

# Beyond Keywords: Evaluating Large Language Model Classification of Nuanced Ableism

Naba Rizvi and Harper Strickland and Saleha Ahmed

and Aekta Kallepalli and Isha Khirwadkar and William Wu and Imani N. S. Munyaka

University of California, San Diego

La Jolla, CA 92093, USA

nrizvi@ucsd.edu

Nedjma Ousidhoum

Cardiff University

## Abstract

Large language models (LLMs) are increasingly used in decision-making tasks like résumé screening and content moderation, giving them the power to amplify or suppress certain perspectives. While previous research has identified disability-related biases in LLMs, little is known about how they conceptualize ableism or detect it in text. We evaluate the ability of four LLMs to identify nuanced ableism directed at autistic individuals. We examine the gap between their understanding of relevant terminology and their effectiveness in recognizing ableist content in context. Our results reveal that LLMs can identify autism-related language but often miss harmful or offensive connotations. Further, we conduct a qualitative comparison of human and LLM explanations. We find that LLMs tend to rely on surface-level keyword matching, leading to context misinterpretations, in contrast to human annotators who consider context, speaker identity, and potential impact. On the other hand, both LLMs and humans agree on the annotation scheme, suggesting that a binary classification is adequate for evaluating LLM performance, which is consistent with findings from prior studies involving human annotators.

**Trigger warning: this paper contains ableist language including explicit slurs and references to violence.**

## 1 Introduction

There is growing interest in using large language models (LLMs) to generate data that reflects human perspectives. However, there remains a significant gap in understanding *which* perspectives LLMs tend to emulate (Long et al., 2024; Rossi et al., 2024; Goyal and Mahmoud, 2025). While LLMs are known to reproduce human biases—particularly those related to disabilities—such biases often emerge in real-world applications, including re-

sume screening (Schramowski et al., 2022; Glazko et al., 2024).

Detecting anti-autistic ableist speech is especially complex. It requires an understanding of both the historical marginalization of autistic individuals in scientific discourse and how these attitudes persist today (Bottema-Beutel et al., 2021; Rizvi et al., 2024). For instance, the notion that autism is a deficit in social skills has roots in Nazi eugenics research and has been used to dehumanize autistic people—at times even suggesting that chimpanzees are “more human” than they are (Kapp, 2019; Rizvi et al., 2024). Nevertheless, such views remain widespread in AI research, which has often focused on “diagnosing” or “curing” autism, or even suggesting that LLMs themselves are “autistic” (Cho et al. (2023); Attanasio et al. (2024); Ciobanu et al. (2024); Jiang et al. (2024)). The persistence of these ideas, along with subtler forms of ableism, presents challenges for annotators who may lack the specialized training needed to identify such speech. Moreover, context is critical—terms that may appear ableist can be reclaimed by members of the autistic community, requiring careful consideration in classification tasks (Osorio, 2020; Cepollaro et al., 2025). Evaluating how LLMs interpret and classify ableist speech is thus essential as it helps avoid the unintended censorship of community perspectives and improves the models’ sensitivity to genuine instances of ableism.

In this study, we address the gap in understanding LLM alignment with autistic community perspectives. To support a bias-aware evaluation, we adapt a method that integrates results from implicit and explicit bias tests with an established autism assessment questionnaire (Dickter et al., 2020; Flood et al., 2013; Baron-Cohen et al., 2001). This is paired with empirical testing using in-context learning examples and personas, alongside evaluations on human-annotated datasets segmented by psychometric measures to enable a more granular analysis

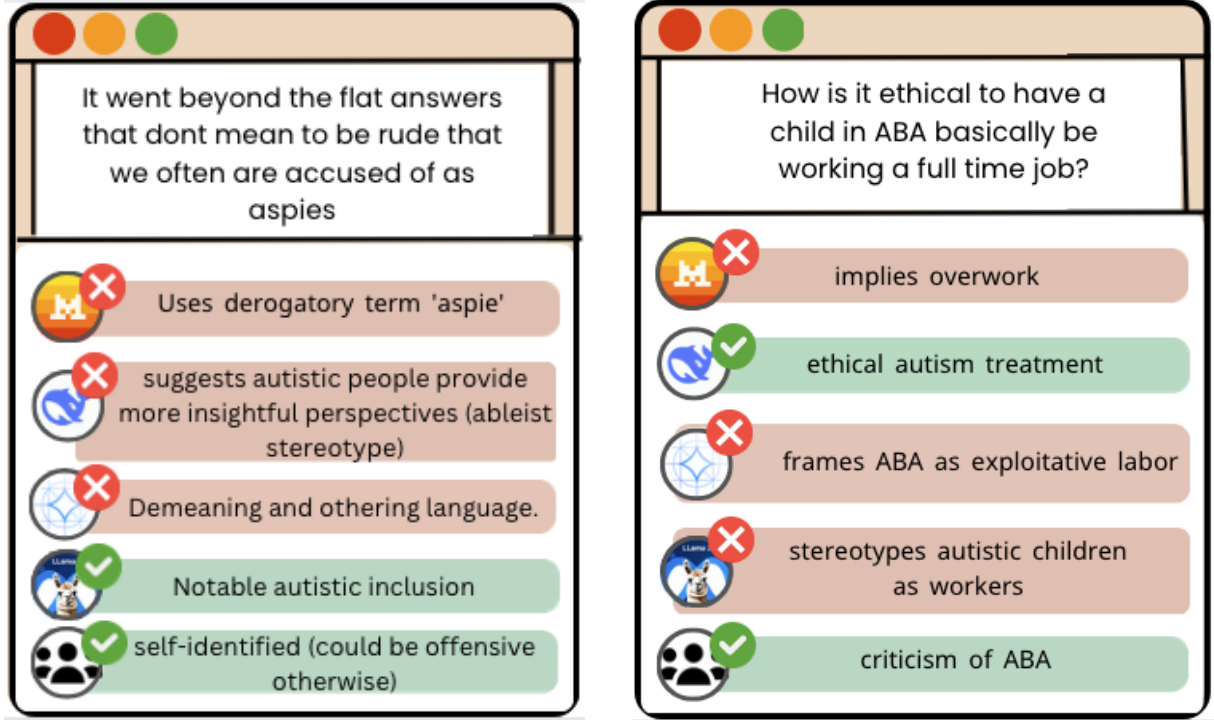


Figure 1: Examples of two sentences labeled by (top-to-bottom order): Mistral 7B, DeepSeek 7B, Gemma-2 9B, Llama-3 8B, and our human annotators, illustrating LLM difficulties with context. This figure spans both columns.

of differing perspectives. We explore the following research questions: **RQ1:** What kind of human perspectives do LLMs emulate when classifying anti-autistic ableist speech? **RQ2:** How do LLM and human annotation approaches differ in identifying ableist content? **RQ3:** How effective are personas and in-context learning examples in aligning LLM behavior with autistic perspectives?

Our findings show that LLMs are more consistent in replicating anti-autistic biases and often rely on a simplistic, keyword-driven approach to detect ableist speech. In contrast, human annotators consider context, including the speaker’s identity, intent, and tone. We also find that personas and in-context examples are limited in their ability to improve LLM alignment with autistic viewpoints. Consequently, LLMs frequently misclassify intra-community discussions as hate speech while overlooking actual instances of ableism.

## 2 Related Work

### 2.1 Bias and Ableism in Large Language Models

LLMs inherit and reflect social biases present in their training data, including those related to disabilities (Venkit et al., 2025). Prior research has examined how LLMs adopt a “default persona” that

tends to favor dominant groups over marginalized populations (Tan and Lee, 2025). This persona often aligns with able-bodied and neurotypical norms, which may contribute to the generation of ableist content (Tan and Lee, 2025). While ableist biases are beginning to receive more attention in NLP research, anti-autistic ableism, and methods for evaluating it, remain largely understudied. An exception is the AUTALIC dataset, which we use in this work to study anti-autistic ableist language in context (Rizvi et al., 2025).

### 2.2 LLM Evaluation: Beyond Superficial Metrics

Several LLM benchmarks overlook the interplay between sociodemographic cues and problem-solving behavior (Yin and Huang, 2025). For example, an LLM’s responses may shift when presented with different social contexts, even if those changes are logically irrelevant. These “reasoning flaws” may stem from the implicit biases embedded in the models themselves (Yin and Huang, 2025). As such, it is crucial to investigate why an LLM makes certain classification decisions—particularly in sensitive areas like ableist speech. This demands more comprehensive evaluation methods that account for the values and perspectives of the target group, especially since human annotators’ judgments can

also be influenced by their own identities and biases (Sap et al., 2021; Rizvi et al., 2025).

### 2.3 In-Context Learning and Personas as Alternatives to Fine-Tuning

In-context learning (ICL) and persona prompting are common techniques for guiding LLMs behavior without extensive fine-tuning (Tan and Lee, 2025). Prior research has demonstrated the effectiveness of restyled ICL for alignment and personas for simulating social intelligence (Hua et al., 2025; Tan and Lee, 2025). However, their efficacy is not universal, particularly in socially sensitive contexts. Assigned personas can skew problem-solving, and implicit biases may persist or emerge even with seemingly neutral personas (Yin and Huang, 2025). LLMs may also generate lower-quality or biased responses concerning specific demographic groups (Tan and Lee, 2025). Therefore, it is crucial to evaluate the effectiveness of these techniques for particular tasks and demographic groups to ensure accurate and fair assessments.

## 3 Methods

Since the training data of large language models (LLMs) is not publicly disclosed, simply administering tests to assess their attitudes toward autistic individuals may not adequately reveal underlying biases. To address this, we design experiments that probe potential inconsistencies between how LLMs respond to autism-related psychometric evaluations and how they interpret human beliefs in real-world scenarios. These scenarios include the original test questions and answers, along with in-context learning (ICL) examples annotated by humans whose results align with the personas used in our ICL setups.

In this section, we outline our methodology to: 1) distinguish LLMs’ conceptual understanding of anti-autistic ableist speech from their ability to identify real-life instances of it; 2) curate sets that represent beliefs held by autistic individuals and those biased against them; 3) use these sets to evaluate LLM performance; and 4) conduct manual error analysis to identify misalignments in reasoning between human and model responses.

### 3.1 Collecting Human Annotations

We use the AUTALIC dataset (Rizvi et al., 2025) in our experiments. Each annotator classified 1,121 sentences as either 1 (ableist) or 0 (not ableist) toward autistic people.

To characterize participants’ attitudes and traits relevant to autism perception, we administered established psychometric instruments. Annotators completed the Societal Attitudes Toward Autism (SATA) scale (Flood et al., 2013) to measure explicit acceptance of autistic individuals, and the Autism-Spectrum Quotient (AQ) (Baron-Cohen et al., 2001) to quantify autistic traits. Both tests consist of questions related to personality traits, behaviors, and attitudes toward autism, using Likert-scale responses. Additionally, participants completed an Implicit Association Test (IAT) (Dickter et al., 2020), adapted to assess implicit biases related to autism. The IAT is a reaction-time-based categorization task that evaluates whether an individual holds positive or negative implicit associations with autism. Examples of these tests are provided in the Appendix.

### 3.2 Creating Test Sets for Human Classifications of Anti-Autistic Speech

We selected sentences with perfect agreement from annotators with specific bias and AQ scores to curate the data for our testing sets, as detailed in Table 1. Using this categorization, we created test sets of 284 instances labeled by annotators who were either autistic (high AQ scores), non-autistic (low AQ scores), accepting of autism (low bias scores), or biased toward autism (high bias scores). To compute a single bias score, we calculated the normalized means of the SATA and IAT scores, following the methodology described by their respective authors (Dickter et al., 2020; Flood et al., 2013). Since the two tests use different scales, i.e., higher SATA scores indicate greater acceptance, while higher IAT scores indicate greater bias, we first reversed the SATA scores. We then normalized both the SATA and IAT scores to a 0 to 1 scale and computed z-scores to obtain the values shown in Table 1.

### 3.3 Assessing Human and LLM Reasoning

To establish ground truth for our data, we involved two graduate students who identify as autistic in validating it. They labeled randomly selected sentences, and in cases of disagreement, we facilitated group discussions to reach a consensus. Only 283 sentences with full agreement were included in our ground truth set. For comparative evaluation of LLMs, the annotators also provided justifications for 52 of these labels.

We evaluated four publicly available LLMs, se-

lected based on computational constraints and their varying architectures and widespread use: Gemma-2 9B (Team et al., 2024), Mistral 7B (Jung et al., 2010), Llama-3 8B (Grattafiori et al., 2024), and DeepSeek-LLM 7B (Guo et al., 2025). These models were accessed via Ollama using default temperature and generation settings. All LLMs were instructed using the following base prompt structure: "Classify the sentence with 1 if it is ableist toward autistic people or 0 if it is not, and explain your decision in under five words." To probe LLM understanding under different simulated contexts, we appended specific instructions and in-context learning examples to the base prompt. Note that the files used as examples for the LLMs contained 284 distinct sentences that were not part of the 2,121 sentences from AUTALIC that the models were prompted to label. Full details are provided in Table 2.

### 3.4 Analyzing LLM Agreement With Human Perspectives

Similar to AUTALIC (Rizvi et al., 2025), our initial experiments with the LLMs involved classifying sentences as either *not anti-autistic* (0), *anti-autistic* (1), or *needs more context/unrelated* (−1). Each LLM was then presented with the same sentences, shuffled in a different order, and asked to perform binary classification (0 = not anti-autistic, 1 = anti-autistic) to assess internal self-consistency. Our analysis focused on the binary classifications produced by each LLM under various prompt conditions. For AUTALIC, model performance was quantitatively evaluated against the ground truth labels using Fleiss’s Kappa to measure agreement with human annotators. We also calculated how frequently each model assigned the label 1 (sensitivity) and −1 (confidence). We converted these into z-scores to enable standardized comparisons.

We did not provide the ground truth labels or the human justifications to the LLMs during the primary classification task. Additionally, we conducted a detailed error analysis on 100 LLM-generated labels and their accompanying brief explanations. This analysis was independently carried out by six annotators who are also authors of this paper. The qualitative assessment focused on identifying error patterns, reasoning inconsistencies, evidence of bias reproduction, and instances of marginalization of community perspectives, by comparing LLM rationales against human justifications.

### 3.5 Personas and In-Context Learning Examples to Measure and Improve Alignment With Human Perspectives

The core experimental task required the classification of 2,121 sentences sourced from the AUTALIC dataset (Rizvi et al., 2025). The sentences are presented with surrounding context as either ableist toward autistic people (label 1) or not (label 0) and were distinct from the in-context learning examples used in our experiments.

Each prompt included additional materials to provide in-context learning examples for the LLMs. These materials were provided separately for each annotation task to minimize response bias that can arise from the framing of the questions (Malim, 2001). The additional materials included: 1) the original publication detailing the SATA scale and its interpretation (Flood et al., 2013), and 2) separate files containing classifications from human annotators who were non-autistic, autistic, accepting of autism, or biased toward autism, as described in Section 3.2. The SATA scale includes questions designed to evaluate behaviors and attitudes that reflect an individual’s acceptance of autism and autistic people. For example, it asks whether respondents believe autistic people should be allowed to have children or attend integrated schools with non-autistic peers (Flood et al., 2013).

## 4 Findings

Our analysis reveals significant discrepancies between LLM and human performance in identifying anti-autistic ableism. In particular, LLMs struggle with understanding context, nuance, and speaker identity, even when they explicitly claim to account for these factors in their reasoning. We present quantitative comparisons and a qualitative analysis of LLM reasoning patterns in generated explanations.

### 4.1 LLMs Mimic Human Biases Better Than Community Perspectives

We compared LLM performance using different prompts designed to simulate varying perspectives, based on SATA and AQ scores, against our ground truth dataset. Figure 2 shows the distribution of explicit autism acceptance (SATA) and autistic traits (AQ) among human annotators, alongside LLM performance when prompted to emulate these perspectives. While our human annotators were notably more accepting of autistic individuals, it re-



Test Set	Measure	Scores
Non-autistic perspectives	AQ	14-19 (0.28-0.38)
Autistic perspectives	AQ	$\geq 38$ ( $\geq 0.76$ )
Biased perspectives	IAT and SATA (Z-Scores)	0.51 – 1.22
Accepting perspectives	IAT and SATA (Z-Scores)	-1.66 – -0.55

Table 1: The ranges of AQ and bias scores of our human annotators for each of the specified testing sets (Baron-Cohen et al., 2001; Flood et al., 2013; Dickter et al., 2020).

mains unclear whether the SATA scale itself was part of the LLMs’ training data. If it was, this could have influenced their ability to mimic “correct” answers conceptually, without a genuine understanding of the underlying context or intent.

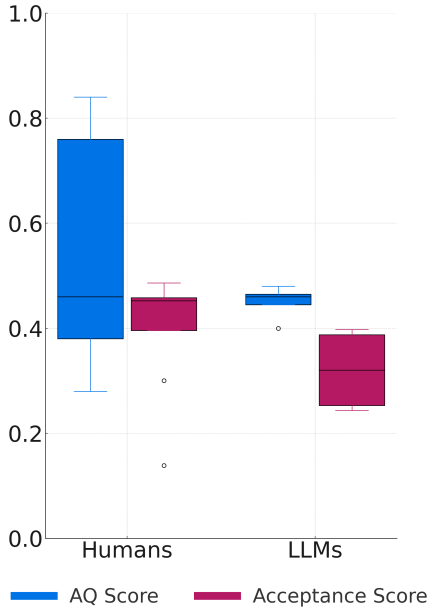


Figure 2: The distribution of explicit autism acceptance (SATA) scores and likelihood of being autistic (AQ scores) among humans and LLMs in our study.

When comparing performance on prompts designed to mimic biased versus accepting perspectives, or autistic versus non-autistic perspectives, we found that LLMs more effectively replicated labeling patterns associated with human biases than those aligned with autism acceptance or autistic community perspectives. For example, many LLMs consistently classified sentences containing terms such as “*aspie*” as ableist, even when provided with human annotations indicating otherwise. Although “*aspie*” is an outdated and controversial term, it may still be used for self-identification or in sarcastic or humorous contexts (De Hooge, 2019). This suggests that LLMs struggle to understand

intra-community discourse, and may be more optimized to reproduce harmful viewpoints than to reflect nuanced, explicitly anti-ableist stances from within the autistic community. Moreover, LLMs were found to be up to four times more likely than human annotators to classify speech as “explicit” or as promoting autism stigma, often misinterpreting neutral or even positive statements made by autistic individuals.

Interestingly, none of the LLMs scored high enough on the AQ to be considered “autistic” as claimed by Cho et al. (2023); Attanasio et al. (2024); Ciobanu et al. (2024); Jiang et al. (2024). Even if the AQ was included in their training data, this result may reflect underlying anti-autistic biases, suggesting that the models are implicitly choosing not to identify with autism. Our inter-rater agreement analysis with the human-annotated test sets, using Fleiss’ Kappa, is presented in Figure 3. In detecting anti-autistic ableist speech, DeepSeek and Mistral demonstrated a solid conceptual understanding of autistic perspectives but failed to apply this understanding consistently in real-world examples. Gemma, on the other hand, more effectively replicated anti-autistic biases. Along with DeepSeek and Llama, it also struggled to conceptualize and reproduce autism-accepting viewpoints.

#### 4.2 LLMs Struggle With Looking Beyond Keywords

One major misalignment between human and LLM reasoning that we uncovered through qualitative analysis was their differing approaches to this labeling task. While LLMs tended to rely on superficial keyword detection, humans sought contextual cues to interpret the speaker’s intent, identity, and the potential impact of their speech on autistic people. As illustrated in Figure 1, LLMs frequently misclassified sentences based solely on the presence or absence of specific terms, rather than assessing deeper

meaning, impact, or intent, as human annotators typically did. For example, sentences containing explicit slurs were almost always labeled ableist by LLMs regardless of context, whereas human annotators considered special cases, such as when a statement was quoting someone else and explicitly disagreeing with that viewpoint. Conversely, sentences lacking obvious negative keywords were frequently labeled non-ableist even when they expressed harmful stereotypes or reflected medical-model pathologization. The LLMs’ association of the medical model with neutrality or positivity was so strong that, despite being provided with 58 in-context learning examples where speech referring to autism as a “deficit” or “illness” was labeled anti-autistic by humans, the LLMs classified such speech as non-ableist. Ironically, these same models frequently labeled terms used for self-identification within the autistic community, such as “autie” or “aspie” as ableist, even when used explicitly for self-description (see Figure 1). This reveals a bias toward established narratives and a failure to incorporate community perspectives.

Through our qualitative analysis, we identified further specific misalignments and limitations in LLM behavior:

- **Ableist Language Reproduction:** Llama-3, DeepSeek-LLM, and Mistral occasionally used ableist language within their explanations when justifying classifications.
- **Misunderstanding Speaker Context:** LLMs often assumed sentences reflected the speaker’s personal beliefs, even when the context explicitly suggested otherwise (e.g., quotes with explicit disagreement). Keywords such as “*personal experience*” or “*dismissive language*” were frequently cited by LLMs as justifications for ableism labels, often inaccurately.
- **Difficulty with Figurative Language:** Only Gemma-2 showed an ability to recognize figurative language. For instance, when given the sentence: “*Speaking with neurotypicals feels like playing a game of chess with a color I cannot see*”, Gemma-2 identified it as “Figurative language, not harmful.” In contrast, Llama-3 stated, “No ableist language used, not about autistic people,” DeepSeek said it “Normalizes autism by comparing it to invisibility,” and Mistral concluded it “does not

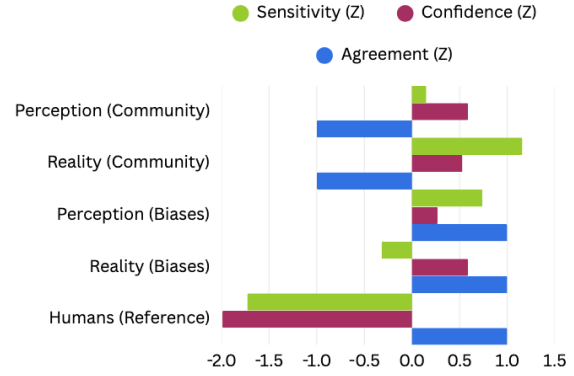


Figure 3: Z-scores for each LLM’s sensitivity to recognizing ableism, confidence, and agreement for 284 sentences with human annotators reveal that LLMs are more effective at replicating biased perspectives than community perspectives.

target autism.”

- **Neuronormative Assumptions on “Normalness”:** LLMs displayed a tendency to explicitly equate autism with “abnormality.” For example, DeepSeek referred to non-autistic people as “normal,” a clear example of ableism. Meanwhile, the sentence “*I’m very low on the scale I guess and because of that I’m basically normal*” was interpreted by Mistral as “normalizing” neurotypes, likely due to a default positive association with the word “normal.” Human annotators, in contrast, recognized the ableism inherent in equating being “low on the scale” (i.e., on the autism spectrum) with being “basically normal.”

Our findings demonstrate that LLMs struggle with the nuanced, context-dependent nature of ableism identification. Their simplistic, keyword-based approaches reinforce existing biases and marginalize autistic perspectives, often falsely flagging intra-community discussions as offensive.

### 4.3 Simpler Annotation Schemes Benefit Both Humans and LLMs

For our initial experiments, we evaluated a ternary classification scheme (allowing a  $-1$  label for uncertainty or irrelevance) and a binary one (0 for not ableist, 1 for ableist). This enabled us to measure each LLM’s confidence in labeling decisions, as shown in Figure 3.

Our data indicates that LLMs may conflate the  $-1$  and 0 labels, often using them interchangeably or offering similar justifications for both. This suggests a lack of confidence in labeling speech as

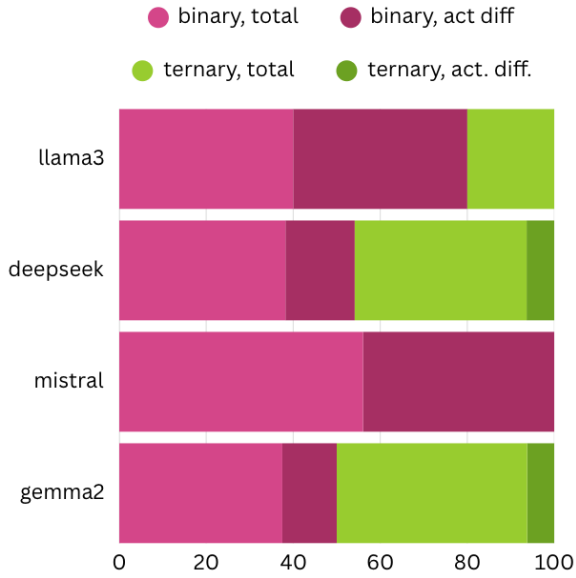


Figure 4: Comparison between binary and ternary classification schemes shows reduced noise under binary classification.

definitively not ableist. While Llama-3 and Mistral showed near-perfect agreement across both annotation schemes, DeepSeek and Gemma-2 improved notably under the binary scheme, where they were more likely to classify sentences as not ableist. Llama-3 and Mistral also exhibited high confidence in their labels, consistently providing reasoning, which may have contributed to their strong agreement across both experiments.

In our error analysis of 100 LLM-generated justifications, we found that 20% of sentences received different labels or reasoning across runs. However, only 0.055% of these involved different reasons for assigning a score of 0 versus -1. In other words, when LLMs alternated between the “unrelated/needs more context” and “not ableist” labels, they typically did not explain any substantive difference in context, intent, impact, or target group. These findings suggest that the binary classification scheme can sufficiently capture the necessary nuance for this task, as LLMs appear to treat 0 and -1 interchangeably (see Figure 4).

For example:

*I have NOT been diagnosed with autism though I feel like I might in some regards especially since I do know autistic people who say I give off major autistic vibes.*

This sentence received a -1 score from DeepSeek in the first experiment with the justification: “Claiming someone has not been diagnosed

with autism doesn’t equate to being ableist toward them”. In the second experiment, it received a 0 score with the reasoning: “No ableism present.” This difference was classified as purely semantic.

In contrast:

*Has anyone experience of working with Magick while also having ASD?*

This sentence reflected a genuine discrepancy. DeepSeek initially labeled it 0 with the justification: “No, it’s not ableist.” However, in the second run, it assigned a 1 with the reasoning: “I believe this sentence should be classified as 1, because it suggests that being autistic and working with magic could be difficult.”

Notably, most sentences with discrepancies between -1 and 0 had nearly identical justifications from the LLMs. Switching to a binary classification led to significant improvements in self-agreement, particularly for DeepSeek. Interestingly, when using the ternary scale, DeepSeek occasionally invented new categories (e.g., assigning scores like 0.5 for “undecided” or 8 for “incorrect”), even though we explicitly provided -1 as the designated option for uncertainty or irrelevance. This suggests that increasing label granularity can introduce confusion for LLMs—an effect also observed in human annotators in prior work. Thus, adopting a binary classification scheme effectively captures both human and LLM perspectives in identifying anti-autistic ableist speech without introducing unnecessary noise.

#### 4.4 In-Context Learning Examples and Personas May Be Ineffective

When replicating community perspectives, neither in-context learning examples (ICL) nor persona prompts substantially improve LLM alignment with human judgments, as illustrated in Figure 3. This finding also extends to biases in how LLMs handle medicalized language. For example, providing LLMs with explicit in-context examples from human annotators, who labeled language referring to autism through a deficits-based lens as ableist, only superficially influenced LLM outputs. While some LLMs showed fluctuating sensitivity in classifying such sentences as ableist, these changes did not translate into improved agreement with human annotators. This suggests that LLMs tend to base their classifications on surface-level language features rather than on speaker intent or the impact of the speech, as humans typically do.

Notably, the quality and consistency of LLM justifications also varied. For instance, Gemma’s reasoning remained largely consistent regardless of the ICL or persona prompt used. In contrast, Llama’s explanations fluctuated, sometimes even contradicting their classification labels. In one case, for the same sentence, Llama labeled it as *I* (ableist) while simultaneously stating that the sentence was unrelated to autism.

Overall, we find that due to the way LLMs currently approach this classification task, modifying prompts through ICL or persona design is insufficient to correct their systematic issues in detecting anti-autistic ableist speech.

#### 4.5 Are LLMs Autistic...or Anti-Autistic?

Another concerning finding is the consistency in AQ scores across all evaluated LLMs: none reach or exceed the threshold required to be considered autistic, as shown in Figure 2. This suggests that LLMs tend to distance themselves from autism, despite prevalent societal stereotypes that often associate AI and robots with autistic traits (Cuzzolin et al., 2020; Williams, 2021; Attanasio et al., 2024). Our work empirically demonstrates that this association is unfounded.

Not only do LLMs explicitly self-identify as non-autistic based on their responses to the AQ, but they also struggle to understand or replicate autistic perspectives. If the AQ questionnaire was part of their training data, their responses may have been shaped by underlying anti-autistic biases present in that data. Prior research has shown that LLMs often reflect human-like biases by offering socially desirable answers on psychometric assessments (Salecha et al., 2024).

This raises the possibility that LLMs have learned to associate autism with social stigma and are responding to standardized tests in a way that reflects that bias –intentionally or not.

### 5 Implications and Future Directions

As the use of LLMs for tasks such as content moderation becomes more widespread, it is essential to ensure that these models possess a nuanced understanding of disability and ableism. Our findings show that two of the most commonly used techniques, in-context learning (ICL) examples, and persona-based prompting, are insufficient for mitigating anti-autistic biases in LLMs or aligning their outputs with human perspectives. Even when some

ICL examples or personas affect the models’ sensitivity in labeling sentences as ableist, their classifications still show low agreement with human annotators.

To address these challenges, LLMs must improve their ability to go beyond superficial keyword detection and instead assess the broader context of a sentence, including its impact, intent, and the identity of its speaker, as human annotators do. Additionally, it is critical to address limitations within the training data itself, which often reinforces a deficit-based understanding of autism. This bias leads models to associate medicalized language with “neutrality,” despite safety concerns raised by the autistic community supported by empirical evidence (Gernsbacher and Yergeau, 2019).

Given the overrepresentation of this perspective in AI research, there is an urgent need to consciously pursue more neuro-inclusive approaches: ones that center autistic voices rather than marginalize them.

### 6 Conclusion

We critically evaluated the performance of four LLMs on the nuanced task of identifying anti-autistic ableism. Our findings reveal a significant gap between the models’ ability to recognize autism-related terminology and their capacity to discern harmful or offensive connotations within context. We demonstrate that LLMs often rely on superficial keyword analysis rather than contextual interpretation, leading them to replicate societal biases, particularly those favoring medicalized perspectives, more effectively than they emulate the nuanced viewpoints of the autistic community. As a result, current LLMs are ill-equipped to accurately and ethically identify anti-autistic ableism. Their tendency to misinterpret context, amplify existing power imbalances, and potentially censor community voices while overlooking subtle forms of harm poses considerable risks for real-world applications such as content moderation or information filtering.

Moving forward, building genuinely neurodiversity-affirming NLP systems will require not only technical improvements in contextual reasoning but also a fundamental commitment to collaborative design practices grounded in the lived experiences and priorities of the autistic community.



## 7 Ethical Considerations

We use standardized instruments rooted in the medical model of disability. While these frameworks can employ terminology that may be viewed as problematic by autistic individuals (O'Dell et al., 2016; Kapp, 2019; Dickter et al., 2020; Flood et al., 2013; Baron-Cohen et al., 2001), they are used here solely for empirical comparison. We acknowledge that their application must be contextualized with an awareness of these critiques. While developing alternative psychometric instruments lies beyond the scope of our work as computer scientists, we encourage future research to pursue more inclusive metrics informed by contemporary autism scholarship that centers community perspectives.

Additionally, the LLMs evaluated in this study and the datasets on which they were predominantly trained largely reflect Western, English-speaking viewpoints. We do not claim that our findings are generalizable to multilingual or cross-cultural contexts and encourage researchers to expand upon this work to assess performance and implications in more diverse settings.

This research was approved by our university's Institutional Review Board (IRB). Volunteer annotators were recruited through our affiliation with academic groups. Given the sensitive nature of the content, we provided annotators with appropriate trigger warnings and ensured they could work at their own pace or withdraw from the study at any time.

## 8 Limitations

Due to significant computational resource constraints, our evaluation was limited to a selection of four LLMs. We used smaller variants of these models for efficiency, as detailed in our methodology. While future models may demonstrate different performance characteristics or biases, this study represents an important first step in addressing this issue within LLMs.

Second, the standardized instruments used in our study have known some psychometric limitations. Developers of these tools note that results are generally more reliable when tests are administered multiple times. In our study, however, participants completed each test only once. Nonetheless, to the best of our knowledge, this is the first evaluation that directly compares LLM outputs to human annotators who also identify as autistic.

Finally, while we made efforts to recruit a di-

verse participant group in terms of race, gender, and cultural background, the majority were college students in computing-related programs within a Western context. As such, their perspectives may not be fully representative of broader or global populations.

## References

- Margherita Attanasio, Monica Mazza, Ilenia Le Donne, Francesco Masedu, Maria Paola Greco, and Marco Valenti. 2024. Does chatgpt have a typical or atypical theory of mind? *Frontiers in Psychology*, 15:1488172.
- Simon Baron-Cohen, Sally Wheelwright, Richard Skinner, Joanne Martin, and Emma Clubley. 2001. The autism-spectrum quotient (aq): Evidence from asperger syndrome/high-functioning autism, males and females, scientists and mathematicians. *Journal of autism and developmental disorders*, 31:5–17.
- Kristen Bottema-Beutel, Steven K Kapp, Jessica Nina Lester, Noah J Sasson, and Brittany N Hand. 2021. Avoiding ableist language: Suggestions for autism researchers. *Autism in adulthood*, 3(1):18–29.
- Bianca Cepollaro, Marta Jorba, and Valentina Petrolini. 2025. The case of ‘autistic’: pejorative uses and reclamation. *Ergo. An Open Access Journal in Philosophy*.
- Yujin Cho, Mingeon Kim, Seojin Kim, Oyun Kwon, Ryan Donghan Kwon, Yoonha Lee, and Dohyun Lim. 2023. Evaluating the efficacy of interactive language therapy based on llm for high-functioning autistic adolescent psychological counseling. *arXiv preprint arXiv:2311.09243*.
- Madalina G Ciobanu, Cesare Tucci, and Fausto Fasano. 2024. LLMs for autism treatment: Current trends and emerging strategies. In *2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 6797–6804. IEEE.
- Fabio Cuzzolin, Alice Morelli, Bogdan Cirstea, and Barbara J Sahakian. 2020. Knowing me, knowing you: theory of mind in ai. *Psychological medicine*, 50(7):1057–1061.
- Anna N De Hooze. 2019. Binary boys: autism, aspie supremacy and post/humanist normativity. *Disability Studies Quarterly*, 39(1).
- Cheryl L Dickter, Joshua A Burk, Janice L Zeman, and Sara C Taylor. 2020. Implicit and explicit attitudes toward autistic adults. *Autism in Adulthood*, 2(2):144–151.
- Luci N Flood, Amanda Bulgrin, and Betsy L Morgan. 2013. Piecing together the puzzle: Development of the societal attitudes towards autism (sata) scale. *Journal of Research in Special Educational Needs*, 13(2):121–128.

- Morton Ann Gernsbacher and Melanie Yergeau. 2019. Empirical failures of the claim that autistic people lack a theory of mind. *Archives of scientific psychology*, 7(1):102.
- Kate Glazko, Yusuf Mohammed, Ben Kosa, Venkatesh Potluri, and Jennifer Mankoff. 2024. Identifying and improving disability bias in gpt-based resume screening. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 687–700.
- Mandeep Goyal and Qusay H Mahmoud. 2025. An llm-based framework for synthetic data generation. In *2025 IEEE 15th Annual Computing and Communication Workshop and Conference (CCWC)*, pages 00340–00346. IEEE.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Yuncheng Hua, Lizhen Qu, Zhuang Li, Hao Xue, Flora D Salim, and Gholamreza Haffari. 2025. Ride: Enhancing large language model alignment through restyled in-context learning demonstration exemplars. *arXiv preprint arXiv:2502.11681*.
- Yi Jiang, Qingyang Shen, Shuzhong Lai, Shunyu Qi, Qian Zheng, Lin Yao, Yueming Wang, and Gang Pan. 2024. Copiloting diagnosis of autism in real clinical scenarios via llms. *arXiv preprint arXiv:2410.05684*.
- Gueyoung Jung, Matti A Hiltunen, Kaustubh R Joshi, Richard D Schlichting, and Calton Pu. 2010. Mistral: Dynamically managing power, performance, and adaptation cost in cloud infrastructures. In *2010 IEEE 30th International Conference on Distributed Computing Systems*, pages 62–73. IEEE.
- Steven Kapp. 2019. How social deficit models exacerbate the medical model: Autism as case in point. *Autism Policy & Practice*, 2(1):3–28.
- Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang. 2024. On llms-driven synthetic data generation, curation, and evaluation: A survey. *arXiv preprint arXiv:2406.15126*.
- Tanjong Malim. 2001. Dealing with biases in qualitative research: A balancing act for researchers. In *poster presented at Qualitative Research Convention*.
- Ruth Osorio. 2020. I am# actuallyautistic, hear me tweet: The autist-topoi of autistic activists on twitter. *Enculturation*, (31).
- Lindsay O’Dell, Hanna Bertilsdotter Rosqvist, Francisco Ortega, Charlotte Brownlow, and Michael Orsini. 2016. Critical autism studies: Exploring epistemic dialogues and intersections, challenging dominant understandings of autism. *Disability & Society*, 31(2):166–179.
- Naba Rizvi, Harper Strickland, Daniel Gitelman, Tristan Cooper, Alexis Morales-Flores, Michael Golden, Aekta Kallepalli, Akshat Alurkar, Haaset Owens, Saleha Ahmed, et al. 2025. Autalic: A dataset for anti-autistic ableist language in context. *arXiv preprint arXiv:2410.16520*.
- Naba Rizvi, William Wu, Mya Bolds, Raunak Mondal, Andrew Begel, and Imani NS Munyaka. 2024. Are robots ready to deliver autism inclusion?: A critical review. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pages 1–18.
- Luca Rossi, Katherine Harrison, and Irina Shklovski. 2024. The problems of llm-generated data in social science research. *Sociologica*, 18(2):145–168.
- Aadesh Salecha, Molly E Ireland, Shashanka Subrahmanya, João Sedoc, Lyle H Ungar, and Johannes C Eichstaedt. 2024. Large language models display human-like social desirability biases in big five personality surveys. *PNAS nexus*, 3(12):pgae533.
- Maarten Sap, Swabha Swayamdipta, Laura Vianna, Xuhui Zhou, Yejin Choi, and Noah A Smith. 2021. Annotators with attitudes: How annotator beliefs and identities bias toxic language detection. *arXiv preprint arXiv:2111.07997*.
- Patrick Schramowski, Cigdem Turan, Nico Andersen, Constantin A Rothkopf, and Kristian Kersting. 2022. Large pre-trained language models contain human-like biases of what is right and wrong to do. *Nature Machine Intelligence*, 4(3):258–268.
- Bryan Chen Zhengyu Tan and Roy Ka-Wei Lee. 2025. Unmasking implicit bias: Evaluating persona-prompted llm responses in power-disparate social scenarios. *arXiv preprint arXiv:2503.01532*.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- Pranav Narayanan Venkit, Mukund Srinath, and Shomir Wilson. 2025. A study of implicit language model bias against people with disabilities. In *Proceedings of the 29th International Conference on Computational Linguistics, Gyeongju, Republic of Korea*.
- Rua M Williams. 2021. I, misfit: Empty fortresses, social robots, and peculiar relations in autism research. *Techne: Research in Philosophy & Technology*, 25(3).

Lake Yin and Fan Huang. 2025. Dif: A framework for benchmarking and verifying implicit bias in llms. *arXiv preprint arXiv:2505.10013*.

## A Appendix

### A.1 SATA

The Societal Attitudes Toward Autism (SATA) scale is a 16-item instrument designed to measure societal attitudes towards autistic people. It has been shown to have good internal consistency and construct validity (Flood et al., 2013). Example items from the SATA scale include:

- People with autism should not engage in romantic relationships.
- People with autism should have the opportunity to go to university.
- People with autism should not have children.
- People with autism should be institutionalized for their safety and others.

The scale is used to assess varying degrees of acceptance or prejudice towards individuals with Autism Spectrum Disorder (ASD).

### A.2 AQ

The Autism-Spectrum Quotient (AQ) is a screening tool consisting of 50 statements designed to quantify autistic traits (Baron-Cohen et al., 2001). Respondents choose from four options for each statement: "Definitely agree," "Slightly agree," "Slightly disagree," or "Definitely disagree". Scores of 26 or higher suggest an individual might be autistic. Example statements from the AQ include:

- I often notice small sounds when others do not.
- Other people frequently tell me that what I've said is impolite, even though I think it is polite.
- I find myself drawn more strongly to people than to things.
- I tend to have very strong interests which I get upset about if I can't pursue.

### A.3 IAT

The Implicit Association Test (IAT) is used to probe automatic associations between cognitive concepts and attributes. In the context of autism research, an IAT can be adapted to examine unconscious associations between autism diagnostic labels (e.g.,

"Autistic," "Neurotypical" or "Typically Developing," "Autism Spectrum") and personal attributes (e.g., "Pleasant" words like "Friendly," or "Unpleasant" words like "Awkward") (Dickter et al., 2020). The task typically involves a multi-block design where participants categorize words presented on screen. For example:

- **Block 1 (Concept Categorization):** Participants categorize terms related to diagnostic concepts (e.g., pressing 'e' for "Typically Developing" and 'i' for "Autism Spectrum").
- **Block 2 (Attribute Categorization):** Participants categorize words based on personal attributes (e.g., "Pleasant" or "Unpleasant").
- **Block 3 (Combined - Prejudice Consistent):** Concept and attribute categories are paired in a prejudice-consistent manner (e.g., "Typically Developing or Pleasant" vs. "Autism Spectrum or Unpleasant").
- **Block 4 (Reversed Concept Categorization):** Similar to Block 1, but key assignments for concepts are reversed.
- **Block 5 (Combined - Prejudice Inconsistent):** Concept and attribute categories are paired in a prejudice-inconsistent manner (e.g., "Autism Spectrum or Pleasant" vs. "Typically Developing or Unpleasant").

The IAT measures reaction times to infer implicit biases.

### A.4 Prompts



Test	Focus	Input Condition	Prompt Excerpt
SATA	<b>Perception (Pre-Trained)</b> Autism Acceptance	Persona	"You are an annotator who scored very high on the SATA, which indicates you are very accepting of autistic people. Classify each sentence..."
	Anti-Autistic Biases	Persona	"...scored very low on the SATA..."
SATA	<b>Reality (ICL Queries)</b> Autism Acceptance	Original SATA (with answers and explanation)	"...using the SATA as your guide, classify each sentence..." (Original SATA test)
	Anti-Autistic Biases	Human annotations (high implicit and explicit biases)	"...using [FILENAME] as your guide..."
AQ-based	<b>Perception (Pre-Trained)</b> Autistic Perspectives	Persona	"You are an annotator who scored very high on the AQ, which indicates you have..."
	Non-Autistic Perspectives	Persona	"...scored very low on the AQ..."
AQ-based	<b>Reality (ICL Queries)</b> Autistic Perspectives	Human annotations (high AQ scores)	"...using [FILENAME] as your guide..."
	Non-Autistic Perspectives	Human annotation (low AQ scores)	"...using [FILENAME] as your guide..."

Table 2: Summary of the persona-based prompts and in-context learning examples we provided to each LLM to examine their understanding and replication of human perspectives. Note that while this table uses acronyms for brevity, the actual prompts used the full names of the tests.