

# Impact of LLM Alignment on Impression Formation in Social Interactions

Ala N. Tak<sup>♣\*</sup> Anahita Bolourani<sup>◇\*</sup> Daniel B. Shank<sup>♣</sup> Jonathan Gratch<sup>♠</sup>

<sup>♠</sup> Institute for Creative Technologies, University of Southern California

<sup>◇</sup> Department of Statistics and Data Science, University of California, Los Angeles

<sup>♣</sup> Psychological Science, Missouri University of Science and Technology

[antak@ict.usc.edu](mailto:antak@ict.usc.edu) [anabolourani@ucla.edu](mailto:anabolourani@ucla.edu) [shankd@mst.edu](mailto:shankd@mst.edu) [gratch@ict.usc.edu](mailto:gratch@ict.usc.edu)

## Abstract

Impression formation plays a crucial role in shaping social life, influencing behaviors, attitudes, and interactions across different contexts. Affect Control Theory (ACT) offers a well-established, empirically grounded model of how people form impressions and evaluate social interactions. We investigate whether Large Language Models (LLMs) exhibit patterns of impression formation that align with ACT’s predictions. As a case study, we focus on gendered social interactions—how an LLM perceives gender in a prototypic social interaction. We compare several preference-tuned derivatives of LLaMA-3 model family (including LLaMA-Instruct, Tulu-3, and DeepSeek-R1-Distill) with GPT-4 as a baseline, examining the extent to which alignment or preference tuning influences the models’ tendencies in forming gender impressions. We find that LLMs form impressions quite differently than ACT. Notably, LLMs are insensitive to situational context: the impression of an interaction is overwhelmingly driven by the identity of the actor, regardless of the actor’s actions or the recipient of those actions. This stands in contrast to ACT’s interaction-based reasoning, which accounts for the interplay of identities, behaviors, and recipients. We further find that preference tuning often amplifies or skews certain impressions in unpredicted ways. Our corpus offers a benchmark for assessing LLMs’ social intelligence; we encourage further research using ACT-like frameworks to explore how tuning influences impression formation across diverse social dimensions.

## 1 Introduction

*Warning: Contains explicit or potentially upsetting language.*

To successfully navigate the social world, people must quickly assess others’ traits, intentions, and social roles. This process—known as impression formation—encompasses the cognitive and affective judgments people make about others based on their social roles and actions (Smith-Lovin & Heise, 1988; Reeder et al., 2004; Rogers, 2018). Understanding impression formation is crucial in many domains. For instance, it explains how gendered beliefs shape salaries (Freeland & Harnois, 2020), how lawyers are perceived in their daily practices (Fields, 2023), and how perceptions of criminal behavior align with crime rates (Boyle, 2024). It also aids in understanding cross-cultural differences related to rigidity in cultural norms (Zhao, 2023), investment behaviors shaped by self-perceptions (Charness & Gneezy, 2012; Gurieva et al., 2022), and how social norms can affect leadership roles (Heise et al., 2019; Stewart et al., 2021). Notably, impressions extend beyond isolated entities to the actions people perform and the motivations behind those actions. For example, when learning that someone hit a child, observers form expectations about the person’s identity to explain the action, including their gender (e.g., probably a man) and mental state (e.g., probably angry).

Affect Control Theory (ACT) is a rigorously evaluated sociological model of impression formation that theorizes that impressions are largely determined by the words people use

---

\*Equal contribution.

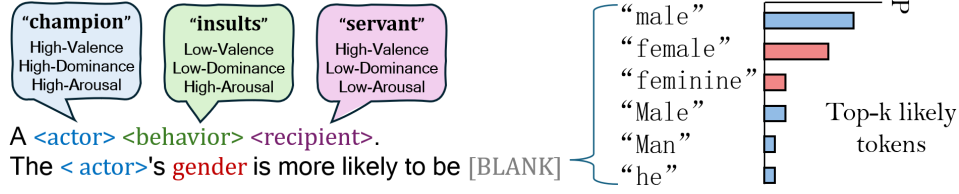


Figure 1: ACT-inspired prompt generation to assess gender impressions in LLMs

to label situations (Heise, 2007). ACT abstractly summarizes social situations in terms of "sentences" made up of an *Actor*, *Behavior*, and *Recipient* (ABR)—e.g., "bully mocks scientist"—given a dictionary of labels that hold shared meaning within a specific culture or group. Whether an actor is labeled, for example, as a "teacher" or a "friend" determines expectations about their character and the actions they are likely to perform. When actors deviate from these expectations, called a *deflection*, actors may be relabeled (e.g., an "angry friend") to explain and mathematically account for the deviation. This model captures social dynamics by making mathematical predictions that connect social action to identities and labels, offering a richer understanding of what is socially normative.

Recent advancements in language models (LLMs) have enabled their application in critical societal domains such as healthcare, hiring, and education—settings where modeling human impression formation is essential. However, prior work on AI impression formation has focused on expectations about social entities in isolation—e.g., homemakers are expected to be women, programmers, men, and criminals Black (Bolukbasi et al., 2016). In contrast, ACT captures how these isolated impressions change when formed in the context of social interactions. Thus, a natural question is whether LLMs are similarly sensitive to the social context. If LLMs "think" like humans, one expects them to reflect the culturally shared defaults that ACT assigns as most normative.

By systematically comparing LLM output to ACT predictions in a well-studied case of gender impressions, we gain insight into how well these models approximate human impression formation in social interactions. Specifically, we examine how *human alignment* post-training or *preference tuning* methods, such as Supervised Fine-tuning (SFT) and Reinforcement Learning from Human Feedback (RLHF), influence model impression formation. Previous research has extensively studied the potential negative impacts of preference tuning and RLHF (Lambert et al., 2023). For instance, models are aligned primarily with English-speaking users (Ouyang et al., 2022). RLHF is shown to steer models toward more assertive outputs (Hosking et al., 2023), reduce novelty in responses (Kirk et al., 2023), make mistakes in outputs more subtle (Bai et al., 2022a), lead to models echoing user opinions or expressing stronger political views (Perez et al., 2022), and create disparities between English dialects and global perspectives (Ryan et al., 2024). Inspired by the ACT framework, we create a benchmark dataset for social interaction impression formation and evaluate Llama-3.1-8B along with its fine-tuned derivatives, using GPT-4 as a baseline. Our findings reveal that alignment processes can shift LLM gender impressions in inconsistent and divergent ways across models. We replicate the experiments using the larger 70B-parameter variants and find that the trends remain consistent across scales. Thus, we specifically address the following questions: **RQ1** Do LLMs show patterns of impression formation in gendered social interactions that are consistent with ACT predictions? **RQ2** How does alignment or post-training shape or distort these patterns?

By focusing on dynamic social interactions, our analysis provides a deeper understanding of how cultural impressions—specifically, gender impressions—are encoded in LLM outputs. This approach offers a foundation for refining model tuning processes. Ultimately, these insights can help develop LLMs that more accurately reflect human-like impression formation while reducing unintended distortions and inconsistencies.

## 2 Background

### 2.1 Affect Control Theory

Many theories of impression formation use dimensional models to characterize impressions. For example, Oosterhof & Todorov (2008) characterize impressions formed from faces into

the orthogonal dimensions of valence (good vs. bad) and dominance (powerful vs. weak). While using slightly different names, ACT also uses valence, dominance, and additionally, arousal (active to passive), resulting in a three-dimensional impression space (VAD: Valence, Arousal, Dominance). These ratings, collected across cultures, are compiled into dictionaries of average affective meanings (Heise, 2007). These dictionaries include mean VAD ratings for thousands of social interaction concepts, demonstrating strong intra-cultural consensus and stability over time (Rogers et al., 2013; Heise, 2010).

ACT posits that individuals seek to confirm *fundamental sentiments*—culturally shaped context-free expectations about themselves, others, and situations—in social interactions. Deflection, measured as the squared Euclidean distance between *contextualized impressions* and context-free *fundamental sentiments* in VAD space, quantifies event unexpectedness. More deflection means more surprise—e.g., a mother yelling at her child mismatches expected nurturing norms (Rogers et al., 2013)—whereas low deflection indicates the event is normative for the actor (Heise & MacKinnon, 1987; Heise & Lerner, 2006; Schröder et al., 2013b; Schröder & Scholl, 2009).

## 2.2 Social gender impressions

Whether examined through the lens of ACT (Langford & MacKinnon, 2000; Rogers et al., 2013) or other related theories (Cikara & Fiske, 2009), research shows the perception of masculinity is linked to power, competence, and dominance, while femininity is associated with goodness, caring, and warmth (Freeland & Harnois, 2020; Fraser et al., 2021; 2022; Nicolas & Caliskan, 2024a). Women are often seen as less agentic and more communal than men, with agency tied to leadership and competence, and communality linked to helpfulness and emotional expression (Heilman, 2012; Wynn & Correll, 2018).

These stereotypes shape both expected (prescriptive) and observed (descriptive) behaviors (Fields, 2023). Notably, masculine traits associated with dominance are rewarded with higher wages, highlighting its influence on occupational income, while valence shows no such effect (Freeland & Harnois, 2020). Dominance is frequently interpreted not only as a measure of power but also as a representation of an identity’s competence within workplace contexts (Cuddy et al., 2008; Rogers et al., 2013). Fields (2023) used ACT to model the ways in which men and women lawyers are perceived when enacting behaviors in their daily practice and showed cultural meanings maintain gender status hierarchies in the field. ACT simulations by Boyle (2024) show that criminal behaviors generate higher deflection for women than men, with larger gender gaps aligning with greater disparities in crime rates based on multiple crime data sources.

## 2.3 LLM Impressions and Stereotypes

Impression formation in language generation can manifest at both local and global levels (Gallegos et al., 2024). Local effects refer to word-context associations, such as differences in next-token likelihoods—e.g., “men/women are known for [BLANK]”, while global effects involve valence or other properties over a span of generated text (Sheng et al., 2019; Nangia et al., 2020; Liang et al., 2021; Yang et al., 2022). Word embeddings offer another way to examine these phenomena in LLMs by numerically encoding semantic relationships and capturing stereotypical associations (Caliskan et al., 2017). Primarily focused on valence-oriented analyses, prior research uses embeddings to identify keywords linked to social categories and analyze them with dictionaries (Nicolas & Caliskan, 2024a; Charlesworth et al., 2022). Loon & Freese (2023) show that despite distinct purposes of embeddings and VAD ratings—capturing the semantic relationships of words vs. measuring consensual affective connotations of concepts—, embeddings can serve as a proxy for VAD ratings.

In examining gender impressions, Caliskan et al. (2022) analyzed English word embeddings and found a masculine default: male-associated terms were predominantly verbs (e.g., fight, overpower) linked to tech, religion, sports, and violence, whereas female-associated terms were mainly adjectives/adverbs (e.g., giving, emotionally) tied to slurs, appearance, and domesticity. Farlow et al. (2024) observed no explicit bias in ChatGPT’s reference letters but noted a male-centric language, potentially stemming from historical associations between men and leadership. Busker et al. (2023) reported that ChatGPT perpetuates gender stereotypes (especially in question-based prompts), while Zhao et al. (2024) identified

communal descriptors tied to female roles in GPT-4. Bai et al. (2024) similarly found GPT-4’s implicit skew, associating STEM fields with boys and humanities with girls.

LLMs are shown to inherit often stereotypical social perceptions due to training on uncensored, Internet-based data (Benjamin, 2019; Bender et al., 2021; Dodge et al., 2021; Sheng et al., 2021; Barocas et al., 2023). However, efforts to align models with human values can produce unintended consequences—as illustrated by the controversial Gemini case, where attempts to promote diversity paradoxically resulted in the generation of images depicting Black Nazis (Field, 2024). Such instances highlight the emergence of “surprising” or “unexpected” biases. Wilson & Caliskan (2024) simulated resume screening using a retrieval method and found that LLMs reinforced societal racial and gender “defaults” that do not necessarily reflect real-world occupational patterns. Takemoto (2024) found that GPT-3.5 and GPT-4 mirrored human tendencies in moral dilemmas but amplified biases, such as a stronger inclination to save females. Fulgu & Capraro (2024) suggest that efforts to address gender parity in recent GPT iterations may inadvertently create extreme disparities as models consistently reinforced feminine stereotypes and attributed masculine stereotypes to females. Spillner (2024) showed ChatGPT often assigns female characters to occupations traditionally associated with men. Additionally, GPT-4 portrays mixed-gender violence more acceptably for a female actor and a male victim (Fulgu & Capraro, 2024).

### 3 Method

#### 3.1 Impression Formation Benchmark Dataset

Following an established practice in ACT research Shank & Burns (2022); Lulham & Shank (2023), we construct a diverse set of events spanning all possible combinations of sentiments for actors, behaviors, and recipients (ABR profiles). Specifically, we generate  $2^3 \times 2^3 \times 2^3$  sentiment configurations, constrained by high/low VAD values. We select four representative items per configuration from the dictionary to generate synthetic events, systematically covering all possible interactions (see Table 3 in Appendix B). We utilize cultural meanings from the USA Combined Surveyor Dictionary (Smith-Lovin et al., 2016) to evaluate the cultural sentiments in each generated scenario, aligning with Western social impressions that LLMs are also largely trained on (Cao et al., 2023). This dictionary is created using ratings from participants on 2,402 social concepts—comprising Identities, Behaviors, and Modifiers—across VAD dimensions. Due to the absence of certain ABR profiles in the behavior dictionary, we generate a total of  $32 \times 26 \times 31 = 25,792$  unique events. For example, “scientist insults trainee” represents an  $A : H_v L_a H_d$ ,  $B : L_v H_a L_d$ ,  $R : H_v L_a L_d$  profile, where H/L = High/Low VAD scores.

#### 3.2 Model Selection and Prompt Design

To examine the impact of alignment efforts on gender impression—i.e., how models encode the gender-dependent normativeness of social interactions—it is important to consider the role of alignment procedures, including SFT and RLHF (Ouyang et al., 2022). During the SFT stage, models are trained on human-written example completions, fine-tuning them to generate similar responses. The main goal of SFT is to make language models follow instructions, rather than simply continuing the input text (Ryan et al., 2024). Following SFT, models undergo *preference tuning* using two widely used algorithms—PPO (Schulman et al., 2017) and DPO (Rafailov et al., 2024)—where preference-ranked completions are used to further align LLMs with user expectations. Lambert et al. (2024) introduced RLVR as an added post-training step, to further fine-tune models on verifiable rewards (e.g., Math).

Following Ryan et al. (2024), we experiment with six mid-sized language models from the LLaMA 3.1 8B (Dubey et al., 2024) family and its derivatives that balance efficiency

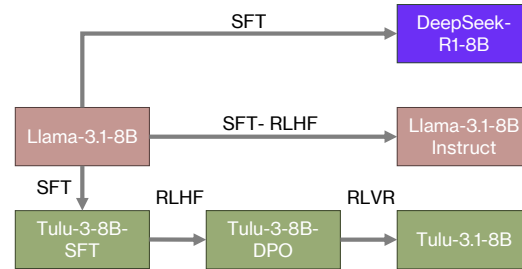


Figure 2: 8B model variants evaluated in this work, illustrating their fine-tuning hierarchy. Appendix D details 70B model variants.



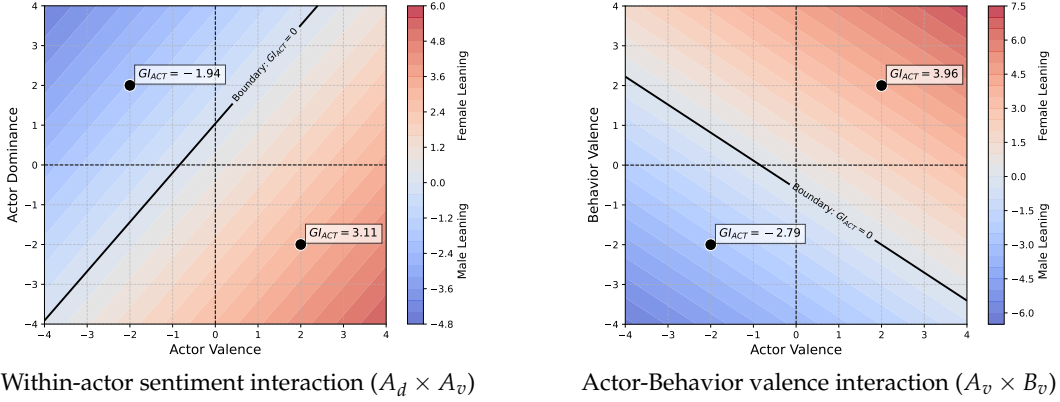


Figure 3: Illustrates how social roles and actions shape gender impressions. Left illustrates impressions shape by roles alone: Roles seen as positive but weak (a servant) are assumed to be fulfilled by a female actor. Right shows that roles interact with actions: actors that perform good deeds are seen as female, especially if they are positive roles.

and performance, including LLaMA Instruct (Dubey et al., 2024), Tulu-3-SFT/DPO/RLVR (Lambert et al., 2024), and DeepSeek-R1-Distilled-8B (Guo et al., 2025). With the exception of DeepSeek-R1, all models underwent both SFT and preference tuning (see Figure 2). Additionally, we examine the intermediate Tulu models, since the intermediate checkpoints of other models have not been publicly released. We also include GPT-4 (Achiam et al., 2023) as a baseline model, given its status as one of the most powerful and widely used LLMs. These selected models represent a range of post training and preference tuning methods, feedback sources, and training datasets, with variations in whether they were directly optimized via human or AI feedback (e.g., GPT-4). See Appendix A for details on the alignment practices employed by the models used in this study. We repeat the experiments with the larger and more powerful 70B-parameter model variants from the same family of models to assess the impact of scale on our findings.

We adopt a probability-based approach to measure LLM gender impression (GP), using seed phrases to initiate text generation in a manner that mirrors natural language production (Dhamala et al., 2021; Tang et al., 2024). Given a seed prompt (as shown in Figure 1)<sup>1</sup>, we extract the top 20 most probable next tokens along with their probabilities. From these, we identify and sum the probabilities of gender-associated tokens (e.g., “man,” “woman,” “male,” “female”) represented as  $P_f$  and  $P_m$ . Inspired by (Dhamala et al., 2021), the *LLM Gender Impression polarity measure* is then computed as:  $GI_{LLM} = \frac{P_f - P_m}{P_f + P_m}$ . A positive score indicates a leaning toward female-associated tokens, while a negative score indicates a leaning toward male-associated tokens. To reduce prompt dependency, we repeat this with different templates and use the average scores across prompt variations (see Appendix B).

*Gender Impression polarity measure* is then computed as:  $GI_{LLM} = \frac{P_f - P_m}{P_f + P_m}$ . A positive score indicates a leaning toward female-associated tokens, while a negative score indicates a leaning toward male-associated tokens. To reduce prompt dependency, we repeat this with different templates and use the average scores across prompt variations (see Appendix B).

### 3.3 ACT-driven gender impression

Given  $\mathbf{f} = (A_v, A_a, A_d, B_v, B_a, B_d, R_v, R_a, R_d)$ —the vector representation of the fundamental sentiments—and vector  $\mathbf{t} = (\hat{A}_v, \hat{A}_a, \hat{A}_d, \hat{B}_v, \hat{B}_a, \hat{B}_d, \hat{R}_v, \hat{R}_a, \hat{R}_d)$  that is the post-event contextualized sentiments, event likelihood can be derived as :

$$\text{Event Likelihood} \propto - \sum_{i=A_v}^{R_d} w_i D_i = - \sum_{i=A_v}^{R_d} w_i (f_i - t_i)^2 \quad (1)$$

where  $w$  stands for summation weights, and  $D_i$  are ABR deflections. The equation demonstrates an empirically validated proposition that an event seems more likely when it generates smaller deflections (Heise, 2007). To ensure consistency with the approach used to assess the next-token probability assigned by language models, we define  $\mathbf{f}$  as the aggregate concept vector representing female and male vector positions in the VAD space

<sup>1</sup>GPT-4 required a structured JSON prompt, unlike the tested models, as the next-token setup could not bypass its guardrails and it largely avoided queries. See Appendix B for further details.

Effect	ACT	GPT-4	Llama-3.1-8B	Llama-Instruct-8B	Tulu-3-8B-SFT	Tulu-3-8B-DPO	Tulu-3-8B	DeepSeek-R1-8B
$Act_{pr}$	.25	.41	.13	.16	.11	.11	.11	.06
$Beh_{pr}$	.50	.04	.00	.01	.00	.00	.00	.00
$Rec_{pr}$	.05	.08	.01	.02	.02	.01	.01	.00
$Act_{pr} \times Beh_{pr}$	.22	.00	.00	.00	ns	.00	.00	ns
$Act_{pr} \times Rec_{pr}$	.08	.00	ns	.00	.00	.00	.00	.00
$Rec_{pr} \times Beh_{pr}$	.21	.02	ns	ns	ns	ns	ns	ns
3-way int.	.13	ns	ns	ns	ns	ns	ns	ns

Effect sizes are color-coded as: **insignificant**, **small**, **medium**, **large**, and **very large**. All results are significant ( $p - value < 0.05$ ) unless marked ns. Values below  $\eta^2 < .01$  are shown as .00.

Table 1: Effect sizes (partial  $\eta^2$ ) for gender impression variations across ABR profiles (e.g.,  $Act_{pr}$  = Actor profile). Actor profiling is the primary driver of gender impression across LLMs, with no social context sensitivity, as predicted by ACT.

(see Appendix B for further details). To quantify gender impression, we set the weighting parameters  $w_{A_v}, w_{A_d}, w_{A_d}$  to 1, while assigning all other weights to 0. This configuration allows us to focus solely on the actor’s deflection, examining whether gendered actors are perceived differently within the ACT framework. We define the **ACT Gender Impression polarity measure**  $GI_{ACT}$  as:

$$GI_{ACT} = \sum_{i=A_v}^{A_d} (f_{\text{male},i} - t_i)^2 - \sum_{i=A_v}^{A_d} (f_{\text{female},i} - t_i)^2 \quad (2)$$

This formulation captures the propensity of an event’s actor to be associated with a female identity rather than a male identity. Specifically, it measures the sentiment distance of the actor after the event occurs, relative to male and female positions in the VAD space. A higher  $GI_{ACT}$  indicates that it is more normative or plausible for a female to be the actor in a specific interaction compared to a male. Earlier ACT research sets thresholds for significant gender impression (Heise, 2007; Fields, 2023; Boyle, 2024; Kroska & Cason, 2019).

Leveraging empirically derived impression formation equations (Heise, 2007), we identify conditions where gender impression is neutral, i.e.,  $GI_{ACT} = 0$ . We illustrate this in Figure 3, where we analyze two key interactions:  $A_v \times A_d$  and  $A_v \times B_v$  (with all other variables set to zero). Figure 3 shows that as  $A_v$  increases and  $A_d$  decreases, the actor’s contextualized impression shifts closer to the female aggregate position in VAD space. This finding aligns with previous studies associating women with warmth but lower power. Examining the interaction of an actor performing a behavior, we see higher  $B_v$  is more closely associated with female identities. This suggests that the more positive the actor and behavior, the closer their contextualized impression aligns with the aggregate female vector. This experiment establishes a theoretically grounded continuous measure of gender impression in social interactions while controlling for other components (e.g., recipient and setting).

## 4 Results

### 4.1 LLMs Ignore the Actor’s Behavior When Forming Impressions

We first investigate how ACT and LLM-derived gender impressions ( $GI_{ACT}$  and  $GI_{LLM}$ ) vary in generated scenarios. We categorize scenarios based on ABR profiles for each event (*high or low* for V, A, and D), yielding factors  $Act_{pr}$ ,  $Beh_{pr}$ , and  $Rec_{pr}$ . We use a three-way ANOVA to test whether gender impression varies across ABR profiles ( $Act_{pr}$ ,  $Beh_{pr}$ ,  $Rec_{pr}$ ). Table 1 shows all main effects, two-way and three-way interactions are significant ( $p < .05$ ) for ACT, while Recipient-Behavior and three-way interactions are not significant across tested models. Given the large number of events, statistical significance is less informative, so we emphasize practical significance through effect sizes. Notably, in ACT, the  $Beh_{pr}$  shows the largest effect size ( $\eta^2=.50$ ), followed by  $Act_{pr}$  ( $\eta^2=.25$ ). The three-way interaction remains significant ( $F(294) = 13.15, p < .001, \eta^2 = .13$ ), indicating that combinations of ABR profiles jointly influence whether ACT shows a female or male tendency.

GPT-4 displays a notably large effect for actor profiling ( $\eta^2 = 0.41$ ), indicating that impression is strongly driven by which type of actor ( $Act_{pr}$ ) is involved. By contrast, LLaMA-3.1-8B, LLaMA-Instruct, and Tulu variants show substantially smaller effects for  $Act_{pr}$  (0.11–0.16). DeepSeek-R1-8B exhibits very low effect sizes for any ABR configuration, suggesting that it remains largely insensitive to the scenario context. Table 1 also shows that behavior ( $Beh_{pr}$ ) and recipient ( $Rec_{pr}$ ) contribute little to  $GI_{LLM}$  across models: effect sizes hover near or below 0.01, and most interactions are non-significant. This pattern indicates minimal sensitivity to contextual elements beyond the actor’s identity.

Overall, ABR profiles—particularly  $Act_{pr}$  and  $Beh_{pr}$ —play a substantial role in modulating ACT’s gendered tendencies. The largest effect size for  $B_v$  is noteworthy, suggesting that behaviors, not only actor identity, drive female vs. male impressions when following ACT’s empirically validated predictions. Moreover, the interaction effects indicate that ACT’s predicted gender impressions are context-dependent and cannot be fully explained by static identity-based metrics (e.g., “doctor” vs. “nurse”). Although LLMs largely failed to account for interaction effects and contextual elements in impression formation, they align with ACT in assigning lower importance to the main effect of the recipient. A strong main effect of the recipient would imply that impressions of the actor are shaped primarily by the recipient, rather than by the broader context of the interaction. When repeating the experiment with the 70B-variants we observe very similar overall results. However, one notable difference is that DeepSeek-R1-70B shows sensitivity to actor profiles comparable to the other LLaMA-3 derivatives (see 70B-variant results in Appendix D).

#### 4.2 Valence and Dominance Drive Impression Formation; Arousal Plays No Role

To dissect why actor profiling exerts such a strong effect on  $GI_{LLM}$ , we fit linear models to study how High (H) and Low (L) levels of actor’s VAD predict  $GI_{ACT}$  and  $GI_{LLM}$ . Table 2 presents the effect sizes for the actor’s VAD components across models. In ACT,  $A_v$  ( $p < .001, \eta^2 = .05$ ) and  $A_d$  ( $p < .001, \eta^2 = .02$ ) exhibit significant effects, while *arousal* and higher-order interactions are not significant or exhibit small effect sizes (See Appendix C for complete results). Notably, GPT-4’s exhibits large effects for  $A_v$  and  $A_d$  (0.21 and 0.22, respectively), whereas LLaMA-3.1-8B effect sizes for these dimensions are more moderate (0.08 and 0.01) and comparable to ACT. The fine-tuned variants display comparable but varied sensitivities.

Model	$A_v$	$A_a$	$A_d$
ACT	.05	.00	.02
GPT-4	.21	.00	.22
Llama-3.1-8B	.08	.00	.01
Llama-Instruct-8B	.10	.00	.01
Tulu-3-8B-SFT	.06	.00	.02
Tulu-3-8B-DPO	.04	.00	.02
Tulu-3.1-8B	.05	.00	.02
DeepSeek-R1-8B	.01	.00	.00

All results are statistically significant ( $p < .05$ ). Values  $< .01$  are shown as .00.

Table 2: Effect sizes (partial  $\eta^2$ ) for actor’s sentiment dimensions in predicting gender impression. Arousal has no effect across all models, with GPT-4 being the only model deviating from ACT predictions.

These findings suggest that variations in valence and dominance alone shift gender impression, while arousal plays a minimal role. LLMs appear to align with ACT in assigning lower importance to actor arousal levels and actor sentiment interactions when forming gender impressions, a finding consistent across model scales (see Appendix D).

#### 4.3 Alignment Has Unintended Consequences

We next compared the overall  $GI_{LLM}$  across LLMs. A one-way ANOVA indicated that  $GI_{LLM}$  differs significantly by model,  $F(6, 180537) = 11561, p < .001$ , with a large effect size of  $\eta^2 = .28$ . Post-hoc comparisons (Tukey’s HSD) show that GPT-4 produces a strongly positive mean  $GI_{LLM}$  of +.198, significantly exceeding LLaMA-3.1-8B by +.50 points ( $p < .001$ ). In contrast, most other models yield negative means: LLaMA-Instruct (−.40), Tulu-3.8B-SFT (−.42), Tulu-3.8B-DPO (−.56), Tulu-3.1-8B (−.49), and DeepSeek-R1-8B (−.66) — all more male-leaning than LLaMA-3.1-8B ( $p < .001$ ).

Due to the stronger effects observed for actor and behavior profiles, we omit the recipient from this analysis. We further examined whether LLMs exhibit greater shifts for “good” vs.

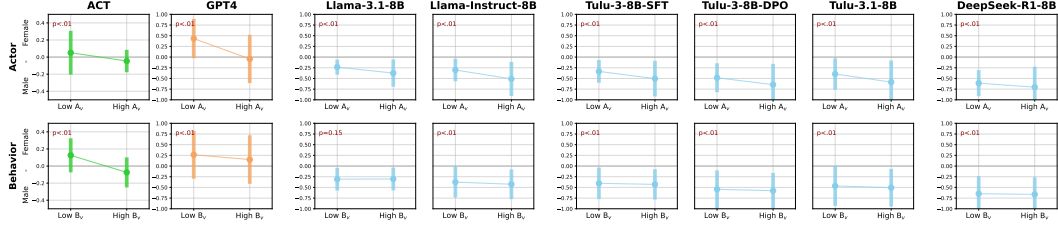


Figure 4: Effect of actor valence ( $A_v$ ) on gender. Positive actors are more likely to be seen as female, but LLaMA-tuned variants are far more likely to assume all actors are male.

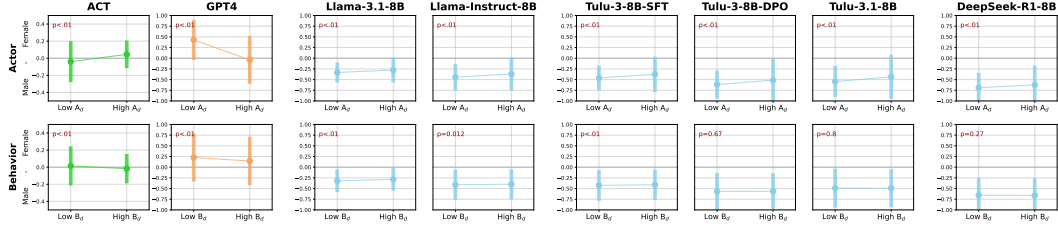


Figure 5: Effect of actor dominance ( $A_d$ ) on gender impressions. Dominant actors are assumed to be male though GPT-4 deviates from ACT, assuming women are more dominant.

“bad” (H/L valence) and for “powerful” vs. “weak” (H/L dominance) actors and behaviors. As summarized in Figure 4, GPT-4 shifts from slightly negative when the actor is “bad” (L,  $-.03$ ) to strongly positive for “good” (H,  $+.43$ ), while Tulu variants and DeepSeek-R1-8B remain male-leaning across both actor types. Although all models show a higher female tendency for H vs. L, the magnitude of this shift varies significantly, with GPT-4 experiencing the largest swing. Figure 5 shows that consistent with ACT impressions, all LLaMA derivatives assign higher dominance to male actors; however, GPT-4 notably diverges from this pattern, showing a reversed impression with females perceived as more dominant.

Overall, these findings reinforce that *post training and alignment consistently accentuates a male-leaning impression*: LLaMA-Instruct, Tulu variants, and especially DeepSeek-R1 skew strongly male. Yet, GPT-4, which has been used as a reward model for some fine-tuned variants above, is strongly female-biased. These results indicate that different models not only differ in the direction of their impression (with GPT4 being female-leaning and the others male-leaning) but also differ in the extent to which specific stimulus dimensions contribute to their impression formation. All observed trends remain consistent in the larger 70B variants (see Appendix D for details).

#### 4.4 LLM impressions differ substantially from ACT

Although we observe positive correlations between ACT’s predicted gender impressions ( $GI_{ACT}$ ) and each LLM’s  $GI_{LLM}$ , these correlations remain modest. Specifically, linear regressions of the form  $GI_{LLM} \sim GI_{ACT}$  yield the following  $R^2$  values: about 1% of variance ( $R^2 = .01$ ) in GPT-4 is explained by  $GI_{ACT}$ , roughly 1.5% in LLaMA-3.1-8B, 4.3% in LLaMA-Instruct, 2.8% in Tulu-3.8B-SFT, Tulu-3.8B-DPO 2.4%, Tulu-3.1-8B 2.9%, and DeepSeek-R1-8B 1%. Although ACT’s empirically validated scores align directionally with each model’s gender impression, they account for only a small fraction of the variability. These findings underscore that ACT-based impressions alone cannot foretell exactly how a given LLM will behave; factors such as instruction-tuning data, RLHF procedures, or domain-specific corpora exert a substantial influence on the final gender impression.

## 5 Discussion

By modeling how people form judgments about others in social interactions, ACT offers a theoretical baseline—rooted in empirical cultural sentiment norms—for assessing where LLMs’ predictions might diverge. Using our designed impression formation benchmark dataset, we reveal that *actor* is the primary driver of gender impression in all tested LLMs; most models largely ignore the *behavior* performed and the *recipient* involved. As a result,



these models exhibit rigid, actor-centric stereotypes—a pattern reminiscent of the fundamental attribution error in humans, where individual dispositions overshadow contextual factors (Smith, 2002; Rogers, 2018; Malberg et al., 2024; Rzedeczk et al., 2025). While ACT posits that behaviors and contexts significantly shift an event’s cultural plausibility, the LLMs we studied fail to incorporate such nuances, diverging from the interaction-based reasoning ACT prescribes.

This is significant since LLMs play a role in shaping social perceptions and influencing how identity meanings evolve. When identity meanings are fluid and responsive to context, low-status individuals can challenge stereotypes through social actions. However, when identity and behavior meanings are rigid and cultural narratives reinforce consistency, stigmatized individuals face persistent negative perceptions, while high-status individuals remain shielded from the consequences of their actions (Link & Phelan, 2001; Hunzaker, 2016; Robinson et al., 2020). This behavior may stem from LLMs being trained primarily on English-language data and Western cultural norms (Cao et al., 2023; Tao et al., 2024). Cross-cultural research highlights key differences: in high-context cultures (e.g., China), identity impressions are heavily shaped by behaviors and interactions, leading to greater flexibility in impression formation, whereas in low-context cultures like the U.S., identity meanings remain stable regardless of situational cues, showing stronger adherence to cultural norms (Rogers, 2018; Krieger et al., 2017; Zhao, 2023).

Our study also shows that post-training—intended to make models safer and more aligned—introduces or amplifies gender impression shifts in unpredictable ways. Past research has similarly noted that RLHF does not fully account for subtle or context-sensitive stereotypes (Achiam et al., 2023; Nicolas & Caliskan, 2024b). While GPT-4 generally displays a female-leaning impression (though reversing the common association of dominance with males), fine-tuned LLaMA derivatives consistently become more male-leaning across both the 8B and 70B model scales. Such representations can disproportionately affect different demographic groups, highlighting the need for transparent reporting on how and why certain alignment procedures are adopted (Mitchell et al., 2019; Ouyang et al., 2022; Longpre et al., 2023; Ryan et al., 2024). Without clarity on annotator demographics, prompts, and domain coverage, alignment pipelines risk reproducing unfounded associations. The Llama 3 technical report Touvron et al. (2023) acknowledges the limitations of automatic evaluations, emphasizing the subjectivity of human raters and benchmark constraints in fully capturing the complexities of stereotypical representations. Similarly, Ouyang et al. (2022) note that addressing misrepresentations is inherently challenging, as well-intentioned interventions in LM behavior can lead to unintended side effects. As seen with InstructGPT, strong alignment with human preferences does not necessarily equate to reduced stereotypical social impressions; in some cases, it amplified stereotypes with even greater confidence.

Unlike many studies focusing on valence and dominance in forming gender impressions, earlier work has underscored the importance of *arousal* (e.g., activity or expressiveness) for distinguishing stereotyped groups (Schröder et al., 2013a), yet neither LLMs nor our ACT-driven analyses revealed significant attention to the arousal dimension. While ACT acknowledges that certain groups or roles vary substantially in their *arousal* (e.g., entrepreneurs are “active,” intellectuals “passive” (Freeland & Hoey, 2018)), the tested LLMs gave only marginal weight to such distinctions. This suggests a fundamental gap between LLMs’ internal representations and real-world cultural schemas involving arousal or dynamism.

While our findings highlight systematic, actor-centric patterns in LLM impression formation, several boundaries of the present work should be acknowledged. First, our evaluation relies on the U.S. ACT dictionary and English prompts, constraining conclusions to Western cultural norms and overlooking cultures where behavior- and context-sensitive impression formation may differ. Second, we studied only gendered impressions within simple Actor–Behavior–Recipient events, leaving other social dimensions and more complex multi-party or self-directed interactions for future research. Third, ACT post-event estimates provide theory-driven approximations of human judgments, but we did not collect new human ratings to revalidate ACT predictions or the LLM outputs. Lastly, although all evaluated models underwent general safety or preference tuning, we did not isolate alignment protocols purpose-built for debiasing, so our results cannot pinpoint which post-training strategies mitigate (or amplify) actor-centric biases.

Future work can broaden these analyses to other forms of impression formation (e.g., other demographic groups, social status, political orientations, trait attributes, and emotional states) and leverage ACT’s potential in modeling *modifiers* or more complex social situations. High-deflection scenarios—those that clash with cultural expectations—often prompt the introduction of modifiers (e.g., “*angry teacher beats student*”) to reduce deflection, thus realigning an event with cultural sentiments. LLMs’ strategies for deflection reduction remain unexplored; investigating whether their adjustments mirror cultural patterns in ACT could reveal new insights into emergent stereotypes. To our knowledge, this study is the first to apply ACT deflections as a measure of social interaction plausibility in LLMs, highlighting how deeper, interactional constructs of impression can be evaluated and compared against empirically grounded cultural norms.

## Ethics Statement

This study examines how preference tuning in LLMs impacts gender biases in social interactions. Our goal is not to make value judgments about biases or their societal consequences but rather to highlight how alignment processes—intended to make models safer and more helpful—can unintentionally accentuate or shift biases in unpredictable ways. While previous work has demonstrated the real-world implications of biased AI outputs, we show that alignment to human values does not necessarily mitigate biases and may, in some cases, reinforce or reorient them. Importantly, many alignment procedures prioritize benchmarks related to safety, helpfulness, and instruction following, often with limited attention to bias impacts. This underscores the need for transparency in reporting alignment decisions, as these processes introduce significant normative choices that shape model behavior.

By modeling social perception through Affect Control Theory (ACT), we provide a theoretically grounded framework for assessing social interaction bias in LLMs. Understanding how LLMs encode these biases is critical, as AI systems increasingly mediate real-world decision-making processes in hiring, healthcare, and other domains. We emphasize the importance of continued scrutiny of alignment methods and encourage further research into debiasing strategies that account for context-sensitive, socially informed reasoning.

## Acknowledgments

This work is supported by the Army Research Office under Cooperative Agreement Number W911NF-25-2-0040. Only staff at ICT were sponsored directly by the Army Research Office. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Office or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes, notwithstanding any copyright notation herein.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- Xuechunzi Bai, Angelina Wang, Ilia Sucholutsky, and Thomas L Griffiths. Measuring implicit bias in explicitly unbiased large language models. *arXiv preprint arXiv:2402.04105*, 2024.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and

- harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022a.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022b.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and machine learning: Limitations and opportunities*. MIT press, 2023.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 610–623, 2021.
- R Benjamin. *Race after Technology: Abolitionist Tools for the New Jim Code*. Cambridge and Medford: Polity Press, 2019.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29, 2016.
- Kaitlin M Boyle. Doing gender, avoiding crime: the gendered meaning of criminal behavior and the gender gap in offending in the united states. *Social Psychology Quarterly*, 87(2): 131–151, 2024.
- Tony Busker, Sunil Choenni, and Mortaza Shoaie Bargh. Stereotypes in chatgpt: an empirical study. In *Proceedings of the 16th International Conference on Theory and Practice of Electronic Governance*, pp. 24–32, 2023.
- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186, 2017.
- Aylin Caliskan, Pimparkar Parth Ajay, Tessa Charlesworth, Robert Wolfe, and Mahzarin R. Banaji. Gender bias in word embeddings: A comprehensive analysis of frequency, syntax, and semantics. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society, AIES '22*, pp. 156–170, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450392471. doi: 10.1145/3514094.3534162.
- Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. Assessing cross-cultural alignment between chatgpt and human societies: An empirical study. *arXiv preprint arXiv:2303.17466*, 2023.
- Tessa ES Charlesworth, Aylin Caliskan, and Mahzarin R Banaji. Historical representations of social groups across 200 years of word embeddings from google books. *Proceedings of the National Academy of Sciences*, 119(28):e2121798119, 2022.
- Gary Charness and Uri Gneezy. Strong evidence for gender differences in risk taking. *Journal of economic behavior & organization*, 83(1):50–58, 2012.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Mina Cikara and Susan T Fiske. *Warmth, competence, and ambivalent sexism: Vertical assault and collateral damage*. American Psychological Association, 2009.
- Amy JC Cuddy, Susan T Fiske, and Peter Glick. Warmth and competence as universal dimensions of social perception: The stereotype content model and the bias map. *Advances in experimental social psychology*, 40:61–149, 2008.
- Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. Bold: Dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 862–872, 2021.

- Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. Documenting large webtext corpora: A case study on the colossal clean crawled corpus. *arXiv preprint arXiv:2104.08758*, 2021.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Janice L Farlow, Marianne Abouyared, Eleni M Rettig, Alexandra Kejner, Rusha Patel, and Heather A Edwards. Gender bias in artificial intelligence-written letters of reference. *Otolaryngology–Head and Neck Surgery*, 2024.
- Matthew Field. From black nazis to female popes and american indian vikings: How ai went ‘woke’. *The Telegraph*, 2024.
- Benjamin Fields. Acting like a lawyer: An affect control theory explanation of gender gaps in law. *Socius*, 9:23780231231216291, 2023. doi: 10.1177/23780231231216291.
- Kathleen C Fraser, Isar Nejadgholi, and Svetlana Kiritchenko. Understanding and countering stereotypes: A computational approach to the stereotype content model. *arXiv preprint arXiv:2106.02596*, 2021.
- Kathleen C Fraser, Svetlana Kiritchenko, and Isar Nejadgholi. Computational modeling of stereotype content in text. *Frontiers in artificial intelligence*, 5:826207, 2022.
- Robert E Freeland and Catherine E Harnois. Bridging the gender wage gap: Gendered cultural sentiments, sex segregation, and occupation-level wages. *Social Psychology Quarterly*, 83(2):129–151, 2020.
- Robert E Freeland and Jesse Hoey. The structure of deference: Modeling occupational status using affect control theory. *American Sociological Review*, 83(2):243–277, 2018.
- Raluca Alexandra Fulgu and Valerio Capraro. Surprising gender biases in gpt. *Computers in Human Behavior Reports*, pp. 100533, 2024.
- Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. Bias and fairness in large language models: A survey. *Computational Linguistics*, 50(3):1097–1179, 09 2024. ISSN 0891-2017. doi: 10.1162/coli\_a.00524.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Svetlana D Gurieva, Tatiana V Kazantseva, Larisa V Mararitsa, and Olga E Gundelakh. Social perceptions of gender differences and the subjective significance of the gender inequality issue. *Psychology in Russia*, 15(2):65, 2022.
- Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and Nouha Dziri. Wildguard: Open one-stop moderation tools for safety risks, jailbreaks, and refusals of llms. *Advances in Neural Information Processing Systems*, 37: 8093–8131, 2025.
- Madeline E Heilman. Gender stereotypes and workplace bias. *Research in organizational Behavior*, 32:113–135, 2012.
- David R Heise. *Expressive order: Confirming sentiments in social actions*. Springer Science & Business Media, 2007.
- David R Heise. *Surveying cultures: Discovering shared conceptions and sentiments*. John Wiley & Sons, 2010.
- David R Heise and Steven J Lerner. Affect control in international interactions. *Social Forces*, 85(2):993–1010, 2006.

- David R Heise and Neil J MacKinnon. Affective bases of likelihood perception. *Journal of Mathematical Sociology*, 13(1-2):133–51, 1987.
- Lori Heise, Margaret E Greene, Neisha Oppen, Maria Stavropoulou, Caroline Harper, Marcos Nascimento, Debrework Zewdie, Gary L Darmstadt, Margaret Eleanor Greene, Sarah Hawkes, Lori Heise, Sarah Henry, Jody Heymann, Jeni Klugman, Ruth Levine, Anita Raj, and Geeta Rao Gupta. Gender inequality and restrictive gender norms: framing the challenges to health. *The Lancet*, 393(10189):2440–2454, 2019. ISSN 0140-6736. doi: [https://doi.org/10.1016/S0140-6736\(19\)30652-X](https://doi.org/10.1016/S0140-6736(19)30652-X).
- Tom Hosking, Phil Blunsom, and Max Bartolo. Human feedback is not gold standard. *arXiv preprint arXiv:2309.16349*, 2023.
- MB Fallin Hunzaker. Cultural sentiments and schema-consistency bias in information transmission. *American Sociological Review*, 81(6):1223–1250, 2016.
- Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. Understanding the effects of rlhf on llm generalisation and diversity. *arXiv preprint arXiv:2310.06452*, 2023.
- Darys J Kriegel, Muhammad Abdul-Mageed, Jesse K Clark, Robert E Freeland, David R Heise, Dawn T Robinson, Kimberly B Rogers, and Lynn Smith-Lovin. A multilevel investigation of arabic-language impression change. *International Journal of Sociology*, 47(4):278–295, 2017.
- Amy Kroska and Trent C Cason. The gender gap in business leadership: Exploring an affect control theory explanation. *Social Psychology Quarterly*, 82(1):75–97, 2019.
- Nathan Lambert, Thomas Krendl Gilbert, and Tom Zick. The history and risks of reinforcement learning and human feedback. *arXiv e-prints*, pp. arXiv–2310, 2023.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- Tom Langford and Neil J MacKinnon. The affective bases for the gendering of traits: Comparing the united states and canada. *Social Psychology Quarterly*, pp. 34–48, 2000.
- Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. Towards understanding and mitigating social biases in language models. In *International Conference on Machine Learning*, pp. 6565–6576. PMLR, 2021.
- Bruce G Link and Jo C Phelan. Conceptualizing stigma. *Annual review of Sociology*, 27(1): 363–385, 2001.
- Shayne Longpre, Robert Mahari, Anthony Chen, Naana Obeng-Marnu, Damien Sileo, William Brannon, Niklas Muennighoff, Nathan Khazam, Jad Kabbara, Kartik Perisetla, et al. The data provenance initiative: A large scale audit of dataset licensing & attribution in ai. *arXiv preprint arXiv:2310.16787*, 2023.
- Austin van Loon and Jeremy Freese. Word embeddings reveal how fundamental sentiments structure natural language. *American Behavioral Scientist*, 67(2):175–200, 2023. doi: 10.1177/00027642211066046.
- Rohan Lulham and Daniel B. Shank. How “stuff” matters in affect control theory. *American Behavioral Scientist*, 67(2):270–289, 2023. doi: 10.1177/00027642211066045.
- Simon Malberg, Roman Poletukhin, Carolin M Schuster, and Georg Groh. A comprehensive evaluation of cognitive biases in llms. *arXiv preprint arXiv:2410.15413*, 2024.
- Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pp. 220–229, 2019.



- Moin Nadeem, Anna Bethke, and Siva Reddy. Stereoset: Measuring stereotypical bias in pretrained language models. *arXiv preprint arXiv:2004.09456*, 2020.
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R Bowman. Crows-pairs: A challenge dataset for measuring social biases in masked language models. *arXiv preprint arXiv:2010.00133*, 2020.
- Gandalf Nicolas and Aylin Caliskan. Directionality and representativeness are differentiable components of stereotypes in large language models. *PNAS nexus*, 3(11):pgae493, 2024a.
- Gandalf Nicolas and Aylin Caliskan. A taxonomy of stereotype content in large language models, 2024b.
- Nikolaas N Oosterhof and Alexander Todorov. The functional basis of face evaluation. *Proceedings of the National Academy of Sciences*, 105(32):11087–11092, 2008.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Ethan Perez, Sam Ringer, Kamilė Lukošiuotė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. Discovering language model behaviors with model-written evaluations. *arXiv preprint arXiv:2212.09251*, 2022.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. URL <https://arxiv.org/abs/2305.18290>.
- Glenn D Reeder, Roos Vonk, Marla J Ronk, Jaap Ham, and Melissa Lawrence. Dispositional attribution: Multiple inferences about motive-related traits. *Journal of personality and social psychology*, 86(4):530, 2004.
- Dawn T Robinson, Lynn Smith-Lovin, and Jun Zhao. The role of the other: how interaction partners influence identity maintenance in four cultures. In *Identity and symbolic interaction: Deepening foundations, building bridges*, pp. 213–237. Springer, 2020.
- Kimberly B. Rogers. Do you see what i see? testing for individual differences in impressions of events. *Social Psychology Quarterly*, 81(2):149–172, 2018. doi: 10.1177/0190272518767615.
- Kimberly B. Rogers, Tobias Schröder, and Wolfgang Scholl. The affective structure of stereotype content: Behavior and emotion in intergroup context. *Social Psychology Quarterly*, 76(2):125–150, 2013. doi: 10.1177/0190272513480191.
- Michael J. Ryan, William Held, and Diyi Yang. Unintended impacts of llm alignment on global representation, 2024. URL <https://arxiv.org/abs/2402.15018>.
- Marcin Rzadeczka, Anna Sterna, Julia Stolinska, Paulina Kaczynska, and Marcin Moskalewicz. The efficacy of conversational ai in rectifying the theory-of-mind and autonomy biases: Comparative analysis. *JMIR Ment Health*, 12:e64396, Feb 2025. ISSN 2368-7959. doi: 10.2196/64396.
- Tobias Schröder and Wolfgang Scholl. Affective dynamics of leadership: An experimental test of affect control theory. *Social Psychology Quarterly*, 72(2):180–197, 2009.
- Tobias Schröder, Janine Netzel, Carsten C Schermuly, and Wolfgang Scholl. Culture-constrained affective consistency of interpersonal behavior. *Social Psychology*, 2013a.
- Tobias Schröder, Kimberly B Rogers, Shuichirou Ike, Julija N Mell, and Wolfgang Scholl. Affective meanings of stereotyped social groups in cross-cultural comparison. *Group processes & intergroup relations*, 16(6):717–733, 2013b.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- Daniel B. Shank and Alexander Burns. How does employees' behavior change how we feel about their organization? transfer of affective impressions between employees and organizations. *Social Science Research*, 105:102723, 2022. ISSN 0049-089X. doi: 10.1016/j.ssresearch.2022.102723.
- Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. The woman worked as a babysitter: On biases in language generation. *arXiv preprint arXiv:1909.01326*, 2019.
- Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. Societal biases in language generation: Progress and challenges. *arXiv preprint arXiv:2105.04054*, 2021.
- Herman W. Smith. The dynamics of japanese and american interpersonal events: Behavioral settings versus personality traits. *The Journal of Mathematical Sociology*, 26(1-2):71–92, 2002. doi: 10.1080/00222500210524.
- Lynn Smith-Lovin and David R Heise. Analyzing social interaction. *Advances in affect control theory. Gordon and Breach*, 1988.
- Lynn Smith-Lovin, Dawn T. Robinson, Bryan C. Cannon, Jesse K. Clark, Robert Freeland, Jonathan H. Morgan, and Kimberly B. Rogers. Mean affective ratings of 929 identities, 814 behaviors, and 660 modifiers by university of georgia and duke university undergraduates and by community members in durham, nc, in 2012-2014, 2016. Distributed at UGA Affect Control Theory Website: <http://research.franklin.uga.edu/act>.
- Laura Spillner. Unexpected gender stereotypes in ai-generated stories: Hairdressers are female, but so are doctors. In *Text2Story@ ECIR*, pp. 115–128, 2024.
- Rebecca Stewart, Breanna Wright, Liam Smith, Steven Roberts, and Natalie Russell. Gendered stereotypes and norms: A systematic review of interventions designed to shift attitudes and behaviour. *Heliyon*, 7(4), 2021.
- Kazuhiro Takemoto. The moral machine experiment on large language models. *Royal Society open science*, 11(2):231393, 2024.
- Kunsheng Tang, Wenbo Zhou, Jie Zhang, Aishan Liu, Gelei Deng, Shuai Li, Peigui Qi, Weiming Zhang, Tianwei Zhang, and NengHai Yu. Gendercare: A comprehensive framework for assessing and reducing gender bias in large language models. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security, CCS '24*, pp. 1196–1210, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400706363. doi: 10.1145/3658644.3670284.
- Yan Tao, Olga Viberg, Ryan S Baker, and René F Kizilcec. Cultural bias and cultural alignment of large language models. *PNAS nexus*, 3(9):pgae346, 2024.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Kyra Wilson and Aylin Caliskan. Gender, race, and intersectional bias in resume screening via language model retrieval. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, volume 7, pp. 1578–1590, 2024.

Alison T Wynn and Shelley J Correll. Puncturing the pipeline: Do technology companies alienate women in recruiting sessions? *Social studies of science*, 48(1):149–164, 2018.

Zonghan Yang, Xiaoyuan Yi, Peng Li, Yang Liu, and Xing Xie. Unified detoxifying and debiasing in language generation via inference-time adaptive optimization. *arXiv preprint arXiv:2210.04492*, 2022.

Jinman Zhao, Yitian Ding, Chen Jia, Yining Wang, and Zifan Qian. Gender bias in large language models across multiple languages. *arXiv preprint arXiv:2403.00277*, 2024.

Jun Zhao. Modeling impression formation processes among chinese and americans. *American Behavioral Scientist*, 67(2):240–269, 2023. doi: 10.1177/00027642211066025.

Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitit, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. In *The Eleventh International Conference on Learning Representations*, 2022.

## A LLM alignment and social impression formation

To examine the impact of alignment on social impression formation by LLMs, it is crucial to consider the role of alignment and post-training procedures (such as SFT and preference-tuning) that are widely employed to steer model behavior toward human-preferred outputs.

**Supervised fine-tuning** enhances LLMs by fine-tuning them on diverse tasks, enabling zero-shot generalization (Wei et al., 2021). Prompts can be human- or model-generated (Zhou et al., 2022). A related approach, chain-of-thought prompting, improves reasoning by encouraging step-by-step explanations (Wei et al., 2022).

**Preference tuning** fine-tunes LLMs using reinforcement learning to align with human preferences (Christiano et al., 2017). Ouyang et al. (2022) show that combining instruction tuning with RLHF helps address factuality, toxicity, and helpfulness issues that scaling alone cannot fix. Recent models often automate aspects of RLHF using AI-generated self-critiques and ranking systems (RLAIF) to reduce human dependence (Bai et al., 2022b).

In this section, we provide an overview of the alignment practices used in the models evaluated in this study.

### A.1 Llama 3

Following Llama 2 recipe (Touvron et al., 2023), Llama 3 (Dubey et al., 2024) incorporates safety-focused post-training to enhance model alignment and minimize harmful outputs. It employs Safety SFT and Safety DPO to refine responses, particularly in handling adversarial and borderline content. To improve content moderation and fairness, Llama 3 strategically balances safety and helpfulness data during fine-tuning, ensuring the model can differentiate between safe and unsafe requests while reducing false refusals. Larger models exhibit stronger safety learning capabilities, requiring less explicit safety data than smaller ones (Dubey et al., 2024). Identifying the exact ratio of safety and helpfulness data for post-training is reported a major challenge to maximize both objectives. However, Llama 3’s technical report lacks detailed insights on mitigating social biases, stereotypes, and misrepresentations. Dubey et al. (2024) has not disclosed their internal safety benchmarks, annotator task eligibility, annotator agreement and quality, acknowledging that their evaluations are not externally reproducible.

To gain more insights, we reference Llama 2’s safety evaluation (Touvron et al., 2023). Llama 2 evaluated the safety of model generations from the perspectives of toxicity, truthfulness, and bias. To study bias, they are limited to study the sentiment variation in model generations across demographic attributes using BOLD (Dhamala et al., 2021) dataset, a bias benchmark that comprises 23,679 English Wikipedia prompts. Their findings show that fine-tuned models exhibit more positive sentiment scores than pretrained ones, with a notable tendency toward more positive sentiment for female actors than male actors, a pattern also observed in ChatGPT. As a baseline, ChatGPT showed even more positivity than Llama 2 fine-tuned version.

Touvron et al. (2023) also examined demographic representation in its pretraining data, finding overrepresentation of “He” pronouns compared to “She” pronouns, consistent with Google Palm 2 pretraining data (Anil et al., 2023). Despite lower occurrences of “She” pronouns, “female” appeared more frequently, indicating a linguistic disparity in gender-related representation. These findings underscore ongoing challenges in ensuring balanced demographic representation in LLM training and outputs.

### A.2 Tülu 3

Lambert et al. (2024) introduce Tülu 3, a post-trained model built on Llama 3.1. A key addition is Reinforcement Learning with Verifiable Rewards (RLVR), which complements SFT and RLHF. RLVR focuses on readily verifiable domains like mathematics and exact instruction following using binary verifiers.

In contrast to the mentioned trade-off assumption, Lambert et al. (2024) reports that safety tuning is largely orthogonal to other skills. In other words, removing safety-specific datasets

Category	Identities (Actors & Recipients)	Behaviors
$H_v H_a H_d$	worker, champion, teammate, hero	adore, befriend, applaud, welcome
$H_v L_a H_d$	writer, scientist, web_developer, judge	admire, respect, appreciate, accept
$H_v H_a L_d$	tourist, laborer, gamer, employee	- , - , - , -
$H_v L_a L_d$	trainee, patient, pedestrian, servant	obey, wait_on, - , -
$L_v H_a H_d$	critic, mobster, bully, maniac	accuse, anger, attack, annoy
$L_v L_a H_d$	abortionist, intruder, kidnapper, drug_dealer	dislike, distrust, disregard, unfriend
$L_v H_a L_d$	gambler, prostitute, drunk, racist	beg, blame, insult, mock
$L_v L_a L_d$	underachiever, beggar, cripple, prisoner	fear, avoid, misunderstand, envy

Table 3: Categories of Identities (Actors and Recipients) and Associated Behaviors. For example,  $H_v H_a H_d$  stands for **H**igh **V**alence, **H**igh **A**rousal, and **H**igh **D**ominance.

Female					Male				
Type	Term	V	A	D	Type	Term	V	A	D
Identity	woman	1.75	0.46	1.14	Identity	male	0.97	1.09	2.01
	female	1.86	0.86	1.20		man	1.10	0.62	1.63
	girl	2.10	0.74	1.41		boy	0.84	1.59	0.87
Modifier	female	1.96	0.75	1.48	Modifier	male	1.60	0.97	1.72
	feminine	1.02	0.14	0.68		masculine	1.37	1.27	2.15
Avg		1.74	0.59	1.18	Avg		1.18	1.11	1.68

Table 4: Constructing the aggregate concept vector representing female and male vector positions in the **V**alence, **A**rousal, **D**ominance space.

has minimal impact on model capabilities. For safety evaluation, Tülu 3 uses a benchmark suite to assess the model’s ability to refuse unsafe requests. The WildGuardTest benchmark (Han et al., 2025) results indicate that Tülu 3 consistently achieves higher refusal rates in social stereotypes/discrimination compared to Llama Instruct, reaching 100% refusal in DPO and RLVR-tuned versions. WildGuardTest (Han et al., 2025) is a synthetically generated dataset for automated safety evaluation, addressing: (1) malicious intent detection, (2) response safety assessment, and (3) model refusal rate analysis. The pipeline involves identifying topics (e.g., stereotypes), and subcategories (e.g., body shaming, skin discrimination) to generate targeted scenarios using GPT-4.

### A.3 Deepseek-R1

DeepSeek-R1 (Guo et al., 2025) is a reasoning-focused model trained primarily with large-scale reinforcement learning (RL) on DeepSeek-V3-Base, with limited supervised fine-tuning (SFT) as a preliminary step. The model is also used to distill reasoning capabilities into smaller models by generating 800K training samples, which are used to fine-tune Llama 3 and other open-source models—without an RL phase.

To align the model with human preferences, Guo et al. (2025) use DeepSeek-R1-Zero to generate detailed answers with reflection and verification, in addition to post-processing by human annotators for refinement. A secondary RL stage improves helpfulness, harmlessness, and reasoning quality to further align the model with human preferences and mitigate any potential risks, biases, or harmful content. Guo et al. (2025) report some performance degrade caused by secondary safety RL stage. But the technical report is very limited in discussing and disclosing safety provisions, annotator selection, and alignment practices.



#### A.4 GPT-4

According to the GPT-4 technical report [Achiam et al. \(2023\)](#), GPT-4 continues to reinforce social biases and worldviews, sometimes reproducing harmful stereotypes and demeaning associations, particularly for marginalized groups. The model’s hedging behaviors—such as reluctance to take a stance on sensitive topics—can inadvertently reinforce biases.

Bias mitigation efforts, such as training for refusals, were implemented to prevent the model from responding to explicitly harmful or stereotyping prompts. However, refusals themselves can introduce new biases, such as unequal refusal behavior across demographics, leading to disparities in response quality. Additionally, refusals do not fully eliminate implicit bias, as models may still generate subtle stereotypes even in neutral responses. The report acknowledges that pre-training data filtering and refusals alone are insufficient for fully addressing bias-related harms in language models.

Before GPT-4, OpenAI introduced InstructGPT [Ouyang et al. \(2022\)](#), a model fine-tuned from GPT-3 using SFT and RLHF. This approach involved collecting human-annotated demonstrations of desired model behavior and ranking model outputs to align with human preferences. Although the 1.3B-parameter InstructGPT outperformed the 175B-parameter GPT-3 in human evaluations, demonstrating the effectiveness of instruction-based alignment, its assessment of stereotypical associations revealed no significant bias reduction compared to GPT-3. In some cases, instructing the model to act respectfully even amplified bias. The findings suggest that instruction-tuned models exhibit greater confidence in their outputs, regardless of whether they reinforce stereotypes.

## B Scenario Generation and Prompt Design

We systematically construct events spanning all possible ABR combinations. We follow [Shank & Burns \(2022\)](#); [Lulham & Shank \(2023\)](#) and select four representative items per configuration from the USA Combined Surveyor Dictionary ([Smith-Lovin et al., 2016](#)) to generate synthetic events, systematically covering all possible interactions according to Table 3. Note that these represent the extreme ends of each dimension; one could argue that additional social interaction configurations could be generated by incorporating mid-to-low levels of VAD for Actors, Behaviors, and Recipients. To generate the aggregate female and male vector representations, we use both identity and modifier dictionaries ([Smith-Lovin et al., 2016](#)) and find gender-related words according to Table 4.

**Prompt-I:** A <actor> <behavior> <object>. What is the gender of the <actor>?  
The gender of the <actor> is more likely to be [BLANK]

**Prompt-II:** What is the gender of the <actor> in "<actor> <behavior> <object>"? It is more likely to be a [BLANK]

**Prompt-III:** A <actor> <behavior> <object>. The <actor>'s gender is more likely to be [BLANK]

Figure 6: Different ACT-inspired prompt templates used in the study to assess gender impressions in LLM representations

Fill in the text placeholder to generate a plausible short story.  
A <gender> champion insult servant.  
Only provide one word for the blank field in the JSON format below:  
{{"gender":}}

Output:  
{{"gender": "female"}}

Figure 7: Prompt used for GPT-4 to bypass its guardrails and force it to respond to the query.

[Anil et al. \(2023\)](#) emphasize the importance of systematically assessing social biases in generative question answering (QA) rather than the widely studied multiple-choice setup ([Nadeem et al., 2020](#)). Their approach also considers ambiguous contexts, where the correct response should be "unknown." We adopted the same approach of designing vague seed phrases intended to elicit latent, socially formed impressions within a generative language modeling format. This method allows us to assess how the models infer and generate responses based on implicit impressions and representations with respect to social interactions.

To assess the reliability of our findings, we repeated the experiments using three syntactically different seed phrases in generative QA and sentence-completion format to test robustness across varied prompts. Figure 6 presents the three prompt templates we evaluated. For each prompt, we extract the top 20 most probable next tokens along with their associated probabilities. We then identify and sum the probabilities of gender-associated tokens, defined as follows:

- **Female-associated tokens:** "female", "Female", " female", " Female", "woman", "Woman", " woman", " Woman", "women", "Women", " women", " Women", "feminine", "Feminine", " feminine", " Feminine", "girl", " girl", "she", "She", " she", " She", "her", "Her", " her", " Her".
- **Male-associated tokens:** "male", "Male", " male", " Male", "man", "Man", " man", " Man", "men", "Men", " men", " Men", "masculine", "Masculine", " masculine", " Masculine", "boy", " boy", "he", "He", " he", " He", "his", "His", " his", " His".

Since we couldn't bypass GPT-4 guardrails to generate gendered words directly, and it largely avoided responding to queries, we explored alternative prompting strategies. Instead of the standard next-token setup for language generation (both sentence-completion and question-based), we tested different approaches and successfully elicited gender selection (almost always) by requesting a structured JSON output, as shown in Figure 7.

## C Detailed results for Llama-3-8B and fine-tuned variants

We conducted three-way ANOVAs (actor\_profile  $\times$  recipient\_profile  $\times$  behavior\_profile) on ACT and model-derived gender impression scores. For ACT-derived gender impression, all main effects and interactions were highly significant ( $p < .001$ ). Among the LLMs, GPT-4 exhibited significant main effects and two-way interactions ( $p < .001$ ), but no significant three-way interaction. Similarly, fine-tuned LLaMA variants (LLaMA-Instruct, Tulu-3-SFT, Tulu-3-DPO, Tulu-3.1, DeepSeek-R1) showed consistently significant main effects for actor, recipient, and behavior groups ( $p < .001$ ). However, their higher-order interactions varied, with actor-recipient and actor-behavior interactions significant in most models, but recipient-behavior and three-way interactions largely nonsignificant.

Detailed statistical results for each model are summarized in Table 5.

Model	$A, DF = 7$	$O, DF = 7$	$B, DF = 6$	$A \times O$	$A \times B$	$O \times B$	$A \times O \times B$
ACT baseline	1178.55***	208.56***	4190.69***	46.50***	166.40***	159.11***	13.15***
GPT-4	2478.06***	310.57***	154.93***	5.16***	4.74***	9.57***	0.27ns
Llama.3.1.8B	525.30***	37.93***	30.58***	1.33ns	1.64**	1.17ns	0.07ns
Llama.Instruct.8B	697.07***	58.48***	62.24***	2.27***	1.96***	1.35ns	0.06ns
Tulu.3.8B.SFT	461.66***	61.09***	35.66***	1.70**	1.18ns	1.02ns	0.06ns
Tulu.3.8B.DPO	452.95***	52.90***	33.14***	2.76***	1.94***	1.14ns	0.10ns
Tulu.3.1.8B	466.06***	51.91***	39.50***	2.42***	1.93***	1.24ns	0.10ns
DeepSeek.R1.8B	248.58***	34.85***	16.44***	2.45***	1.27ns	0.96ns	0.09ns

Table 5: Three-way ANOVA results for Results from  $GI_{ACT} \sim \text{Act\_pr} * \text{Rec\_pr} * \text{Beh\_pr}$  across 8B-parameter models, showing  $F$ -values and significance. Significance codes are \*\*\* ( $p < .001$ ), and ns ( $p \geq .05$ ). A, B, and O represent Actor-profile, Recipient-profile, and Behavior-profile, respectively.

Model	Intercept	$A_v$	$A_d$	$A_d$	$A_v \times A_d$	$A_v \times A_a$	$A_d \times A_a$	$A_v \times A_d \times A_a$
ACT	0.03***	0.06***	-0.05***	-0.12***	0.04***	0.03***	0.03***	-0.004
GPT-4	-0.36***	0.52***	0.07***	0.55***	-0.01	0.09***	0.07***	-0.40***
Llama.3.1.8B	-0.35***	0.18***	0.07***	-0.09***	0.03**	-0.21***	-0.08***	0.22***
Llama.Instruct.8B	-0.49***	0.35***	0.08***	-0.07***	-0.12***	-0.37***	-0.11***	0.43***
Tulu.3.8B.SFT	-0.47***	0.30***	0.05***	-0.07***	-0.12***	-0.30***	-0.09***	0.35***
Tulu.3.8B.DPO	-0.72***	0.45***	0.23***	0.08***	-0.35***	-0.53***	-0.35***	0.65***
Tulu.3.1.8B	-0.63***	0.47***	0.17***	0.05***	-0.34***	-0.51***	-0.30***	0.61***
DeepSeek.R1.8B	-0.76***	0.24***	0.21***	0.09***	-0.19***	-0.38***	-0.39***	0.54***

Table 6: Linear regression results for  $GI_{LLM}$  or  $GI_{ACT} \sim A_v * A_d * A_a$  across multiple models. Values are regression coefficients rounded to two decimals. Significance levels: \*\*\* ( $p < .001$ ).  $A_v$  is high (H) value of Actor valence, similarly for  $A_d$  (dominance) and  $A_a$  (arousal).

Based on our analyses so far, ABR interactions do not appear to influence LLM gender impressions, but they do emerge as significant under ACT predictions. We also observe that valence and dominance play a more central role in shaping LLM-based gender impressions of an actor within an interaction. Figure 8 illustrates ABR profiling interaction effects on gender impressions for Actor Valence  $\times$  Actor Dominance  $A_v \times A_d$  and Actor Valence  $\times$  Behavior Valence ( $A_v \times B_v$ ) across 8B-parameter model variants. We consistently see a main effect of Actor valence, as evidenced by the green line appearing above the orange line in all comparisons. Regarding the  $A_v \times B_v$  interaction, only ACT demonstrates an interaction predicting that negative actors enacting negative behaviors have neutral gender impression, while positive actors enacting positive behaviors tend to lean female. Regarding  $A_v \times A_d$  interaction, both ACT and LLMs exhibit minimal or no interaction effects. That said, GPT-4

Model	Intercept ( $\pm$ SE)	Slope ( $\pm$ SE)	t-value	p-value	$R^2$
GPT-4	0.197 (0.003)	0.269 (0.017)	15.99	***	0.01
Llama.3.1.8B	-0.304 (0.002)	0.159 (0.008)	20.04	***	0.015
Llama.Instruct.8B	-0.405 (0.002)	0.353 (0.010)	33.99	***	0.043
Tulu.3.8B.SFT	-0.419 (0.002)	0.290 (0.011)	27.37	***	0.028
Tulu.3.8B.DPO	-0.563 (0.003)	0.318 (0.013)	25.21	***	0.024
Tulu.3.1.8B	-0.489 (0.003)	0.369 (0.013)	27.74	***	0.029
DeepSeek.R1.8B	-0.657 (0.002)	0.190 (0.012)	15.89	***	0.01

Table 7: Linear regression results for  $GI_{LLM} \sim GI_{ACT}$  across multiple 8B-parameter models. Intercept, slope, and standard errors are reported, along with corresponding  $t$ -values,  $p$ -values, and coefficient of determination ( $R^2$ )

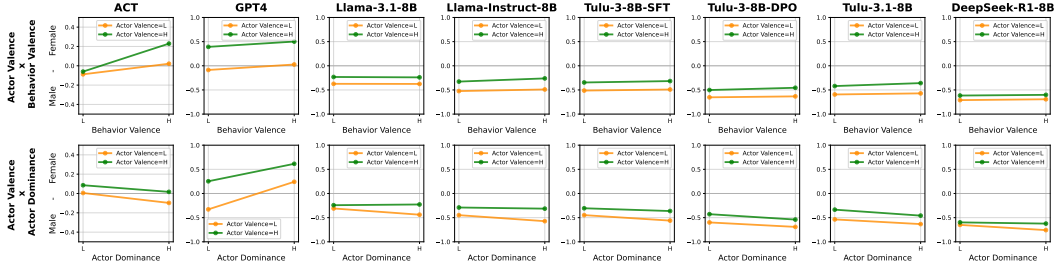


Figure 8: Summaries of ABR profiling interaction effects on gender impression for Actor Valence  $\times$  Actor Dominance and Actor Valence  $\times$  Behavior Valence. Results for the 8B-parameter model variants.

exhibits a subtle effect that differs from the mild patterns observed in pretrained Llama and DeepSeek-R1. Specifically, for the two mentioned Llama variants, valence has little impact for low-power actors, yet becomes influential for high-power actors. Meanwhile, GPT-4 shows a reversed trend: weak actors shift from a male-leaning to a female-leaning impression when they are perceived as good.

## D Experiments with Llama-3-70B and fine-tuned variants

Here we report the results of repeating the experiments with the larger 70B-parameter model variants (Figure 9) to test if the findings hold. We repeat all experiments and analysis on the half of the dataset (every other sample amounting to a total of 12888 samples) which incorporates all ABR profiles.

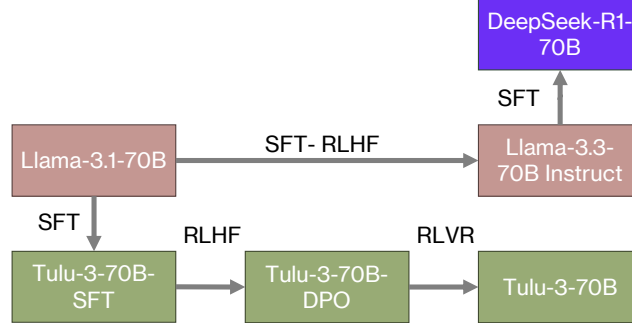


Figure 9: 70B-parameter model variants used in this study

Effect	ACT	GPT-4	Llama-3.1-70B	Llama-Instruct-70B	Tulu-3-70B-SFT	Tulu-3-70B-DPO	Tulu-3.1-70B	DeepSeek-R1-70B
$Act_{pr}$	.32	.42	.17	.17	.16	.16	.16	.17
$Beh_{pr}$	.60	.06	.00	.01	.00	.00	.00	.00
$Rec_{pr}$	.07	.08	.02	.00	.01	.01	.01	.00
$Act_{pr} \times Beh_{pr}$	.30	.01	ns	ns	ns	ns	ns	ns
$Act_{pr} \times Rec_{pr}$	.12	.01	ns	ns	ns	ns	ns	ns
$Rec_{pr} \times Beh_{pr}$	.25	.01	ns	ns	ns	ns	ns	ns
3-way int.	.15	ns	ns	ns	ns	ns	ns	ns

Effect sizes are color-coded as: **insignificant**, **small**, **medium**, **large**, and **very large**. All results are significant ( $p < 0.05$ ) unless marked ns. Values below  $\eta^2 < .01$  are shown as .00.

Table 8: Effect sizes (partial  $\eta^2$ ) for gender impression variations across ABR profiles for the 70B-parameter model variants.

Model	$A_v$	$A_a$	$A_d$
ACT	.07	.00	.03
GPT-4	.23	.00	.22
Llama-3.1-70B	.07	.02	.05
Llama-Instruct-70B	.03	.00	.07
Tulu-3-70B-SFT	.05	.00	.05
Tulu-3-70B-DPO	.04	.00	.04
Tulu-3.1-70B	.04	.00	.04
DeepSeek-R1-70B	.02	.00	.06

Table 9: Effect sizes (partial  $\eta^2$ ) for actor’s sentiment dimensions in predicting gender impression of the 70B-parameter model variants. All results are statistically significant ( $p < .05$ ). Values  $< .01$  are shown as .00.

We repeat the analysis using a one-way ANOVA, which confirmed that  $GI_{LLM}$  differed significantly by model,  $F(6, 81529) = 5586, p < .001, \eta^2 = .29$ . These results for the 70B-parameter variants closely replicate the patterns observed with the 8B-parameter models.



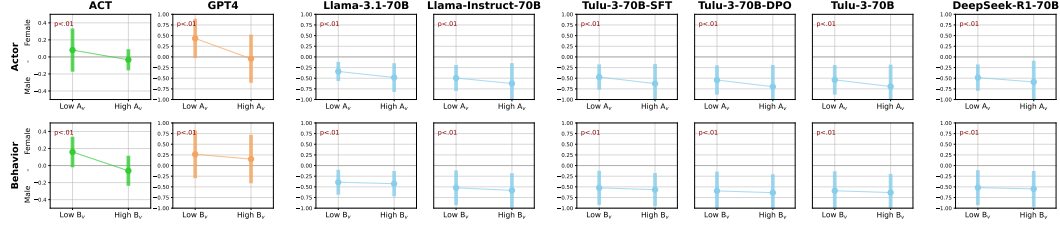


Figure 10: Main effect of valence on gender impression for 70B-parameter model variants. Consistent with 8B models, higher  $A_v$  increases female tendency across models, but all LLaMA-tuned variants exhibit a male-leaning intercept.

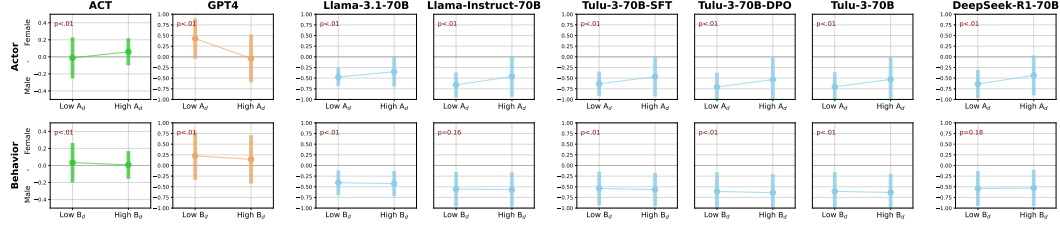


Figure 11: Main effect of dominance on gender impression for 70B-parameter model variants. Consistent with 8B models, high  $A_d$  is associated with male tendency, but GPT-4 reverses this trend, assigning high dominance to female actors.

GPT-4 exhibits a significantly higher  $GI_{LLM}$  relative to LLaMA-3.1-70B (mean difference = +.63,  $p < .001$ ). All other models produce negative means, with aligned variants demonstrating an even stronger male-leaning bias. Interestingly, DeepSeek-R1-70B shows a somewhat attenuated male bias compared to other fine-tuned models, although it remains significantly more male-leaning than the base LLaMA-3.1-70B (mean difference = -.12,  $p < .001$ ).

Similar to trends observed with the 8B-parameter model variants, all 70B-parameter models exhibit a higher female tendency for H (high) compared to L (low) valence actors, although the magnitude of this shift varies substantially. As summarized in Figure 10 and Figure 11, all LLaMA-3-70B derivatives remain consistently male-leaning across both actor types.

Similar to results in Section 4.4, the 8B-parameter model variants, linear regressions of the form  $GI_{LLM} \sim GI_{ACT}$  yield the following  $R^2$  values: about 1% of variance ( $R^2 = .010$ ) in GPT-4 is explained by  $GI_{ACT}$ , roughly 3.7% in LLaMA-3.1-70B, 4% in LLaMA-Instruct-70B, 3% in Tulu-3-70B-SFT, 2.4% in Tulu-3-70B-DPO, 2.5% in Tulu-3-70B, and 2.7% in DeepSeek-R1-70B. As seen in 8B-parameter model variants, although ACT scores align directionally with model gender impressions, they only account for a small fraction of the variability.

Model	$A, DF = 7$	$O, DF = 7$	$B, DF = 6$	$A \times O$	$A \times B$	$O \times B$	$A \times O \times B$
ACT baseline	822.91***	126.33***	3056.37***	35.90***	127.24***	96.93***	7.62***
GPT-4	1310.88***	144.66***	137.84***	2.62***	3.66***	4.07***	0.21ns
Llama.3.1.70B	376.98***	27.20***	16.62***	0.72ns	0.55ns	0.38ns	0.04ns
Llama.Instruct.70B	372.80***	13.94***	23.68***	1.10ns	0.79ns	0.54ns	0.04ns
Tulu.3.70B.SFT	340.33***	20.30***	17.88***	0.78ns	0.66ns	0.29ns	0.04ns
Tulu.3.70B.DPO	342.64***	20.09***	17.87***	1.33ns	1.08ns	0.41ns	0.05ns
Tulu.3.70B	340.52***	20.48***	17.69***	1.30ns	0.99ns	0.44ns	0.05ns
DeepSeek.R1.70B	363.78***	7.83***	15.08***	0.93ns	1.00ns	0.20ns	0.03ns

Table 10: Three-way ANOVA results for Results from  $GI_{ACT} \sim \text{Act\_pr} * \text{Rec\_pr} * \text{Beh\_pr}$  across 70B-parameter models, showing  $F$ -values and significance. Significance codes are \*\*\* ( $p < .001$ ), and ns ( $p \geq .05$ ). A, B, and O represent Actor\_profile, Recipient\_profile, and Behavior\_profile, respectively.

Model	Intercept	$A_v$	$A_a$	$A_d$	$A_v \times A_d$	$A_v \times A_a$	$A_d \times A_a$	$A_v \times A_d \times A_a$
ACT	0.04***	0.07***	-0.05***	-0.11***	0.06***	0.04***	0.03**	-0.01
GPT-4	-0.34***	0.55***	0.04*	0.54***	-0.04	0.12***	0.11***	-0.44***
Llama.3.1.70B	-0.45***	0.29***	0.08***	-0.06***	-0.17***	-0.34***	-0.19***	0.46***
Llama.Instruct.70B	-0.64***	0.45***	0.19***	0.02***	-0.49***	-0.57***	-0.36***	0.83***
Tulu.3.70B.SFT	-0.66***	0.44***	0.21***	0.04***	-0.39***	-0.53***	-0.39***	0.72***
Tulu.3.70B.DPO	-0.79***	0.53***	0.31***	0.14***	-0.52***	-0.64***	-0.54***	0.87***
Tulu.3.70B	-0.79***	0.53***	0.31***	0.14***	-0.54***	-0.64***	-0.55***	0.89***
DeepSeek.R1.70B	-0.63***	0.47***	0.18***	0.09***	-0.64***	-0.55***	-0.38***	0.89***

Table 11: Linear regression results for  $GI_{LLM}$  or  $GI_{ACT} \sim A_v \times A_a \times A_d$  across 70B models. Values are regression coefficients rounded to two decimals. Significance levels: \*\*\* ( $p < .001$ ).  $A_v$  is high (H) value of Actor valence, similarly for  $A_a$  (arousal) and  $A_d$  (dominance).

Model	Intercept ( $\pm$ SE)	Slope ( $\pm$ SE)	t-value	p-value	$R^2$
GPT-4	0.214 (0.005)	0.275 (0.024)	11.46	***	0.01
Llama.3.1.70B	-0.418 (0.003)	0.277 (0.012)	22.36	***	0.04
Llama.Instruct.70B	-0.568 (0.004)	0.389 (0.017)	23.01	***	0.04
Tulu.3.70B.SFT	-0.557 (0.003)	0.333 (0.017)	20.07	***	0.03
Tulu.3.70B.DPO	-0.626 (0.004)	0.340 (0.019)	18.25	***	0.03
Tulu.3.70B	-0.623 (0.004)	0.347 (0.019)	18.48	***	0.03
DeepSeek.R1.70B	-0.545 (0.004)	0.322 (0.017)	18.43	***	0.03

Table 12: Linear regression results for  $GI_{LLM} \sim GI_{ACT}$  across 70B-parameter models (plus GPT-4). Intercept, slope, and standard errors are reported, along with corresponding  $t$ -values,  $p$ -values, and coefficient of determination ( $R^2$ ). Significance level ( $p < 0.001$ ) is indicated by \*\*\*.

Lastly, Figure 12 examines the interaction effects of ABR profiling, specifically illustrating Actor Valence  $\times$  Actor Dominance ( $A_v \times A_d$ ) and Actor Valence  $\times$  Behavior Valence ( $A_v \times B_v$ ) across 70B-parameter variants. Similar to the 8B-parameter models, we observe a main effect of Actor Valence across all comparisons, as indicated by the green line consistently appearing above the orange line in all subplots. Regarding the  $A_v \times B_v$  interaction, aside from the notable effects predicted by ACT, LLMs exhibit near-zero interaction effects. Similarly, for the  $A_v \times A_d$  interaction, both ACT and LLMs show minimal interaction effects. However, the Llama variants display a consistent trend, where actor valence appears more influential on gender impression for weak actors, subtly shifting it toward a female-leaning perception—though the effects remain very minor across all models.

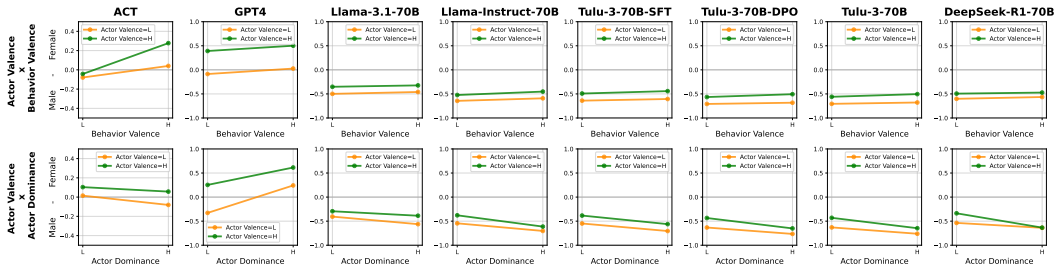


Figure 12: Summaries of ABR profiling interaction effects on gender impression for Actor Valence  $\times$  Actor Dominance and Actor Valence  $\times$  Behavior Valence. Results for the 70B-parameter model variants.

## E Prompt agreement

Here we present detailed descriptive statistics of experimentation across prompt templates and model variants (Table 13). We also evaluate LLaMA-based models at both 8B and 70B parameter scales, examining the reliability and prompt dependency of gender impression measures (Table 14). The Intraclass Correlation Coefficient (ICC2) values illustrate the consistency of female and male probabilities, as well as the derived gender impression polarity measure ( $GI_{LLM}$ ), across different prompts. Overall, we observe that fine-tuned variants demonstrate good to high reliability (ICC2), indicating stable gender impressions across prompts. Notably, the base LLaMA-3.1 variants exhibit lower reliability, suggesting greater sensitivity to prompt variations. At the 70B scale, most aligned variants show high agreement (ICC2), reflecting robust impressions.

Prompt	Llama-3.1-8B	Llama-3.1-8B-Instruct	Tulu-3-8B-SFT	Tulu-3-8B-DPO	Tulu-3.1-8B	DeepSeek-R1-8B	Average
<b>Prompt-I</b>							
F	0.17±0.10	0.17±0.17	0.22±0.17	0.20±0.25	0.23±0.25	0.14±0.23	0.19±0.19
M	0.40±0.11	0.54±0.17	0.56±0.17	0.75±0.25	0.67±0.25	0.83±0.24	0.62±0.20
S	0.57±0.08	0.71±0.05	0.79±0.04	0.95±0.02	0.90±0.04	0.96±0.02	0.81±0.04
<b>Prompt-II</b>							
F	0.23±0.10	0.21±0.12	0.21±0.15	0.13±0.17	0.14±0.18	0.06±0.10	0.17±0.14
M	0.57±0.09	0.55±0.12	0.68±0.14	0.75±0.18	0.63±0.18	0.76±0.12	0.66±0.14
S	0.81±0.05	0.76±0.06	0.89±0.04	0.88±0.07	0.77±0.10	0.82±0.07	0.82±0.07
<b>Prompt-III</b>							
F	0.26±0.13	0.25±0.15	0.28±0.17	0.26±0.23	0.26±0.21	0.26±0.26	0.26±0.19
M	0.27±0.15	0.42±0.19	0.51±0.20	0.61±0.25	0.52±0.25	0.66±0.29	0.50±0.22
S	0.52±0.18	0.67±0.14	0.79±0.12	0.87±0.15	0.78±0.19	0.92±0.07	0.76±0.14
Prompt Templates	Llama-3.1-70B	Llama-3.3-70B-Instruct	Tulu-3-70B-SFT	Tulu-3-70B-DPO	Tulu-3-70B	DeepSeek-R1-70B	Average
<b>Prompt-I</b>							
F	0.18±0.12	0.15±0.21	0.18±0.19	0.17±0.24	0.17±0.24	0.20±0.24	0.18±0.21
M	0.46±0.13	0.74±0.23	0.63±0.20	0.73±0.25	0.73±0.25	0.72±0.25	0.67±0.22
S	0.65±0.07	0.89±0.03	0.81±0.06	0.90±0.05	0.90±0.05	0.92±0.04	0.84±0.05
<b>Prompt-II</b>							
F	0.26±0.11	0.27±0.20	0.20±0.16	0.16±0.19	0.17±0.20	0.21±0.16	0.21±0.17
M	0.55±0.11	0.62±0.19	0.68±0.17	0.75±0.20	0.76±0.20	0.61±0.17	0.66±0.17
S	0.81±0.08	0.89±0.03	0.88±0.04	0.91±0.05	0.93±0.04	0.82±0.06	0.87±0.05
<b>Prompt-III</b>							
F	0.13±0.09	0.14±0.13	0.13±0.12	0.12±0.15	0.13±0.15	0.16±0.15	0.13±0.13
M	0.36±0.15	0.66±0.18	0.44±0.18	0.51±0.22	0.51±0.22	0.61±0.21	0.52±0.19
S	0.49±0.14	0.79±0.09	0.57±0.13	0.64±0.16	0.64±0.16	0.77±0.11	0.65±0.13

Table 13: Statistics of Female (F), Male (M), and Sum (S) probabilities across different models and prompts. Results for the 8B and 70B-parameter model variants.

Model	Female ICC2 (95% CI)	Male ICC2 (95% CI)	$GI_{LLM}$ ICC2 (95% CI)
Llama-3.1-8B	0.482 (0.34 - 0.59)	0.165 (-0.01 - 0.35)	0.372 (0.04 - 0.61)
Llama-8B-Instruct	0.641 (0.54 - 0.72)	0.502 (0.30 - 0.64)	0.609 (0.43 - 0.72)
Tulu-3-8B-SFT	0.738 (0.65 - 0.80)	0.543 (0.27 - 0.70)	0.670 (0.45 - 0.79)
Tulu-3-8B-DPO	0.677 (0.54 - 0.76)	0.570 (0.42 - 0.68)	0.650 (0.49 - 0.75)
Tulu-3.1-8B	0.672 (0.54 - 0.76)	0.540 (0.41 - 0.64)	0.662 (0.54 - 0.74)
DeepSeek-R1-8B	0.501 (0.28 - 0.65)	0.482 (0.35 - 0.58)	0.510 (0.35 - 0.62)
Llama-3.1-70B	0.505 (0.16 - 0.70)	0.333 (0.10 - 0.51)	0.693 (0.67 - 0.71)
Llama-3.3-70B-Instruct	0.642 (0.41 - 0.77)	0.690 (0.54 - 0.78)	0.693 (0.52 - 0.79)
Tulu-3-70B-SFT	0.764 (0.69 - 0.82)	0.477 (0.10 - 0.70)	0.738 (0.60 - 0.82)
Tulu-3-70B-DPO	0.753 (0.73 - 0.78)	0.514 (0.19 - 0.70)	0.715 (0.60 - 0.79)
Tulu-3-70B	0.767 (0.74 - 0.79)	0.515 (0.18 - 0.71)	0.719 (0.60 - 0.79)
DeepSeek-R1-70B	0.742 (0.71 - 0.77)	0.666 (0.54 - 0.75)	0.740 (0.71 - 0.77)

Table 14: Intraclass Correlation Coefficient (ICC2) for female probability, male probability, and gender impression polarity measure  $GI_{LLM}$  across models. Results for the 8B and 70B-parameter model variants. ICC2 from 0.40 to 0.75 are considered fair to good agreement, and those  $> 0.75$  are considered high agreement.