/mitpress/11307.001.0001
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., & Iwasawa, Y. (2022). Large Language Models are Zero-Shot Reasoners. *arXiv*. doi: 10.48550/arxiv.2205.11916
- Lai, V., Chen, C., Smith-Renner, A., Liao, Q. V., & Tan, C. (2023, June). Towards a Science of Human-AI Decision Making: An Overview of Design Space in Empirical Human-Subject Studies. In *2023 ACM Conference on Fairness, Accountability, and Transparency* (pp. 1369–1385). Chicago IL USA: ACM. Retrieved 2024-08-22, from <https://dl.acm.org/doi/10.1145/3593013.3594087> doi: 10.1145/3593013.3594087
- Laverghetta Jr., A., & Licato, J. (2022). Developmental negation processing in transformer language models. In *Proceedings from acl 2022*.
- Laverghetta Jr., A., & Licato, J. (2023a). Automatic generation of cognitive test items using large language models. In *Proceedings of the 2023 international meeting of the psychometric society (imps)*.
- Laverghetta Jr., A., & Licato, J. (2023b). Generating better items for cognitive assessments using large language models. In *Proceedings of the acl 18th workshop on innovative use of nlp for building educational applications*.
- Laverghetta Jr., A., Nighojkar, A., Mirzakhalov, J., & Licato, J. (2021). Can Transformer Language Models Predict Psychometric Properties? In *Proceedings of the 10th joint conference on lexical and computational semantics (*sem 2021)*. Bangkok, Thailand.
- Laverghetta Jr., A., Nighojkar, A., Mirzakhalov, J., & Licato, J. (2022). Predicting human psychometric properties using computational language models. In M. Wilberg, D. Molenaar, J. González, J.-S. Kim, & H. Hwang (Eds.), *Quantitative psychology* (Vol. Forthcoming). Springer.
- Lieber, O., Lenz, B., Bata, H., Cohen, G., Osin, J., Dalmedigos, I., ... Shoham, Y. (2024, July). *Jamba: A Hybrid Transformer-Mamba Language Model*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2403.19887> (arXiv:2403.19887 [cs])
- Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., & Neubig, G. (2023, September). Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. *ACM Computing Surveys*, 55(9), 1–35. Retrieved 2024-08-20, from <https://dl.acm.org/doi/10.1145/3560815> doi: 10.1145/3560815
- Ma, H., Zhang, C., Bian, Y., Liu, L., Zhang, Z., Zhao, P., ... Wu, B. (2023, March). *Fairness-guided Few-shot Prompting for Large Language Models*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2303.13217> (arXiv:2303.13217 [cs])
- Maclure, J. (2021, September). AI, Explainability and Public Reason: The Argument from the Limitations of the Human Mind. *Minds and Machines*, 31(3), 421–438. Retrieved 2024-08-20, from <https://link.springer.com/10.1007/s11023-021-09570-x> doi: 10.1007/s11023-021-09570-x
- Meissner, J. M., Thumwanit, N., Sugawara, S., & Aizawa, A. (2021). Embracing Ambiguity: Shifting the Training Target of NLI Models. *Proceedings of the 59th Annual Meeting*

- of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, 862–869. doi: 10.18653/v1/2021.acl-short.109
- Mitchell, M. (1996). *An introduction to genetic algorithms*. Cambridge, Mass.: MIT Press. (OCLC: 42854439)
- Mittelstadt, B. (2019, November). Principles alone cannot guarantee ethical AI. *Nature Machine Intelligence*, 1(11), 501–507. Retrieved 2024-08-22, from <https://www.nature.com/articles/s42256-019-0114-4> doi: 10.1038/s42256-019-0114-4
- Murdock, B. B. (1962, November). The serial position effect of free recall. *Journal of Experimental Psychology*, 64(5), 482–488. Retrieved 2024-08-20, from <http://doi.apa.org/getdoi.cfm?doi=10.1037/h0045106> doi: 10.1037/h0045106
- Nadeem, M., Bethke, A., & Reddy, S. (2021). StereoSet: Measuring stereotypical bias in pretrained language models. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 5356–5371. doi: 10.18653/v1/2021.acl-long.416
- Nie, Y., Williams, A., Dinan, E., Bansal, M., Weston, J., & Kiela, D. (2020). Adversarial NLI: A New Benchmark for Natural Language Understanding. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 4885–4901. doi: 10.18653/v1/2020.acl-main.441
- Nie, Y., Zhou, X., & Bansal, M. (2020). What Can We Learn from Collective Human Opinions on Natural Language Inference Data? *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 9131–9143. doi: 10.18653/v1/2020.emnlp-main.734
- Nighojkar, A. (2024). *An Inference-Centric Approach to Natural Language Processing and Cognitive Modeling* (Ph.D., University of South Florida, United States – Florida). Retrieved from <https://www.proquest.com/dissertations-theses/inference-centric-approach-natural-language/docview/3084657113/se-2?accountid=14745> (ISBN: 9798383569337 Publication Title: ProQuest Dissertations and Theses 31334853)
- Nighojkar, A., Jr, A. L., & Licato, J. (2023). No Strong Feelings One Way or Another: Re-operationalizing Neutrality in Natural Language Inference. *Proceedings of the 17th Linguistic Annotation Workshop (LAW-XVII)*, 199–210. doi: 10.18653/v1/2023.law-1.20
- Nighojkar, A., Khlyzova, A., & Licato, J. (2022, August). *Cognitive Modeling of Semantic Fluency Using Transformers*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2208.09719> [cs]
- Nighojkar, A., & Licato, J. (2021a). Improving paraphrase detection with the adversarial paraphrasing task..
- Nighojkar, A., & Licato, J. (2021b). Mutual Implication as a Measure of Textual Equivalence. *The International FLAIRS Conference Proceedings*, 34(1). doi: 10.32473/flairs.v34i1.128519
- Nisbett, R. E., & Wilson, T. D. (1977, May). Telling more than we can know: Verbal reports on mental processes. *Psychological Review*, 84(3), 231–259. Retrieved 2024-08-22, from <https://doi.apa.org/doi/10.1037/0033-295X.84.3.231> doi: 10.1037/

0033-295X.84.3.231

- Nye, M., Tessler, M., Tenenbaum, J., & Lake, B. M. (2021). Improving coherence and consistency in neural sequence models with dual-system, neuro-symbolic reasoning. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, & J. W. Vaughan (Eds.), *Advances in neural information processing systems* (Vol. 34, pp. 25192–25204). Curran Associates, Inc. Retrieved from <https://proceedings.neurips.cc/paper/2021/file/d3e2e8f631bd9336ed25b8162aef8782-Paper.pdf>
- O'Neil, C. (2017). *Weapons of math destruction: how big data increases inequality and threatens democracy* (First paperback edition ed.). New York: B/D/W/Y Broadway Books.
- OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., ... Zoph, B. (2023). *GPT-4 Technical Report*. arXiv. Retrieved 2024-08-12, from <https://arxiv.org/abs/2303.08774> (Version Number: 6) doi: 10.48550/ARXIV.2303.08774
- Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2001). BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics - ACL '02* (p. 311). Philadelphia, Pennsylvania: Association for Computational Linguistics. Retrieved 2024-08-15, from <http://portal.acm.org/citation.cfm?doid=1073083.1073135> doi: 10.3115/1073083.1073135
- Pavlick, E., & Kwiatkowski, T. (2019). Inherent Disagreements in Human Textual Inferences. *Transactions of the Association for Computational Linguistics*, 7, 677–694. doi: 10.1162/tacl_a_00293
- Peer, E., & Gamlie, E. (2013, 01). Heuristics and biases in judicial decisions. *Court Review*, 49, 114–118.
- Pellert, M., Lechner, C. M., Wagner, C., Rammstedt, B., & Strohmaier, M. (2024, January). AI Psychometrics: Assessing the Psychological Profiles of Large Language Models Through Psychometric Inventories. *Perspectives on Psychological Science*, 17456916231214460. Retrieved 2024-09-01, from <http://journals.sagepub.com/doi/10.1177/17456916231214460> doi: 10.1177/17456916231214460
- Penedo, G., Malartic, Q., Hesslow, D., Cojocaru, R., Cappelli, A., Alobeidli, H., ... Launay, J. (2023, June). *The RefinedWeb Dataset for Falcon LLM: Outperforming Curated Corpora with Web Data, and Web Data Only*. arXiv. Retrieved 2024-08-20, from <https://arxiv.org/abs/2306.01116> [cs]
- Prescott, T. J., & Wilson, S. P. (2023, May). Understanding brain functional architecture through robotics. *Science Robotics*, 8(78), eadg6014. Retrieved 2024-08-20, from <https://www.science.org/doi/10.1126/scirobotics.adg6014> doi: 10.1126/scirobotics.adg6014
- Quinlan, J. R. (1986, March). Induction of decision trees. *Machine Learning*, 1(1), 81–106. Retrieved 2024-08-20, from <http://link.springer.com/10.1007/BF00116251> doi: 10.1007/BF00116251
- Rachlinski, J. J., Wistrich, A. J., & Guthrie, C. (2015). Can judges make reliable numeric judgments? distorted damages and skewed sentences. *Indiana Law Journal*, 90.
- Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C. D., & Finn, C. (2023). *Direct Preference Optimization: Your Language Model is Secretly a Reward Model*. arXiv. Retrieved 2024-08-12, from <https://arxiv.org/abs/2305.18290> (Version

- Number: 3) doi: 10.48550/ARXIV.2305.18290
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... Liu, P. J. (2023, September). *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/1910.10683> [arXiv:1910.10683 [cs, stat]]
- Robertson, S. (2004, October). Understanding inverse document frequency: on theoretical arguments for IDF. *Journal of Documentation*, 60(5), 503–520. Retrieved 2024-08-22, from [https://www.emerald.com/insight/content/doi/10.1108/00220410410560582](https://www.emerald.com/insight/content/doi/10.1108/00220410410560582/full/html) doi: 10.1108/00220410410560582
- Roccas, S., Sagiv, L., Schwartz, S. H., & Knafo, A. (2002, June). The Big Five Personality Factors and Personal Values. *Personality and Social Psychology Bulletin*, 28(6), 789–801. Retrieved 2024-08-11, from <http://journals.sagepub.com/doi/10.1177/0146167202289008> doi: 10.1177/0146167202289008
- Saha, S., Prasad, A., Chen, J. C.-Y., Hase, P., Stengel-Eskin, E., & Bansal, M. (2024, July). *System-1.x: Learning to Balance Fast and Slow Planning with Language Models*. arXiv. Retrieved 2024-08-11, from <http://arxiv.org/abs/2407.14414> [arXiv:2407.14414 [cs]]
- Sahoo, P., Singh, A. K., Saha, S., Jain, V., Mondal, S., & Chadha, A. (2024, February). *A Systematic Survey of Prompt Engineering in Large Language Models: Techniques and Applications*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2402.07927> [arXiv:2402.07927 [cs]]
- Sainz, O., Campos, J., García-Ferrero, I., Etxaniz, J., De Lacalle, O. L., & Agirre, E. (2023). NLP Evaluation in trouble: On the Need to Measure LLM Data Contamination for each Benchmark. In *Findings of the Association for Computational Linguistics: EMNLP 2023* (pp. 10776–10787). Singapore: Association for Computational Linguistics. Retrieved 2024-09-01, from <https://aclanthology.org/2023.findings-emnlp.722> doi: 10.18653/v1/2023.findings-emnlp.722
- Sap, M., Swayamdipta, S., Vianna, L., Zhou, X., Choi, Y., & Smith, N. (2022). Annotators with Attitudes: How Annotator Beliefs And Identities Bias Toxic Language Detection. *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 5884–5906. doi: 10.18653/v1/2022.naacl-main.431
- Schmidhuber, J. (2015, January). Deep learning in neural networks: An overview. *Neural Networks*, 61, 85–117. Retrieved 2024-08-20, from <https://linkinghub.elsevier.com/retrieve/pii/S0893608014002135> doi: 10.1016/j.neunet.2014.09.003
- Stanovich, K. E., & West, R. F. (2000, October). Individual differences in reasoning: Implications for the rationality debate? *Behavioral and Brain Sciences*, 23(5), 645–665. Retrieved 2024-08-15, from https://www.cambridge.org/core/product/identifier/S0140525X00003435/type/journal_article doi: 10.1017/S0140525X00003435
- Stanovich, K. E., West, R. F., & Toplak, M. E. (2018). *The rationality quotient: toward a test of rational thinking* (First paperback edition ed.). Cambridge, Massachusetts London, England: The MIT Press.
- Stone, M., Aravopoulou, E., Ekinci, Y., Evans, G., Hobbs, M., Labib, A., ... Machtyngier, L. (2020, April). Artificial intelligence (AI) in strategic marketing decision-making: a research agenda. *The Bottom Line*, 33(2), 183–200. Retrieved 2024-08-

- 22, from <https://www.emerald.com/insight/content/doi/10.1108/BL-03-2020-0022/full/html> doi: 10.1108/BL-03-2020-0022
- Sullivan, G. M., & Artino, A. R. (2013, December). Analyzing and Interpreting Data From Likert-Type Scales. *Journal of Graduate Medical Education*, 5(4), 541–542. Retrieved 2024-08-22, from <https://meridian.allenpress.com/jgme/article/5/4/541/34037/Analyzing-and-Interpreting-Data-From-LikertType> doi: 10.4300/JGME-5-4-18
- Sun, J., Zheng, C., Xie, E., Liu, Z., Chu, R., Qiu, J., ... Li, Z. (2024, January). *A Survey of Reasoning with Foundation Models*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2312.11562> (arXiv:2312.11562 [cs])
- Team, G. (2024). Gemma. *Online*. Retrieved from <https://www.kaggle.com/m/3301> doi: 10.34740/KAGGLE/M/3301
- Teknium, theemozilla, karan4d, & huemin_art. (2024). *Nous Hermes 2 Mistral 7B DPO*. Retrieved from <https://huggingface.co/NousResearch/Nous-Hermes-2-Mistral-7B-DPO>
- Templeton, A. (2024). *Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet*. Anthropic.
- Tin Kam Ho. (1998, August). The random subspace method for constructing decision forests. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(8), 832–844. Retrieved 2024-08-20, from <http://ieeexplore.ieee.org/document/709601/> doi: 10.1109/34.709601
- Vatsal, S., & Dubey, H. (2024, July). *A Survey of Prompt Engineering Methods in Large Language Models for Different NLP Tasks*. arXiv. Retrieved 2024-08-15, from <http://arxiv.org/abs/2407.12994> (arXiv:2407.12994 [cs])
- Walton, D. (1985). *Arguer's position: A pragmatic study of ad hominem attack, criticism, refutation, and fallacy*. Greenwood Press.
- Walton, D., Reed, C., & Macagno, F. (2008). *Argumentation schemes*. Cambridge University Press.
- Walton, D. N. (1990). What is reasoning? what is an argument? *The Journal of Philosophy*, 87(8), 399–419. Retrieved from <http://www.jstor.org/stable/2026735>
- Wang, L., Xu, W., Lan, Y., Hu, Z., Lan, Y., Lee, R. K.-W., & Lim, E.-P. (2023). Plan-and-Solve Prompting: Improving Zero-Shot Chain-of-Thought Reasoning by Large Language Models. *arXiv*.
- Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., ... Zhou, D. (2023, March). *Self-Consistency Improves Chain of Thought Reasoning in Language Models*. arXiv. Retrieved 2024-08-15, from <http://arxiv.org/abs/2203.11171> (arXiv:2203.11171 [cs])
- Wason, P. C. (1966). Reasoning. In P. C. Wason (Ed.), *New horizons in psychology* (pp. 135–151). Penguin Books.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., ... Zhou, D. (2022). Chain of Thought Prompting Elicits Reasoning in Large Language Models. *arXiv*. doi: 10.48550/arxiv.2201.11903
- Wen, Z., Cao, J., Yang, Y., Yang, R., & Liu, S. (2024, March). Affective- NLI: Towards Accurate and Interpretable Personality Recognition in Conversation. In *2024 IEEE International Conference on Pervasive Computing and Communications (PerCom)* (pp.

- 184–193). Biarritz, France: IEEE. Retrieved 2024-09-01, from <https://ieeexplore.ieee.org/document/10494487/> doi: 10.1109/PerCom59722.2024.10494487
- Weston, J., & Sukhbaatar, S. (2023, November). *System 2 Attention (is something you might need too)*. arXiv. Retrieved 2023-11-21, from <http://arxiv.org/abs/2311.11829> (arXiv:2311.11829 [cs])
- Williams, A., Nangia, N., & Bowman, S. (2018). A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 1112–1122. doi: 10.18653/v1/n18-1101
- Williams, A., Thrush, T., & Kiela, D. (2022, February). ANLIzing the adversarial natural language inference dataset. In A. Ettinger, T. Hunter, & B. Prickett (Eds.), *Proceedings of the society for computation in linguistics 2022* (pp. 23–54). online: Association for Computational Linguistics. Retrieved from <https://aclanthology.org/2022.scil-1.3>
- Wistrich, A. J., Rachlinski, J. J., & Guthrie, C. (2015). Heart versus head: Do judges follow the law or follow their feelings?. *Texas Law Review*, 93(4), 855 - 923.
- Workshop, B., Scao, T. L., Fan, A., Akiki, C., Pavlick, E., Ilić, S., ... Wolf, T. (2023, June). *BLOOM: A 176B-Parameter Open-Access Multilingual Language Model*. arXiv. Retrieved 2024-08-20, from <http://arxiv.org/abs/2211.05100> (arXiv:2211.05100 [cs])
- Yang, J., Hu, Z., Zhu, D., & Nie, D. (2023). Belief bias, conflict detection, and logical complexity. *Current Psychology*, 1–9. doi: 10.1007/s12144-023-04562-9
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y., & Narasimhan, K. (2023). Tree of Thoughts: Deliberate Problem Solving with Large Language Models. *arXiv*. doi: 10.48550/arxiv.2305.10601
- Yu, P., Xu, J., Weston, J., & Kulikov, I. (2024, July). *Distilling System 2 into System 1*. arXiv. Retrieved 2024-08-11, from <http://arxiv.org/abs/2407.06023> (arXiv:2407.06023 [cs])
- Zdaniuk, B. (2014). Ordinary Least-Squares (OLS) Model. In A. C. Michalos (Ed.), *Encyclopedia of Quality of Life and Well-Being Research* (pp. 4515–4517). Dordrecht: Springer Netherlands. Retrieved 2024-08-20, from http://link.springer.com/10.1007/978-94-007-0753-5_2008 doi: 10.1007/978-94-007-0753-5_2008
- Zhang, S., Gong, C., & Choi, E. (2021, February). Capturing Label Distribution: A Case Study in NLI. *arXiv*. Retrieved from <http://arxiv.org/abs/2102.06859>
- Zhang, S., Rudinger, R., Duh, K., & Durme, B. V. (2017). Ordinal Common-sense Inference. *Transactions of the Association for Computational Linguistics*, 5, 379–395. doi: 10.1162/tacl_a_00068
- Zhang, X. F., & Marneffe, M.-C. d. (2021). Identifying inherent disagreement in natural language inference. *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4908–4915. doi: 10.18653/v1/2021.naacl-main.390
- Zhang, Y., Yang, J., Yuan, Y., & Yao, A. C.-C. (2024, April). *Cumulative Reasoning with Large Language Models*. arXiv. Retrieved 2024-08-15, from <http://arxiv.org/abs/2308.04371> (arXiv:2308.04371 [cs])

- Zhou, X., Nie, Y., & Bansal, M. (2022). Distributed NLI: Learning to Predict Human Opinion Distributions for Language Reasoning. In *Findings of the Association for Computational Linguistics: ACL 2022* (pp. 972–987). Dublin, Ireland: Association for Computational Linguistics. Retrieved 2024-08-20, from <https://aclanthology.org/2022.findings-acl.79> doi: 10.18653/v1/2022.findings-acl.79
- Zollman, D. A., Sirnoorkar, A., & Laverty, J. T. (2023, October). Analyzing AI and student responses through the lens of sensemaking and mechanistic reasoning. In *2023 Physics Education Research Conference Proceedings* (pp. 415–420). Sacramento, CA: American Association of Physics Teachers. Retrieved 2024-08-20, from <https://www.per-central.org/items/detail.cfm?ID=16619> doi: 10.1119/perc.2023.pr.Zollman