

Rethinking Research on Stereotypes: An Analysis through Social Psychological and Computational Perspectives

Anonymous ACL submission

Abstract

Stereotypes are very harmful social constructs shaping human perception and behavior. Recent work shows that large language models (LLMs) may inherit and amplify such social harms. However, most existing research often focuses on stereotypical biases and overlooks stereotypes and the rich social-psychological literature on them, resulting in resource wastage and slowed progress in stereotype research. We argue that meaningful progress in mitigating stereotypes in LLMs requires tighter integration between social psychology and computational research. To address this gap, we review core social-psychological theories and frameworks and analyze their computational operationalization, highlighting substantial open opportunities. We also analyze computational progress across media narratives, body imaging, and multilingual, multicultural, and multimodal contexts, identifying key gaps and limitations in each domain. We also present a unified analysis of challenges in stereotype research. We further discuss implications for responsible AI, highlighting stereotypes as a root source of downstream harms, and briefly examine the limitations of current mitigation approaches along with potential improvements via explainability and interpretability. We frame stereotypes in AI as socio-technical phenomena and urge further research in responsible AI, informed by the perspectives and future directions presented in this paper.

1 Introduction

Humans are quite good at identifying patterns and forming clusters. They naturally construct conceptual groupings based on shared features, even under uncertainty (Bruner et al., 1956; Shepard et al., 1961), though this process does not always follow strictly logical rules (Rosch, 2024). Humans tend to classify those similar to themselves as the “in-group” and those perceived as different as the

“out-group” (Brewer, 1999; Linville et al., 1989; Mullen et al., 1992; Fiske et al., 2002). Theories such as the similarity–attraction hypothesis (Byrne, 1971) and social identity theory (Tajfel and Turner, 1979) suggest that people are more attracted to others who share similar attitudes, values, and traits (i.e., in-groups). This gives rise to in-group favoritism and can also produce varied emotional responses toward out-groups, such as hate, pity, or respect, as studied by (Brewer, 1999; Fiske et al., 2002; Turner and Reynolds, 2003; Cuddy et al., 2004, 2008). These feelings toward out-groups are gradually translated into thoughts, which then solidify into beliefs: what we call “*stereotypes*”.

The human brain is evolutionarily tuned to respond rapidly to stimuli perceived as critical for survival, such as predators or fire, thereby prioritizing System 1 processing¹ (Harari, 2014; Mobbs et al., 2015; LeDoux, 2012; Wise, 2009). Social competition and interpersonal relationships likewise constitute fundamental survival-relevant contexts (Lakoff, 2024; Rosenbaum, 2014; Griffin, 2001), and therefore tend to elicit fast responses governed by System 1. The origin of stereotypes can thus be attributed to System 1 thinking and the cognitive distinctions between “in-groups” and “out-groups”.

LLMs are increasingly adopted across a wide range of domains, ranging from educational applications such as teaching assistants (Liu et al., 2025) to medical settings, including clinical report generation (Busch et al., 2025), to mention a few; their societal impact continues to expand. Recent work shows that LLMs inherit and sometimes amplify these stereotypes as they learn them from their large-scale pre-training corpora (Pagano et al., 2023; Jeoung et al., 2023; Guo et al., 2024). To mitigate these challenges, research has initially focused on assessing bias in LLMs by exploiting var-

¹System 1 refers to fast, automatic, intuitive, and emotion-driven cognition, in contrast to System 2, which is slower, deliberate, and analytical (Kahneman, 2011).

ious tasks, including Natural Language Generation (NLG) (Nadeem et al., 2021; Felkner et al., 2023), counterfactual reasoning (Nangia et al., 2020; Sahoo et al., 2024), Natural Language Inference (NLI) (Baldini et al., 2023), Question Answering (Parrish et al., 2022; Tomar et al., 2025b), and prompt completion (Gehman et al., 2020; Dhamala et al., 2021). Gallegos et al. (2024) provides a detailed analysis of dataset and metrics designed for assessing bias in LLMs. But fundamentally, bias is a different concept from stereotypes, and confusing biases with stereotypes can give rise to inefficient benchmarks, resulting in substantial resource waste (Shejole and Bhattacharyya, 2025). These concepts are well-studied in social psychology; however, only a few papers draw on social-psychological insights, limiting progress in this domain. Stereotypes are the primary origin of inter-group relations and should therefore be studied separately to understand their effect in Responsible AI. For example, Tomar et al. (2025a) found that incorporating stereotype detection can improve bias detection accuracy. Thus, stereotypes hold considerable potential, which, if systematically explored, can significantly advance Responsible AI research.

To address this gap, this paper focuses primarily on stereotype research and analyzes it from social-psychological and computational perspectives. Our contributions are:

1. A systematic review of social-psychological theories and frameworks on stereotypes that will guide future computational research (§ 2). We also review the computational operationalization of these frameworks and theories, highlighting open opportunities. We analyze computational progress and gaps across domains such as narrative, media, and body imaging, and provide future directions (§ 3).
2. A multimodal, linguistic, and geographic analysis of stereotype research, identifying key gaps and underexplored requirements (§ 4).
3. A unified analysis of challenges in stereotype research by integrating social-psychological and computational perspectives (§ 5).
4. An analysis of implications for Responsible AI, framing stereotypes as foundational to downstream harms, and briefly examining existing mitigation approaches' failures, while suggesting potential improvements through explainability and interpretability (§ 6).

2 Social Psychological Perspectives on Stereotypes

In this section, we review key social psychological theories (§ 2.1) and frameworks (§ 2.2) on the formation, structure, and function of stereotypes.

2.1 Foundational Theories

1. *Similarity–Attraction and Social Identity Theory*: As discussed in the introduction, similarity-attraction theory (Byrne, 1971) and Social Identity Theory (Tajfel and Turner, 1979) posit systematic *in-group* favoritism, whereby individuals favor in-groups over out-groups to enhance self-esteem (Turner and Reynolds, 2003). Self-esteem comprises personal and social identity, the latter derived from group memberships based on attributes such as nationality or age. According to Social Identity Theory, threats to self-esteem intensify in-group favoritism, which in turn restores self-worth, a prediction supported empirically (Ellemers and Haslam, 2012; Postmes and Branscombe, 2010). From this perspective, stereotypes function as mechanisms for self-esteem maintenance, emerging through in-group favoritism and out-group derogation when out-groups are perceived as threatening, thereby conceptualizing stereotypes as *self-esteem protectors*.
2. *Social Role Theory*: This theory (Eagly, 1987) focuses on socialization processes and posits that stereotypes are shaped by the social roles people occupy, such as lower-status versus higher-status jobs. Media plays a direct role in shaping stereotypes, often without individuals being consciously aware of its influence (Ward and Friedman, 2006). In particular, media representations strongly affect body image by promoting stereotypical ideals, such as muscular and lean bodies for males, and fashionable, thin bodies for females (Gauntlett, 2008; Bartlett et al., 2013). Social Role Theory is closely related to Social Learning Theory (Bandura and Walters, 1977), as both emphasize learning through observation and social reinforcement. These theories conceptualize stereotypes as *social representations* representing existing social roles.
3. *Social Categorization Theory*: This theory states that group-based perception is as funda-

181	mental as individual-based perception (Turner et al., 1987). It argues that stereotyping and categorization are the two central components of perception. It states that both the process of stereotyping and the content of stereotypes are fluid and dynamic, varying across social contexts. Social context determines the nature of <i>self–other</i> comparisons and shapes how group boundaries are constructed. It considers that stereotypes reflect the emergent properties of social groups. It conceptualizes stereotypes as <i>psychologically valid representations</i> (Augoustinos and Walker, 1998), grounded in group-based cognition.	231
182		232
183		233
184		234
185		235
186		236
187		237
188		238
189		239
190		240
191		241
192		242
193		243
194		244
195		245
196		246
197		247
198		248
199		249
200		250
201		251
202		252
203		253
204		254
205		255
206		256
207		257
208		258
209		259
210		260
211		261
212		262
213		263
214		264
215		265
216		266
217		267
218		268
219		269
220		270
221		271
222		272
223		273
224		274
225		275
226		276
227		277
228		278
229		279
230		280

281
282
283
284
285
286

287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313

314
315
316
317
318
319
320
321
322
323
324
325

326
327

distinct emotional and behavioral tendencies, ranging from active facilitation to active harm, enabling the SCM to predict real-world social behaviors such as inclusion, neglect, or discrimination (Fiske et al., 2002; Cuddy et al., 2011).

2. *Agency–Beliefs–Communion (ABC) Model:* The ABC model² (Koch et al., 2016) re-frames stereotype content by positing that social perception is fundamentally organized around *Agency* (socioeconomic power) and *Beliefs* (ideological orientation), rather than the warmth-competence dimensions central to the SCM. Developed as a critique of SCM, it challenges its theory-driven structure and reliance on predefined social groups, which may limit the discovery of naturally salient dimensions. Adopting a bottom-up approach, the ABC model shows that *Communion* (including warmth and morality) is not a primary dimension but an emergent construct arising from combinations of Agency and Beliefs. Empirical evidence across multiple studies indicates that spontaneous group categorization aligns most strongly with these two dimensions: Agency shapes power-related judgments, while Beliefs capture ideological alignment. Notably, groups at extreme levels of Agency are perceived as low in communion, whereas moderate Agency is associated with higher communal attributions, suggesting that warmth-based judgments are secondary rather than foundational.
3. *Dual-Perspective Model:* The SCM proposed by Fiske et al. (2002) considers competence as Agency (A) and warmth as Communion (C). Abele et al. (2016) observed that A and C contain multiple components; for example, masculinity (e.g., “assertive” or “decisive”) is also part of Agency, while morality (e.g., “fair,” “honest”) is part of Communion. They proposed a facet model that differentiates A into assertiveness (AA) and competence (AC), and C into warmth (CW) and morality (CM), and reported a good model fit.
4. *Five-Tuple Framework:* Both Davani et al. (2025) and Shejole and Bhattacharyya (2025)

²The terms Agency (A) and Communion (C) were coined by Bakan (1966).

converge on a five-tuple framework for characterizing stereotypes, consisting of the *target group* (T), *relationship characteristics* (R), *associated attributes* (A), the *perceiving group or community* in which the stereotype is held (C), and the *context or time interval* (I) in which it emerges. Both works emphasize that stereotypes are inherently dynamic, varying across social groups and evolving over time, rather than being static representations. This perspective aligns with earlier social psychological theories highlighting the context-dependent and socially constructed nature of stereotyping (Turner et al., 1987). This framework is particularly valuable for computational modeling of stereotypes, as it enables the integration of diverse methodological approaches, such as knowledge graph-based representations, to support structured and systematic analysis.

Table 3 (Appendix C) summarizes these theories and frameworks and Table 2 (Appendix C) contrasts the theoretical assumptions and perceptual mechanisms of the SCM and ABC models.

3 Computational Research on Stereotypes

3.1 Operationalizing Social-Psychological Frameworks

Fraser et al. (2021) computationally operationalized the SCM by deriving warmth and competence directions from lexicon-based word embeddings (Nicolas et al., 2021) and projecting social groups into this space. They also modeled anti-stereotypes³ and validated their findings against survey data. Extending this approach, Fraser et al. (2022) used sentence embeddings and demonstrated strong alignment with human judgments through empirical validation and case studies on gender- and age-related stereotypes. Beyond stereotype measurement, SCM has been applied to assess disability bias (Herold et al., 2022) and bias mitigation (Ungless et al., 2022; Omrani et al., 2023). Cao et al. (2022) operationalized the ABC model as a computational framework to identify group–trait associations in language models, demonstrating moderate alignment with human judgments, supporting intersectional analysis, and evaluating the

³Anti-stereotypes refer to attributes strongly counter to commonly held beliefs about a social group (e.g., football players being weak).

328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347

348
349
350
351

352

353
354

355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373

approach in a U.S.-centric context. These works underscores the multidimensional structure of stereotypes. Building on this view, Fraser et al. (2024) analyzed stereotypes across six psychologically grounded dimensions⁴ for ten occupational groups, showing that while correlations with survey measures vary by dimension, free-text data capture fine-grained and contextually grounded trait associations. Kim and Johnson (2025) extended SCM resources beyond English by constructing and validating a Korean warmth–competence lexicon and a labeled Korean sentence dataset, representing the first SCM-based lexical resource for Korean. There is a need for more research that leverages social-psychological theories and frameworks across multiple languages and cultural contexts.

3.2 Narrative and Media-Based Analyses

As discussed in Section 2, Social Role Theory (Eagly, 1987) posits that media plays a central role in shaping and reinforcing societal stereotypes. A substantial body of work has examined stereotypical portrayals in cartoons, films, and broader media narratives (Schweinitz et al., 2010; Kumar et al., 2022; Xu et al., 2019; Gallego et al., 2025; Shehata, 2020; Atillah et al., 2020; Ji, 2021; Madaan et al., 2017a,b, 2018). More recently, Wang and Lin (2024) used LLMs to extract stereotypes from storytelling content. These studies demonstrate that stereotypes are deeply embedded in media narratives, lending empirical support to Social Role Theory. They further highlight the role of media in amplifying stereotypical beliefs. We believe that greater emphasis should be placed on developing techniques for proactively identifying such stereotypes and assessing their potential social harms before media content is disseminated to the public.

3.3 Body-Image Stereotypes

Body-image stereotypes play a significant role in shaping social norms, although systematic research in this area remains at an early stage. Media representations strongly shape body-image ideals, often reinforcing culturally specific preferences. For example, thin body types are frequently idealized for women in the United States (Lelwica, 2011), whereas medium-sized bodies are more socially preferred in some Middle Eastern contexts (Khalaf et al., 2015; Musaiyer et al., 2000). Such norms

⁴These dimensions were Sociability, Morality, Ability, Assertiveness, Beliefs, and Status.

can generate psychological and behavioral pressure, including the use of weight-altering drugs with potential health risks, highlighting the need for sustained research on body-image stereotypes and their societal consequences. Bias benchmarks such as StereoSet (Nadeem et al., 2021), CrowS-Pairs (Nangia et al., 2020), and BBQ (Parrish et al., 2021) provide limited coverage of body-imaging stereotypes. While they include attributes such as “dark-skinned” or “short,” these representations remain narrow and insufficient to capture the multidimensional nature of body image. Recent efforts such as BISTereo (Asad et al., 2025) advance this line of work by incorporating appearance-related attributes⁵ and using NLI to evaluate bias in LLMs. Automatic modeling of body-image stereotypes from media and narratives remains an important open problem. Future work should quantify body-image bias across LLMs and assess the extent to which their outputs reflect such stereotypes.

4 Analyzing Multimodal, Linguistic and Geographic Coverage

4.1 Multimodal Representations

Stereotypes manifest across multiple modalities, including text, images, video, and audio. However, advances in NLP have led most prior work to focus on textual representations, resulting in a proliferation of text-based benchmarks. More recently, images have received increased attention. Studies such as Fraser and Kiritchenko (2024) reveal substantial gender and racial biases in large vision-language models (VLMs), while Jha et al. (2024) introduce ViSAGE, a dataset evaluating nationality-based stereotypes across 135 countries, showing that stereotypical attributes are nearly three times more likely to appear in generated images and are more offensive for identities from the Global South. A growing body of work (Lee et al., 2025; Zhou et al., 2022; Pang, 2025; Srinivasan and Bisk, 2022; Hamidieh et al., 2024; Malik and Johansson, 2022) further confirms the prevalence of stereotypical biases in VLMs, underscoring a critical and under-explored challenge for multimodal AI. In contrast, research on speech and video remains limited; for example, Kurinec and Weaver III (2021) show that vocal cues alone can activate racial stereotypes. These findings highlight the need for broader investigation into stereotype detection and mitigation in

⁵Skin complexion, body shape, height, attire, hair texture, and eye color.

469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519

conversational audio and video modalities.

4.2 Linguistic and Geographic Coverage

SeeGULL (Jha et al., 2023) and *Visage* (Jha et al., 2024) examine geographic variation in stereotypes, but primarily operationalize geography through nationality. *WinoQueer* (Felkner et al., 2023) focuses on stereotypes related to LGBTQ+ identities, providing a dedicated resource for studying sexual and gender minority representation. Benchmarks such as *EMGSD* (King et al., 2024) and *MGSD* (Zekun et al., 2023), inspired by earlier bias datasets including *StereoSet* (Nadeem et al., 2021) and *CrowS-Pairs* (Nangia et al., 2020), span dimensions such as race, religion, gender, and age. However, they inherit key conceptual limitations, notably ambiguous or inconsistent targets of stereotyping, conflating social groups with non-human or geopolitical entities (e.g., “Norwegian salmon” or “Norway”) and uneven representation of religions (Blodgett et al., 2021). *StereoDetect* (Shejole and Bhattacharyya, 2025) addresses these issues by grounding dataset design in social-psychological distinctions between bias and stereotypes, but remains limited to English and a U.S.-centric context.

These gaps underscore the need for more conceptually grounded and multilingual benchmarks. Recent efforts include datasets for *Korean* (*KoBBQ* (Jin et al., 2024), *KOLD* (Jeong et al., 2022)), *French* (*French-CrowS-Pairs* (Névéol et al., 2022)), *Hindi* (*IndiBias* (Sahoo et al., 2024), *BharatBBQ* (Tomar et al., 2025b)), and *Italian* (*FB-Stereotypes* (Bosco et al., 2023), *QueeroTypes* (Cignarella et al., 2024), *StereoHoax-IT* (Schmeisser-Nieto et al., 2024)). The multilingual *MRHC* dataset (Bourgeade et al., 2023), covering Italian, Spanish, and French, examines racial stereotypes in social media. More recently, *SHADES* (Mitchell et al., 2025) advances the field by curating over 300 stereotypes across 37 regions, translated into 16 languages and annotated with multiple attributes to enable fine-grained multilingual analysis. Despite these efforts, substantial gaps remain in linguistic and cultural inclusion. Global coverage is uneven, with limited resources for many low- to middle-resource languages, including several *Dravidian* and *North-East Indian* languages, as well as Arabic and African languages such as Swahili. Moreover, existing work often underrepresents critical sociocultural dimensions such as caste, region, religion, race, and ethnicity, constraining the representational breadth and

equity of current evaluations. Future research should explicitly incorporate these factors to enable more comprehensive and equitable assessments of LLMs.

5 Challenges in Stereotype Research

5.1 The Problem of Generalization

Social psychological theories, such as Social Role Theory and Social Categorization Theory, clearly require the specification of a social target group for a given stereotype; that is, stereotypes vary depending on the target group under consideration. Consequently, when datasets have limited coverage, any model trained to detect stereotypes will possess knowledge only about those target groups explicitly represented in the training data. Therefore, it is not reliable to use such models to predict stereotypes for unseen target groups, as these models lack the broader social knowledge embedded within a community. Shejole and Bhattacharyya (2025) proposed a solution to this problem through Retrieval-Augmented Generation (RAG). However, extracting context-specific information that is relevant to a particular society and temporal setting remains highly challenging, and the reliability of the sources used also plays a critical role. Future research on more efficient methods for social analysis may contribute to addressing this challenge.

5.2 Annotation and Labeling Challenges

Stereotypes are embedded in a community (Section 2). Therefore, when constructing benchmarks, it is essential to select a representative subset of annotators reflecting the target community. Skewed selection may lead to inefficient or biased benchmarks. Datasets examining more nuanced aspects, such as the effect of language or regional state, as in the case of India or the USA, require a substantial number of annotators, since each state may hold differing perceptions of individuals from other states. Accordingly, annotators must be carefully chosen for each dimension to ensure they appropriately represent the context in which stereotype data is being collected. Obtaining skilled annotators poses a significant challenge. Another important concern relates to labeling quality: annotators may be insufficiently informed or may submit random responses for compensation. Thus, continuous monitoring and guidance of less-informed annotators, as well as the identification and removal of spammers, is necessary to maintain data reliability.

520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568

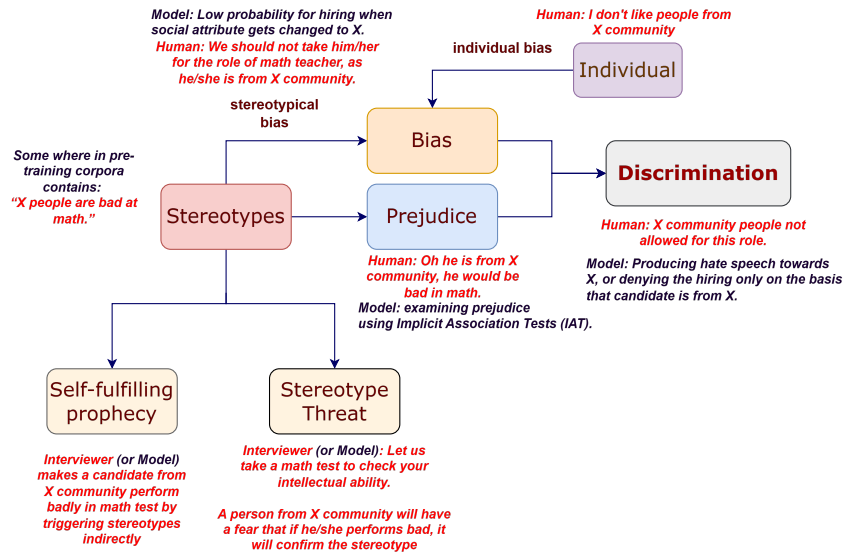


Figure 1: Inter-relationship between of concepts of social psychology and connecting it with Responsible AI scenarios.

5.3 Scalability Constraints

As discussed in previous sections, achieving comprehensive global coverage across languages, cultures, and social dimensions requires substantial, coordinated effort. Ensuring that a language model is globally fair is therefore essential. One possible approach is to evaluate multilingual models separately within each context, as demonstrated in studies such as Singh et al. (2025); Nie et al. (2024); Gamboa et al. (2025). Global representation has consistently posed a significant challenge in research on stereotypes and bias.

5.4 The Dynamic Nature of Stereotypes

Social Categorization Theory and the Five-tuple Framework highlight the fluid and dynamic nature of stereotypes, i.e., that “stereotypes change over time.” Fortunately, this change tends to be gradual and slow. Davani et al. (2025) propose the use of knowledge graphs to model phenomena such as stereotype shifts over time. We further emphasize that stereotype shifts should be systematically studied through efficient modeling approaches, drawing insights from social-psychological theories and frameworks to address this issue.

6 Implications for Responsible AI

6.1 Stereotype is the root cause

Steele and Aronson (1995) showed that black students performed worse on a test framed as measuring intellectual ability due to fear of

confirming negative stereotypes, a phenomenon known as **stereotype threat**. This highlights the risks of LLMs making judgments based on such stereotypes, as observed in computational studies (Shrawgi et al., 2024). Given the serious social-psychological consequences, *AI systems must avoid inheriting the risk of stereotype threat*. Another factor is the role of **confirmation bias** in stereotypes, where people tend to notice information that supports their preconceptions. More concerning is when members of stereotyped groups are led to behave in ways that confirm these stereotypes, a phenomenon called *self-fulfilling prophecies* (Merton, 1948; Jussim, 1986). These occur when a perceiver’s false expectations cause a person to act in ways that confirm them. In Responsible AI, for example in settings where models act as teaching assistants, it is crucial to monitor and prevent self-fulfilling prophecies. If models exhibit implicit bias, stereotypes could trigger these effects, so *models must be both fair and aware of psychological factors to mitigate them*. Bias, prejudice, and discrimination are core components of social harm (see Appendix A). Figure 1 shows the inter-relationship of stereotype with social harms connecting social psychology and computational perspectives. It can be seen that stereotypes are the origin of many problems because of their presence in pre-training corpora. To assess bias, techniques often analyze probability distributions, while prejudice is commonly measured using simulated implicit association tests (see Section 6.2). Discrim-

ination is the most evident, with numerous computational studies demonstrating biased behavior in hiring scenarios (Peña et al., 2025; Anzenberg et al., 2025; Wang et al., 2024; An et al., 2024; Armstrong et al., 2024). Further research is needed to determine whether LLMs and AI models exhibit personal biases similar to humans and to understand the underlying causes.

6.2 Does the Absence of Stereotypical Outputs Imply Fairness?

These questions have been extensively studied in social psychology, where individuals may not explicitly admit bias yet exhibit it in practice. Such bias, termed implicit bias (Greenwald and Krieger, 2006), is a key contributor to prejudice (Kahn, 2017; Payne et al., 2017) and is shaped by automatic cognitive processes, as described in Social Cognition Theory (Section 2). The Implicit Association Test (IAT) (Greenwald et al., 1998, 2009) was developed to measure this phenomenon. Similar tests applied to LLMs (Bai et al., 2024; Mhatre, 2023) reveal that, despite producing non-stereotypical outputs, models may implicitly rely on stereotypes, indicating latent prejudice. It highlights the importance of social-psychology for uncovering hidden prejudice in LLMs.

6.3 Mitigation, Interpretability, and Explainability

We provide a brief analysis for failure of bias mitigation strategies in Appendix D. From a social-psychological perspective, most mitigation strategies target explicit social harms, yet addressing implicit model biases remains essential (Section 6.2). Evidence that anti-stereotypes reduce human prejudice (Cuddy et al., 2008; Fraser et al., 2021) suggests their promise for future stereotype mitigation in LLMs. In Responsible AI, explainability and interpretability techniques offer promising directions for addressing a wide range of challenges. Recent studies, such as work on attention-head pruning (Yang et al., 2025; Ma et al., 2025; Zayed et al., 2024; Hossain et al., 2025), show that selectively modifying internal components of LLMs can reduce bias to some extent. These approaches can be promising for identifying stereotype subspaces in LLMs, namely regions of the parameter space that contains the knowledge of stereotypes prevalent in society. Interpretability methods can play an important role in locating and characterizing these subspaces. In parallel, explainability

techniques such as SHAP (Lundberg, 2017) and LIME (Ribeiro et al., 2016) can be used to analyze the attributions produced by stereotype detectors. These attributions can be analyzed through established social-psychological theories, enhancing theoretical rigor and interpretability in stereotype research. Future work could investigate how modifying stereotype-related subspaces impacts other harms and model’s original efficiency contributing to the transparency of LLMs.

7 Conclusion

Stereotypes have been extensively studied in social psychology; however, computational research has yet to fully leverage this body of knowledge. In this paper, we first reviewed key social-psychological theories and frameworks on stereotype formation and persistence, and examined how they have been operationalized computationally, highlighting that existing work has only scratched the surface and that substantial opportunities remain for deeper computational engagement with these theories. We also analyzed computational progress across media narratives, body imaging, and multilingual, multi-cultural, and multimodal contexts, identifying key gaps and limitations in each domain. We presented a unified analysis of challenges in stereotype research by jointly considering social-psychological and computational perspectives. Finally, we discussed implications for responsible AI, positioning stereotypes as a root cause of downstream harms, connecting them to broader social-psychological constructs, and examining their impact from both AI model and human perspectives. We also briefly reflected on the failures of existing bias mitigation approaches and highlighted some points on how explainability and interpretability techniques can help in solving these issues. We position stereotypes in AI as socio-technical phenomena and argue for a reframing of how responsible AI research conceptualizes and addresses stereotype-related harms. We contend that advancing fairness and reducing social harms in responsible AI requires a shift in perspective. We summarize future research directions discussed in this paper in Table 1 (Appendix B). By grounding future computational research in established social-psychological underpinnings and by pursuing the future research directions outlined in this paper, responsible AI systems can move toward more principled, culturally grounded, and effective interventions.

730 Limitations

731 This paper integrates insights from social psychol-
732 ogy and computational research to provide a com-
733 prehensive view of stereotyping in large language
734 models (LLMs), but several limitations should
735 be noted. First, our focus on combining social-
736 psychological and computational perspectives may
737 limit discussion of other relevant factors, such as
738 technical optimization or purely algorithmic in-
739 terventions, which are beyond the scope of this
740 work. Second, although we review computational
741 progress across multimodal, linguistic, and cultural
742 domains, practical challenges remain. Achieving
743 global inclusivity requires substantial resources and
744 skilled annotators, which can constrain scalability
745 and coverage. While we suggest potential strate-
746 gies, such as modeling contexts separately, these ap-
747 proaches remain aspirational. Overall, our analysis
748 highlights the importance of a joint computational
749 and social-psychological perspective for ground-
750 ing stereotype evaluation in linguistic, social, and
751 historical contexts. Future work should continue
752 bridging these perspectives while addressing prac-
753 tical constraints in data collection, annotation, and
754 model design.

755 References

756 Andrea E Abele, Nicole Hauke, Kim Peters, Eva
757 Louvet, Aleksandra Szymkow, and Yanping Duan.
758 2016. Facets of the fundamental content dimen-
759 sions: Agency with competence and assertive-
760 ness—communion with warmth and morality. *Frontiers in psychology*, 7:1810.
761

762 Gordon W. Allport. 1954. *The Nature of Prejudice*.
763 Addison-Wesley, Reading, MA.

764 Haozhe An, Christabel Acquaye, Colin Wang, Zongxia
765 Li, and Rachel Rudinger. 2024. Do large language
766 models discriminate in hiring decisions on the ba-
767 sis of race, ethnicity, and gender? *arXiv preprint*
768 *arXiv:2406.10486*.

769 Eitan Anzenberg, Arunava Samajpati, Sivasankaran
770 Chandrasekar, and Varun Kacholia. 2025. Evaluating
771 the promise and pitfalls of llms in hiring decisions.
772 *arXiv preprint arXiv:2507.02087*.

773 Lena Armstrong, Abbey Liu, Stephen MacNeil, and
774 Danaë Metaxa. 2024. The silicon ceiling: Auditing
775 gpt’s race and gender biases in hiring. In *Proceedings*
776 *of the 4th ACM Conference on Equity and Access in*
777 *Algorithms, Mechanisms, and Optimization*, pages
778 1–18.

Kenneth J. Arrow. 1973. The theory of discrimination. 779
In Orley Ashenfelter and Albert Rees, editors, *Dis-* 780
crimination in Labor Markets, pages 3–33. Princeton 781
University Press, Princeton, NJ. 782

Narjis Asad, Nihar Ranjan Sahoo, Rudra Murthy, 783
Swaprava Nath, and Pushpak Bhattacharyya. 2025. 784
“you are beautiful, body image stereotypes are ugly!” 785
BIStereo: A benchmark to measure body image 786
stereotypes in language models. In *Findings of* 787
the Association for Computational Linguistics: ACL 788
2025, pages 24471–24496, Vienna, Austria. Associa- 789
tion for Computational Linguistics. 790

Widya Atillah, M Bahri Arifin, and Nita Maya 791
Valiantien. 2020. An analysis of stereotype in 792
zootopia movie. *Ilmu Budaya: Jurnal Bahasa, Sas-* 793
tra, Seni, dan Budaya, 4(1):49–62. 794

Martha Augoustinos and Iain Walker. 1998. The con- 795
struction of stereotypes within social psychology: 796
From social cognition to ideology. *Theory & Psy-* 797
chology, 8(5):629–652. 798

Xuechunzi Bai, Angelina Wang, Ilia Sucholutsky, and 799
Thomas L Griffiths. 2024. Measuring implicit bias 800
in explicitly unbiased large language models. *arXiv* 801
preprint arXiv:2402.04105. 802

David Bakan. 1966. The duality of human existence: 803
An essay on psychology and religion. 804

Ioana Baldini, Chhavi Yadav, Manish Nagireddy, 805
Payel Das, and Kush R Varshney. 2023. Keeping 806
up with the language models: Systematic bench- 807
mark extension for bias auditing. *arXiv preprint* 808
arXiv:2305.12620. 809

Mahzarin R Banaji. 2002. Stereotypes, social psychol- 810
ogy of. *International encyclopedia of the social and* 811
behavioral sciences, pages 15100–15104. 812

Albert Bandura and Richard H Walters. 1977. *Social* 813
learning theory, volume 1. Prentice hall Englewood 814
Cliffs, NJ. 815

Djurdja Bartlett, Agnes Rocamora, and Shaun Cole. 816
2013. Fashion media. 817

Gary S. Becker. 1957. *The Economics of Discrimina-* 818
tion. University of Chicago Press, Chicago. 819

Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, 820
Robert Sim, and Hanna Wallach. 2021. *Stereotyping* 821
Norwegian salmon: An inventory of pitfalls in fair- 822
ness benchmark datasets. In *Proceedings of the 59th* 823
Annual Meeting of the Association for Computational 824
Linguistics and the 11th International Joint Confer- 825
ence on Natural Language Processing (Volume 1: Long 826
Papers), pages 1004–1015, Online. Association 827
for Computational Linguistics. 828

Galen V. Bodenhausen and Jennifer A. Richeson. 2010. 829
Prejudice, stereotyping, and discrimination. In Roy F. 830
Baumeister and Eli J. Finkel, editors, *Advanced So-* 831
cial Psychology: The State of the Science, pages 832
350–380. Oxford University Press. 833

834	Cristina Bosco, Viviana Patti, Simona Frenda, Alessandra Teresa Cignarella, Marinella Paciello, and Francesca D'Errico. 2023. Detecting racial stereotypes: An italian social media corpus where psychology meets nlp. <i>Information Processing & Management</i> , 60(1):103118.	Amy J. C. Cuddy, Susan T. Fiske, and Peter Glick. 2011. Stereotype content model across cultures: Towards universal similarities and some differences . <i>British Journal of Social Psychology</i> , 50(3):472–486.	889
835			890
836			891
837			892
838			
839			
840	Tom Bourgeade, Alessandra Teresa Cignarella, Simona Frenda, Mario Laurent, Wolfgang Schmeisser-Nieto, Farah Benamara, Cristina Bosco, Véronique Moriceau, Viviana Patti, and Mariona Taulé. 2023. A multilingual dataset of racial stereotypes in social media conversational threads. In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pages 686–696.	Amy JC Cuddy, Susan T Fiske, and Peter Glick. 2004. When professionals become mothers, warmth doesn't cut the ice. <i>Journal of Social issues</i> , 60(4):701–718.	893
841			894
842			895
843			
844			
845			
846			
847			
848	Marilynn B Brewer. 1999. The psychology of prejudice: Ingroup love and outgroup hate? <i>Journal of social issues</i> , 55(3):429–444.	Amy JC Cuddy, Susan T Fiske, and Peter Glick. 2008. Warmth and competence as universal dimensions of social perception: The stereotype content model and the bias map. <i>Advances in experimental social psychology</i> , 40:61–149.	896
849			897
850			898
851	Jerome S Bruner, Jacqueline J Goodnow, and A George. 1956. Austin. a study of thinking. <i>New York: John Wiley & Sons, Inc</i> , 14:330.	Aida Mostafazadeh Davani, Sunipa Dev, Héctor Pérez-Urbina, and Vinodkumar Prabhakaran. 2025. A comprehensive framework to operationalize social stereotypes for responsible ai evaluations. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 30018–30031.	901
852			902
853			903
854	Felix Busch, Lena Hoffmann, Daniel Pinto Dos Santos, Marcus R Makowski, Luca Saba, Philipp Prucker, Martin Hadamitzky, Nassir Navab, Jakob Nikolas Kather, Daniel Truhn, and 1 others. 2025. Large language models for structured reporting in radiology: past, present, and future. <i>European Radiology</i> , 35(5):2589–2602.	Patricia G. Devine. 1989. Stereotypes and prejudice: Their automatic and controlled components . <i>Journal of Personality and Social Psychology</i> , 56(1):5–18.	904
855			905
856			906
857			
858			
859			
860			
861	D Byrne. 1971. The attraction paradigm academic press. <i>New York, NY, USA</i> .	Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. 2021. Bold: Dataset and metrics for measuring biases in open-ended language generation. In <i>Proceedings of the 2021 ACM conference on fairness, accountability, and transparency</i> , pages 862–872.	910
862			911
863	Yang Trista Cao, Anna Sotnikova, Hal Daumé III, Rachel Rudinger, and Linda Zou. 2022. Theory-grounded measurement of us social stereotypes in english language models. <i>arXiv preprint arXiv:2206.11684</i> .	John F. Dovidio, Miles Hewstone, Peter Glick, and Victoria M. Esses. 2010. Prejudice, stereotyping and discrimination: Theoretical and empirical overview. In John F. Dovidio, Miles Hewstone, Peter Glick, and Victoria M. Esses, editors, <i>The SAGE Handbook of Prejudice, Stereotyping and Discrimination</i> , pages 3–28. SAGE Publications.	912
864			913
865			914
866			915
867			916
868	Anna Carastathis. 2014. The concept of intersectionality in feminist theory. <i>Philosophy compass</i> , 9(5):304–314.	Alice H. Eagly. 1987. <i>Sex Differences in Social Behavior: A Social-role Interpretation</i> . Lawrence Erlbaum Associates, Hillsdale, NJ.	917
869			918
870			919
871	Sumi Cho, Kimberlé Williams Crenshaw, and Leslie McCall. 2013. Toward a field of intersectionality studies: Theory, applications, and praxis. <i>Signs: Journal of women in culture and society</i> , 38(4):785–810.	Derek Edwards. 1991. Categories are for talking: On the cognitive and discursive bases of categorization. <i>Theory & psychology</i> , 1(4):515–542.	920
872			921
873			922
874			923
875			
876	Alessandra Teresa Cignarella, Manuela Sanguinetti, Simona Frenda, Andrea Marra, Cristina Bosco, and Valerio Basile. 2024. Queereotypes: A multi-source italian corpus of stereotypes towards lgbtqia+ community members. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 13429–13441.	Naomi Ellemers and S Alexander Haslam. 2012. Social identity theory. <i>Handbook of theories of social psychology</i> , 2:379–398.	924
877			925
878			926
879			
880			
881			
882			
883			
884	Kimberlé Crenshaw. 2013. Demarginalizing the intersection of race and sex: A black feminist critique of antidiscrimination doctrine, feminist theory and antiracist politics. In <i>Feminist legal theories</i> , pages 23–51. Routledge.	Virginia K Felkner, Ho-Chun Herbert Chang, Eugene Jang, and Jonathan May. 2023. Winoqueer: A community-in-the-loop benchmark for anti-lgbtq+ bias in large language models. <i>arXiv preprint arXiv:2306.15087</i> .	933
885			934
886			935
887			936
888			937
		Alan P Fiske and Nick Haslam. 1996. Social cognition is thinking about relationships. <i>Current directions in psychological science</i> , 5(5):143–148.	938
			939
			940

941	Susan T Fiske. 1992. Thinking is for doing: portraits of social cognition from daguerreotype to laser-photo. <i>Journal of personality and social psychology</i> , 63(6):877.	995
942		996
943		997
944		998
945	Susan T. Fiske, Amy J. C. Cuddy, Peter Glick, and Jun Xu. 2002. A model of (often mixed) stereotype content: Competence and warmth respectively follow from perceived status and competition. <i>Journal of Personality and Social Psychology</i> , 82(6):878–902.	999
946		1000
947		
948		1001
949		1002
950	Susan T Fiske and 1 others. 1993. Social cognition and social perception. <i>Annual review of psychology</i> , 44(1):155–194.	1003
951		
952		
953	Susan T Tufts Fiske and Shelley E Taylor. 2020. Social cognition: From brains to culture.	1004
954		1005
955	Kathleen Fraser and Svetlana Kiritchenko. 2024. Examining gender and racial bias in large vision–language models using a novel dataset of parallel images. In <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 690–713, St. Julian’s, Malta. Association for Computational Linguistics.	1006
956		1007
957		1008
958		1009
959		1010
960		1011
961		1012
962		1013
963	Kathleen Fraser, Svetlana Kiritchenko, and Isar Nejadgholi. 2024. How does stereotype content differ across data sources? In <i>Proceedings of the 13th Joint Conference on Lexical and Computational Semantics (*SEM 2024)</i> , pages 18–34, Mexico City, Mexico. Association for Computational Linguistics.	1014
964		1015
965		
966		
967		
968		
969	Kathleen C Fraser, Svetlana Kiritchenko, and Isar Nejadgholi. 2022. Computational modeling of stereotype content in text. <i>Frontiers in artificial intelligence</i> , 5:826207.	1016
970		1017
971		1018
972		1019
973	Kathleen C Fraser, Isar Nejadgholi, and Svetlana Kiritchenko. 2021. Understanding and countering stereotypes: A computational approach to the stereotype content model. <i>arXiv preprint arXiv:2106.02596</i> .	1020
974		
975		
976		
977		
978	Ana Guadalupe Gallego, Camino Ferreira, and Ana Rosa Arias-Gago. 2025. Stereotyped representations of disability in film and television: A critical review of narrative media. <i>Disabilities</i> , 5(4):1–25.	1021
979		1022
980		1023
981		1024
982	Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2024. Bias and fairness in large language models: A survey. <i>Computational Linguistics</i> , 50(3):1097–1179.	1025
983		1026
984		
985		
986		
987		
988	Lance Calvin Lim Gamboa, Yue Feng, and Mark Lee. 2025. Social bias in multilingual language models: A survey. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 27845–27868.	1027
989		1028
990		
991		
992		
993	David Gauntlett. 2008. <i>Media, gender and identity: An introduction</i> . Routledge.	1029
994		1030
		1031
		1032
		1033
		1034
		1035
		1036
		1037
		1038
		1039
		1040
		1041
		1042
		1043
		1044
		1045
		1046
		1047

1048	Sullam Jeoung, Yubin Ge, and Jana Diesner. 2023.	1103
1049	StereoMap: Quantifying the awareness of human-	1104
1050	like stereotypes in large language models. In <i>Pro-</i>	
1051	<i>ceedings of the 2023 Conference on Empirical Meth-</i>	1105
1052	<i>ods in Natural Language Processing</i> , pages 12236–	1106
1053	12256, Singapore. Association for Computational	1107
1054	Linguistics.	1108
		1109
1055	Akshita Jha, Aida Mostafazadeh Davani, Chandan K	
1056	Reddy, Shachi Dave, Vinodkumar Prabhakaran, and	1110
1057	Sunipa Dev. 2023. SeeGULL: A stereotype bench-	1111
1058	mark with broad geo-cultural coverage leveraging	1112
1059	generative models. In <i>Proceedings of the 61st Annual</i>	1113
1060	<i>Meeting of the Association for Computational Lin-</i>	1114
1061	<i>guistics (Volume 1: Long Papers)</i> , pages 9851–9870,	
1062	Toronto, Canada. Association for Computational Lin-	1115
1063	guistics.	1116
		1117
1064	Akshita Jha, Vinodkumar Prabhakaran, Remi Denton,	1118
1065	Sarah Laszlo, Shachi Dave, Rida Qadri, Chandan	1119
1066	Reddy, and Sunipa Dev. 2024. Visage: A global-	1120
1067	scale analysis of visual stereotypes in text-to-image	
1068	generation. In <i>Proceedings of the 62nd Annual Meet-</i>	1121
1069	<i>ing of the Association for Computational Linguistics</i>	1122
1070	<i>(Volume 1: Long Papers)</i> , pages 12333–12347.	1123
		1124
1071	Jiaxin Ji. 2021. Analysis of gender stereotypes in disney	1125
1072	female characters. In <i>2021 3rd International Con-</i>	
1073	<i>ference on Literature, Art and Human Development</i>	1126
1074	<i>(ICLAHD 2021)</i> , pages 451–454. Atlantis Press.	1127
		1128
1075	Jiho Jin, Jiseon Kim, Nayeon Lee, Haneul Yoo, Al-	1129
1076	ice Oh, and Hwaran Lee. 2024. Kobbq: Korean	
1077	bias benchmark for question answering. <i>Transac-</i>	1130
1078	<i>tions of the Association for Computational Linguis-</i>	1131
1079	<i>tics</i> , 12:507–524.	1132
1080	John T Jost. 2019. A quarter century of system justi-	1133
1081	fication theory: Questions, answers, criticisms, and	1134
1082	societal applications. <i>British Journal of Social Psy-</i>	
1083	<i>chology</i> , 58(2):263–314.	1135
		1136
1084	John T Jost, Mahzarin R Banaji, and Brian A Nosek.	1137
1085	2004. A decade of system justification theory: Ac-	1138
1086	cumulated evidence of conscious and unconscious	
1087	bolstering of the status quo. <i>Political psychology</i> ,	1139
1088	25(6):881–919.	1140
1089	John T Jost and Jojanneke Van der Toorn. 2012. System	1141
1090	justification theory. <i>Handbook of theories of social</i>	1142
1091	<i>psychology</i> , 2:313–343.	1143
		1144
1092	Lee Jussim. 1986. Self-fulfilling prophecies: A theo-	1145
1093	retical and integrative review. <i>Psychological review</i> ,	
1094	93(4):429.	1146
		1147
1095	Jonathan Kahn. 2017. Pills for prejudice: implicit bias	1148
1096	and technical fix for racism. <i>American Journal of</i>	1149
1097	<i>Law & Medicine</i> , 43(2-3):263–278.	
		1150
1098	Daniel Kahneman. 2011. Thinking, fast and slow. <i>Far-</i>	1151
1099	<i>rar, Straus and Giroux.</i>	1152
1100	Atika Khalaf, Albert Westergren, Vanja Berggren, Ör-	1153
1101	jan Ekblom, and Hazzaa M. Al-Hazzaa. 2015. Per-	1154
1102	ceived and ideal body image in young women in	1155
	south western saudi arabia. <i>Journal of Obesity</i> ,	1156
	2015(1):697163.	1157
	Michelle YoungJin Kim and Kristen Johnson. 2025. Ko-	
	rean stereotype content model: Translating stereo-	
	types across cultures. In <i>Proceedings of the 3rd</i>	
	<i>Workshop on Cross-Cultural Considerations in NLP</i>	
	<i>(C3NLP 2025)</i> , pages 59–70.	
	Theo King, Zekun Wu, Adriano Koshiyama, Emre	
	Kazim, and Philip Treleaven. 2024. Hearts: A	
	holistic framework for explainable, sustainable and	
	robust text stereotype detection. <i>arXiv preprint</i>	
	<i>arXiv:2409.11579.</i>	
	Alex Koch, Roland Imhoff, Ron Dotsch, Christian	
	Unkelbach, and Hans Alves. 2016. The abc of	
	stereotypes about groups: Agency/socioeconomic	
	success, conservative-progressive beliefs, and com-	
	munion. <i>Journal of Personality and Social Psychol-</i>	
	<i>ogy</i> , 110(5):675–709.	
	Arjun M Kumar, Jasmine YQ Goh, Tiffany HH Tan, and	
	Cynthia SQ Siew. 2022. Gender stereotypes in holly-	
	wood movies and their evolution over time: Insights	
	from network analysis. <i>Big Data and Cognitive Com-</i>	
	<i>puting</i> , 6(2):50.	
	Courtney A Kurinec and Charles A Weaver III. 2021.	
	“sounding black”: Speech stereotypicality activates	
	racial stereotypes and expectations about appearance.	
	<i>Frontiers in psychology</i> , 12:785283.	
	George Lakoff. 2024. Women, fire, and dangerous	
	things: What categories reveal about the mind. Uni-	
	versity of Chicago press.	
	Joseph LeDoux. 2012. Rethinking the emotional brain.	
	<i>Neuron</i> , 73(4):653–676.	
	Messi HJ Lee, Soyeon Jeon, Jacob M Montgomery, and	
	Calvin K Lai. 2025. Visual cues of gender and race	
	are associated with stereotyping in vision-language	
	models. <i>arXiv preprint arXiv:2503.05093.</i>	
	Michelle Lelwica. 2011. The religion of thinness.	
	<i>Scripta Instituti Donneriani Aboensis</i> , 23:257–285.	
	Patricia W Linville, Gregory W Fischer, and Peter Sa-	
	lovev. 1989. Perceived distributions of the character-	
	istics of in-group and out-group members: empirical	
	evidence and a computer simulation. <i>Journal of per-</i>	
	<i>sonality and social psychology</i> , 57(2):165.	
	Jiayi Liu, Bo Jiang, and Yu’ang Wei. 2025. Llms as	
	promising personalized teaching assistants: How do	
	they ease teaching work? <i>ECNU Review of Educa-</i>	
	<i>tion</i> , 8(2):343–348.	
	Scott Lundberg. 2017. A unified approach to	
	interpreting model predictions. <i>arXiv preprint</i>	
	<i>arXiv:1705.07874.</i>	
	Sibo Ma, Alejandro Salinas, Julian Nyarko, and Peter	
	Henderson. 2025. Breaking down bias: On the limits	
	of generalizable pruning strategies. In <i>Proceedings</i>	
	<i>of the 2025 ACM Conference on Fairness, Account-</i>	
	<i>ability, and Transparency</i> , pages 2437–2450.	

1158	Nishtha Madaan, Sameep Mehta, Tanea Agrawaal, Vrinda Malhotra, Aditi Aggarwal, Yatin Gupta, and Mayank Saxena. 2018. Analyze, detect and remove gender stereotyping from bollywood movies. In <i>Conference on fairness, accountability and transparency</i> , pages 92–105. PMLR.	1215
1159		1216
1160		1217
1161		1218
1162		1219
1163		1220
1164	Nishtha Madaan, Sameep Mehta, Tanea S Agrawaal, Vrinda Malhotra, Aditi Aggarwal, and Mayank Saxena. 2017a. Analyzing gender stereotyping in bollywood movies. <i>arXiv preprint arXiv:1710.04117</i> .	1222
1165		1223
1166		1224
1167		1225
1168	Nishtha Madaan, Sameep Mehta, Mayank Saxena, Aditi Aggarwal, Tanea S Agrawaal, and Vrinda Malhotra. 2017b. Bollywood movie corpus for text, images and videos. <i>arXiv preprint arXiv:1710.04142</i> .	1226
1169		1227
1170		1228
1171		
1172	Manuj Malik and Richard Johansson. 2022. Controlling for stereotypes in multimodal language model evaluation . In <i>Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP</i> , pages 263–271, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	1229
1173		1230
1174		1231
1175		1232
1176		
1177		
1178		
1179	Robert K Merton. 1948. The self-fulfilling prophecy. <i>The antioch review</i> , 8(2):193–210.	1233
1180		1234
1181	Aatmaj Mhatre. 2023. Detecting the presence of social bias in gpt-3.5 using association tests. In <i>2023 international conference on advanced computing technologies and applications (ICACTA)</i> , pages 1–6. IEEE.	1235
1182		1236
1183		1237
1184		1238
1185		1239
1186	Margaret Mitchell, John Smith, Alice Lee, Ravi Kumar, and Li Wang. 2025. Shades: Towards a multilingual assessment of stereotypes in language models . In <i>Proceedings of the 2025 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 11995–12041. Association for Computational Linguistics.	1240
1187		1241
1188		1242
1189		1243
1190		1244
1191		1245
1192		1246
1193		
1194	Dean Mobbs, Cindy C Hagan, Tim Dalgleish, Brian Silston, and Charlotte Prévost. 2015. The ecology of human fear: survival optimization and the nervous system. <i>Frontiers in neuroscience</i> , 9:121062.	1247
1195		1248
1196		1249
1197		1250
1198	Brian Mullen, John F Dovidio, Craig Johnson, and Carolyn Copper. 1992. In-group-out-group differences in social projection. <i>Journal of Experimental Social Psychology</i> , 28(5):422–440.	1251
1199		1252
1200		1253
1201		1254
1202	Abdulrahman O Musaiger, Abdul-hai A Al-Awadi, and Mariam A Al-Mannai. 2000. Lifestyle and social factors associated with obesity among the bahraini adult population. <i>Ecology of food and nutrition</i> , 39(2):121–133.	1255
1203		1256
1204		1257
1205		1258
1206		
1207	Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained language models . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 5356–5371, Online. Association for Computational Linguistics.	1259
1208		1260
1209		1261
1210		
1211		
1212		
1213		
1214		
	Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020. CrowS-pairs: A challenge dataset for measuring social biases in masked language models . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1953–1967, Online. Association for Computational Linguistics.	1262
		1263
		1264
		1265
		1266
		1267
	Aurélie Névéol, Yoann Dupont, Julien Bezançon, and Karèn Fort. 2022. French crows-pairs: Extending a challenge dataset for measuring social bias in masked language models to a language other than english. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8521–8531.	1268
		1269
		1270
		1271
		1272
	Gandalf Nicolas, Xuechunzi Bai, and Susan T Fiske. 2021. Comprehensive stereotype content dictionaries using a semi-automated method. <i>European Journal of Social Psychology</i> , 51(1):178–196.	
	Shangrui Nie, Michael Fromm, Charles Welch, Rebekka Görgé, Akbar Karimi, Joan Plepi, Nazia Mowmita, Nicolas Flores-Herr, Mehdi Ali, and Lucie Flek. 2024. Do multilingual large language models mitigate stereotype bias? In <i>Proceedings of the 2nd Workshop on Cross-Cultural Considerations in NLP</i> , pages 65–83.	
	Ali Omrani, Alireza Salkhordeh Ziabari, Charles Yu, Preni Golazizian, Brendan Kennedy, Mohammad Atari, Heng Ji, and Morteza Dehghani. 2023. Social-group-agnostic bias mitigation via the stereotype content model. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4123–4139.	
	Tiago P Pagano, Rafael B Loureiro, Fernanda VN Lisboa, Rodrigo M Peixoto, Guilherme AS Guimarães, Gustavo OR Cruz, Maira M Araujo, Lucas L Santos, Marco AS Cruz, Ewerton LS Oliveira, and 1 others. 2023. Bias and unfairness in machine learning models: a systematic review on datasets, tools, fairness metrics, and identification and mitigation methods. <i>Big data and cognitive computing</i> , 7(1):15.	
	Devah Pager and Hana Shepherd. 2008. The sociology of discrimination: Racial discrimination in employment, housing, credit, and consumer markets. <i>Annual Review of Sociology</i> , 34:181–209.	
	Bo Pang. 2025. <i>Investigating Stereotypical Bias in Large Language and Vision-Language Models</i> . Ph.D. thesis, University of Auckland New Zealand.	
	Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. 2022. Bbq: A hand-built bias benchmark for question answering. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2086–2105.	
	Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel R Bowman. 2021. Bbq: A hand-built bias benchmark for question answering. <i>arXiv preprint arXiv:2110.08193</i> .	

1273	B Keith Payne, Heidi A Vuletich, and Kristjen B Lundberg. 2017. The bias of crowds: How implicit bias bridges personal and systemic prejudice. <i>Psychological Inquiry</i> , 28(4):233–248.	1326
1274		1327
1275		1328
1276		1329
1277	Alejandro Peña, Julian Fierrez, Aythami Morales, Gonzalo Mancera, Miguel Lopez-Duran, and Ruben Tolosana. 2025. Addressing bias in llms: Strategies and application to fair ai-based recruitment. In <i>Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society</i> , volume 8, pages 1976–1987.	1330
1278		1331
1279		1332
1280		1333
1281		1334
1282		1335
1283	Thomas F. Pettigrew. 1998. Intergroup contact theory . <i>Annual Review of Psychology</i> , 49:65–85.	1336
1284		1337
1285	Thomas F. Pettigrew and Linda R. Tropp. 2006. How does intergroup contact reduce prejudice? meta-analytic tests of three mediators .	1338
1286		1339
1287		1340
1288	Edmund S. Phelps. 1972. The statistical theory of racism and sexism. <i>The American Economic Review</i> , 62(4):659–661.	1341
1289		1342
1290		1343
1291	Tom Ed Postmes and Nyla R Branscombe. 2010. <i>Rediscovering social identity</i> . Psychology Press.	1344
1292		1345
1293	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. " why should i trust you?" explaining the predictions of any classifier. In <i>Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining</i> , pages 1135–1144.	1346
1294		1347
1295		1348
1296		1349
1297		1350
1298		1351
1299	Eleanor Rosch. 2024. Principles of categorization. In <i>Cognition and categorization</i> , pages 27–48. Routledge.	1352
1300		1353
1301		1354
1302	David A Rosenbaum. 2014. <i>It's a jungle in there: How competition and cooperation in the brain shape the mind</i> . Oxford University Press.	1355
1303		1356
1304		1357
1305	Nihar Sahoo, Pranamy Kulkarni, Arif Ahmad, Tanu Goyal, Narjis Asad, Aparna Garimella, and Pushpak Bhattacharyya. 2024. IndiBias: A benchmark dataset to measure social biases in language models for Indian context . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 8786–8806, Mexico City, Mexico. Association for Computational Linguistics.	1358
1306		1359
1307		1360
1308		1361
1309		1362
1310		1363
1311		1364
1312		1365
1313		1366
1314		1367
1315	Wolfgang S Schmeisser-Nieto, Alessandra Teresa Cignarella, Tom Bourgeade, Simona Frenda, Alejandro Ariza-Casabona, Laurent Mario, Paolo Giovanni Cicirelli, Andrea Marra, Giuseppe Corbelli, Farah Benamara, and 1 others. 2024. Stereohoax: a multilingual corpus of racial hoaxes and social media reactions annotated for stereotypes. <i>Language Resources and Evaluation</i> , pages 1–39.	1368
1316		1369
1317		1370
1318		1371
1319		1372
1320		1373
1321		1374
1322		1375
1323	Jörg Schweinitz, Johanna Eder, Fotis Jannidis, and Ralf Schneider. 2010. Stereotypes and the narratological analysis of film characters. <i>Revisionen</i> , (3):276–289.	1376
1324		1377
1325		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
		1406
		1407
		1408
		1409
		1410
		1411
		1412
		1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
		1421
		1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
		1431
		1432
		1433
		1434
		1435
		1436
		1437
		1438
		1439
		1440
		1441
		1442
		1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
		1458
		1459
		1460
		1461
		1462
		1463
		1464
		1465
		1466
		1467
		1468
		1469
		1470
		1471
		1472
		1473
		1474
		1475
		1476
		1477
		1478
		1479
		1480
		1481
		1482
		1483
		1484
		1485
		1486
		1487
		1488
		1489
		1490
		1491
		1492
		1493
		1494
		1495
		1496
		1497
		1498
		1499
		1500

1382	Aditya Tomar, Nihar Ranjan Sahoo, and Pushpak Bhat-	Wu Zekun, Sahan Bulathwela, and Adriano Soares	1436
1383	tacharyya. 2025b. Bharatbbq: A multilingual bias	Koshiyama. 2023. Towards auditing large language	1437
1384	benchmark for question answering in the indian con-	models: Improving text-based stereotype detection.	1438
1385	text. <i>Transactions of the Association for Computa-</i>	<i>ArXiv</i> , abs/2311.14126.	1439
1386	<i>tional Linguistics</i> , 13:1672–1692.		
1387	John C Turner, Michael A Hogg, Penelope J Oakes,	Kankan Zhou, Eason Lai, and Jing Jiang. 2022. VI-	1440
1388	Stephen D Reicher, and Margaret S Wetherell. 1987.	stereoset: A study of stereotypical bias in pre-trained	1441
1389	<i>Rediscovering the social group: A self-categorization</i>	vision-language models. In <i>Proceedings of the 2nd</i>	1442
1390	<i>theory</i> . basil Blackwell.	<i>Conference of the Asia-Pacific Chapter of the Asso-</i>	1443
1391	John C Turner and Katherine J Reynolds. 2003. The so-	<i>ciation for Computational Linguistics and the 12th</i>	1444
1392	cial identity perspective in intergroup relations: Theo-	<i>International Joint Conference on Natural Language</i>	1445
1393	ries, themes, and controversies. <i>Blackwell handbook</i>	<i>Processing (Volume 1: Long Papers)</i> , pages 527–538.	1446
1394	<i>of social psychology: Intergroup processes</i> , pages		
1395	133–152.		
1396	Eddie Ungless, Amy Rafferty, Hrichika Nag, and Björn		
1397	Ross. 2022. A robust bias mitigation procedure based		
1398	on the stereotype content model. In <i>Proceedings of</i>		
1399	<i>the Fifth Workshop on Natural Language Process-</i>		
1400	<i>ing and Computational Social Science (NLP+ CSS)</i> ,		
1401	pages 207–217.		
1402	Yilin Wang and Chujun Lin. 2024. Stereotypes at the		
1403	intersection of perceivers, situations, and identities:		
1404	analyzing stereotypes from storytelling using natural		
1405	language processing.		
1406	Ze Wang, Zekun Wu, Xin Guan, Michael Thaler, Adri-		
1407	ano Koshiyama, Skylar Lu, Sachin Beepath, Ediz		
1408	Ertekin, and Maria Perez-Ortiz. 2024. Jobfair: A		
1409	framework for benchmarking gender hiring bias in		
1410	large language models. In <i>Findings of the associ-</i>		
1411	<i>ation for computational linguistics: EMNLP 2024</i> ,		
1412	pages 3227–3246.		
1413	L Monique Ward and Kimberly Friedman. 2006. Using		
1414	tv as a guide: Associations between television view-		
1415	ing and adolescents’ sexual attitudes and behavior.		
1416	<i>Journal of research on adolescence</i> , 16(1):133–156.		
1417	Margaret Wetherell and Jonathan Potter. 1992. <i>Mapping</i>		
1418	<i>the language of racism: Discourse and the legitima-</i>		
1419	<i>tion of exploitation</i> . Columbia University Press.		
1420	Jeff Wise. 2009. <i>Extreme fear: The science of your</i>		
1421	<i>mind in danger</i> . St. Martin’s Press.		
1422	Huimin Xu, Zhang Zhang, Lingfei Wu, and Cheng-		
1423	Jun Wang. 2019. The cinderella complex: Word		
1424	embeddings reveal gender stereotypes in movies and		
1425	books. <i>PloS one</i> , 14(11):e0225385.		
1426	Yi Yang, Hanyu Duan, Ahmed Abbasi, John P Lalor,		
1427	and Kar Yan Tam. 2025. Bias a-head? analyzing bias		
1428	in transformer-based language model attention heads.		
1429	In <i>Proceedings of the 5th Workshop on Trustworthy</i>		
1430	<i>NLP (TrustNLP 2025)</i> , pages 276–290.		
1431	Abdelrahman Zayed, Gonçalo Mordido, Samira Sha-		
1432	banian, Ioana Baldini, and Sarath Chandar. 2024.		
1433	Fairness-aware structured pruning in transformers.		
1434	In <i>Proceedings of the AAAI Conference on Artificial</i>		
1435	<i>Intelligence</i> , volume 38, pages 22484–22492.		

1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493

A Stereotypes, Bias, Prejudice and Discrimination

In Section 6.1, we discuss bias, prejudice, and discrimination as core components of social harm. Below, we briefly elaborate on each of these concepts to clarify their roles and interrelationships in the formation and perpetuation of social harm.

A.1 Stereotypes

Stereotypes are overgeneralized beliefs about members of a social category, attributing uniform traits to all individuals and ignoring individual differences (Devine, 1989). For example, “Asians are good at math” overlooks variation within the group and can lead to unfair assumptions (Snyder and Swann, 1978; Steele and Aronson, 1995). Anti-stereotypes—beliefs that counter prevailing stereotypes—can reduce biased thinking (Devine, 1989). Only beliefs about social categories, not general truths, qualify as stereotypes.

A.2 Bias

Bias refers to inclinations or partiality favoring or disadvantaging certain groups, which can be explicit (conscious) or implicit (automatic) (Dovidio et al., 2010; Bodenhausen and Richeson, 2010). Explicit bias underlies overt discrimination, whereas implicit bias operates unconsciously, subtly influencing perceptions and behaviors (Fiske et al., 2002). Bias is distinct from stereotypes, though stereotypical biases arise from underlying stereotypical beliefs (Gallegos et al., 2024).

A.3 Prejudice

Prejudice is an affective attitude toward individuals based solely on their social category, reflecting emotions such as fear, contempt, or dislike (Allport, 1954; Devine, 1989). It often arises automatically (System 1) but can be mitigated through deliberate reflection (System 2) (Kahneman, 2011). Prejudice forms the emotional basis for discriminatory behavior and can be reduced by positive intergroup contact (Pettigrew and Tropp, 2006; Pettigrew, 1998).

A.4 Discrimination

Discrimination is the behavioral enactment of biased attitudes, leading to unfair treatment of individuals or groups (Allport, 1954). It can be:

- **Direct:** overt actions such as refusing service or workplace harassment (Becker, 1957; Dovidio et al., 2010).

- **Indirect:** neutral-appearing policies or practices that disproportionately disadvantage certain groups, e.g., standardized tests or institutional barriers (Phelps, 1972; Arrow, 1973; Pager and Shepherd, 2008).

A.5 Distinguishing Stereotypes from Bias

Recent work (Shejole and Bhattacharyya, 2025) highlights persistent conceptual confusion between stereotypes and biases, which has led to the construction of benchmarks having inconsistencies for stereotypes (e.g., MGSD (Zekun et al., 2023), EMGSD King et al. (2024)). This confusion limits the validity and generalizability of stereotype-detection models, as they often capture surface-level biases rather than the underlying social structures that define stereotypes. Bias is a distinct concept and should not be confused with stereotypes. While stereotypical bias refers to biases that originate from underlying stereotypes, stereotypes themselves are not equivalent to bias. For a detailed discussion of bias in the context of LLMs, we refer the reader to Gallegos et al. (2024).

B Summarizing Future Directions

In Section 7, we argued that grounding future computational research in established social psychological foundations, together with the research directions outlined in this paper, can enable the development of more principled, culturally grounded, and effective Responsible AI interventions. We summarize these future research directions in Table 1. The table provides a systematic synthesis of the research scope and open opportunities identified throughout this review, explicitly mapping them to the paper’s sections and subsections. It highlights key directions for bridging social psychological theories, such as the Five-Tuple Framework, with computational research areas including multimodal narrative analysis and broader global linguistic coverage. In addition, the table identifies challenges related to scalability and the evolving nature of stereotypes, while situating these issues within Responsible AI efforts focused on implicit bias detection and model interpretability. By organizing these gaps and opportunities, Table 1 offers a structured roadmap for future interdisciplinary work aimed at understanding and mitigating stereotypes in LLMs.

Table 1: Future Research Scope and Opportunities: Bridging Social Psychological and Computational Perspectives.

Section	Subsection	Future Research Scope & Opportunities
Section 2	Major Frameworks (§ 2.2)	Leverage the Five-Tuple Framework (Target, Relation, Attributes, Community, Time) to enable structured computational analysis, such as through knowledge graph-based representations.
	Computational Operationalization (§ 3.1)	Focus on using social-psychological theories to guide the development of robust techniques for measuring and operationalizing stereotypes; address gaps in multilingual and multicultural contexts.
Section 3	Narrative/Media (§ 3.2)	Implement proactive identification of stereotypes in media narratives to assess and mitigate potential social harms before dissemination.
	Body-Image (§ 3.3)	Systematically quantify body-image bias in LLMs and develop automatic modeling from media representations to monitor stereotypical ideals.
Section 4	Multimodality (§ 4.1)	Expand investigations into stereotype detection and mitigation beyond text and images to include conversational audio and video.
	Linguistic/Geographic Coverage (§ 4.2)	Create conceptually grounded, multilingual benchmarks moving beyond English/US-centric data; include complex dimensions like caste and regional state-level perceptions (e.g., India or USA).
Section 5	Generalization (§ 5.1)	Research more efficient methods for social analysis to help models handle unseen target groups and extract context-specific information.
	Annotation (§ 5.2)	Select representative annotator subsets reflecting the target community to ensure unbiased benchmarks and avoid skewed selections.
	Scalability (§ 5.3)	Explore strategies for modeling contexts separately to achieve global inclusivity despite current resource and scalability constraints.
	Dynamic Nature (§ 5.4)	Systematically study the dynamic nature of stereotype shifts through efficient modeling approaches, drawing insights from social-psychological theories and frameworks.
Section 6	Stereotype as the origin (§ 6.1)	Monitor and prevent self-fulfilling prophecies and stereotype threat; investigate whether LLMs and AI models exhibit personal biases similar to humans and understand underlying causes.
	Implicit Bias (§ 6.2)	Conduct more research revealing implicit bias through measures like simulated implicit association tests and other psychological frameworks.
	Mitigation, Interpretability and Explainability (§ 6.3)	Removing Implicit Bias for mitigation; Anti-stereotypes for mitigation; Identify stereotype subspaces in LLMs; use explainability techniques (e.g., SHAP, LIME) to analyze model attributions through established theories; investigate impacts on original task efficiency.

Aspect	Stereotype Content Model (SCM)	Agency-Beliefs-Communion (ABC) Model
Core dimensions	Warmth and competence	Agency and beliefs; communion is emergent
Methodological stance	Theory-driven; predefined groups and traits	Data-driven; dimensions emerge from spontaneous judgments
Conceptual focus	Intentions (warmth) and ability (competence)	Socioeconomic power (agency) and ideology (beliefs)
Role of communion	Fundamental evaluative dimension	Derived from combinations of agency and beliefs
Group perception	Warmth and competence vary independently	Extreme agency predicts lower perceived communion

Table 2: Comparison of the Stereotype Content Model (SCM) and the Agency-Beliefs-Communion (ABC) Model.

Theory / Framework	Core Assumptions	View of Stereotypes	Key References
Similarity-Attraction & Social Identity Theory	Individuals derive self-esteem from group memberships; intergroup comparison motivates ingroup favoritism and outgroup derogation. Social identity is shaped by perceived group belonging.	Stereotypes function as self-esteem regulators that maintain positive social identity and reinforce ingroup–outgroup boundaries.	Byrne (1971); Tajfel and Turner (1979); Turner and Reynolds (2003); Ellemers and Haslam (2012)
Social Role Theory	Social structures and role distributions shape expectations about groups; repeated exposure normalizes role-based differences.	Stereotypes emerge as reflections of socially assigned roles and are reinforced through cultural and media representations.	Eagly (1987); Ward and Friedman (2006); Gauntlett (2008); Bartlett et al. (2013)
Social Categorization Theory	Humans perceive the social world through group-based categorization; context determines which identities become salient.	Stereotypes are fluid, context-dependent representations emerging from group-level perception rather than fixed beliefs.	Turner et al. (1987); Augoustinos and Walker (1998)
Social Cognition Theories	Cognitive efficiency drives humans to rely on schemas and heuristics to manage informational complexity.	Stereotypes are cognitive shortcuts—functional yet potentially biasing mental representations.	Fiske (1992); Fiske and Haslam (1996); Fiske and Taylor (2020)
System Justification Theory	Individuals are motivated to preserve existing social hierarchies, even when personally disadvantaged by them.	Stereotypes serve ideological functions by legitimizing and stabilizing unequal social systems.	Jost et al. (2004); Jost and Van der Toorn (2012); Jost (2019); Banaji (2002)
Discursive Approaches to Categorization	Social reality is constructed through language and discourse rather than fixed cognitive representations.	Stereotypes are discursive resources—contextual, flexible, and rhetorically constructed in interaction.	Wetherell and Potter (1992); Edwards (1991); Augoustinos and Walker (1998)
Intersectionality Theory	Social identities are interdependent and mutually constitutive rather than additive.	Stereotypes emerge at the intersections of multiple identities, producing context-specific and compounded forms of marginalization.	Crenshaw (2013); Cho et al. (2013); Carastathis (2014)
Stereotype Content Model (SCM)	Group perception is structured along warmth and competence dimensions shaped by competition and status.	Stereotypes map onto predictable emotional and behavioral responses (e.g., admiration, pity, contempt).	Fiske et al. (2002); Cuddy et al. (2011)
Agency-Beliefs-Communion (ABC) Model	Social perception is organized around agency and ideological beliefs, with communion emerging secondarily.	Stereotypes reflect perceived power relations and ideological alignment rather than intrinsic warmth.	Koch et al. (2016)
Dual-Perspective (Facet) Model	Agency and communion each consist of multiple sub-dimensions (e.g., assertiveness, morality).	Stereotypes operate through fine-grained evaluative dimensions rather than coarse traits.	Abele et al. (2016)
Five-Tuple Framework	Stereotypes are relational, contextual, and temporally grounded phenomena.	Stereotypes are structured as (Target, Relation, Attributes, Community, Time Interval), enabling computational modeling.	Davani et al. (2025); Shejole and Bhattacharyya (2025)

Table 3: Summary of major theories and frameworks explaining the formation, function, and structure of stereotypes across social psychology and computational social science.

1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589

C Summarizing Social-Psychological Theories and Frameworks

In Section 2, we discussed various theories and frameworks related to stereotypes. We summarize them in Table 3. We compare the SCM Model and the ABC Model in Table 2.

D Briefly Analyzing Failure of Bias Mitigation Strategies

In Section 6.3, we note that current bias mitigation techniques exhibit notable limitations, which are briefly discussed in this section.

There are various techniques for bias mitigation (Gallegos et al., 2024). From a computational perspective, it becomes clear that many existing techniques fail because they focus on surface-level symptoms, including words, tokens, or decoding heuristics, rather than the underlying causes of harm. These root drivers include biased data collection practices, entangled social identities, model inductive biases, and poorly specified objectives. Consequently, interventions based on limited word lists, proxy attributes, or simple reweighting often miss substantial forms of harm or introduce new distortions, such as erasure, reduced representational diversity, and unintended distribution shifts. Many approaches rely on strong but implicit assumptions, including binary or immutable social categories, the interchangeability of harms across groups, or the preservation of meaning under surface-level substitutions. Such assumptions rarely hold in realistic linguistic and social contexts. In addition, mitigation methods frequently optimize inappropriate metrics, for example token-level parity, rather than outcomes tied to downstream social impact. Together with computational constraints and the brittleness of classifiers used to identify harmful content, these limitations result in mitigation strategies that appear effective on narrow benchmarks but fail when evaluated with real users and within existing power structures. Meaningful progress therefore requires approaches that target root causes through careful attention to data provenance and representational choices, articulate explicit fairness objectives linked to concrete harms, and employ rigorous, human-centered evaluation guided by social-psychological principles. We refer the reader to Gallegos et al. (2024) for a detailed discussion of bias mitigation techniques.

From a social-psychological perspective, many mitigation strategies primarily target the explicit

components of social harm, such as overtly toxic or abusive outputs. However, as discussed in Section 6.2, addressing social harm also requires confronting implicit bias in models. It also guides that Anti-stereotypes can be used for stereotype mitigation in reducing human prejudice (Cuddy et al., 2008; Fraser et al., 2021), hence this techniques can also be explored for mitigation in the future.

E Use of AI Assistants

We used Gemini and ChatGPT to assist with minor writing refinements and grammatical corrections.

1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600