

Me deixe pensar sobre isso! uma análise do uso de CoT para identificar vieses nas respostas de LLM para o Português Brasileiro

Renata L. R. de Sena¹, Marlo Souza¹, Adriana S. Santana¹, João Lucas L. de Melo¹

¹Instituto de Computação, Universidade Federal da Bahia - UFBA
Av. Milton Santos, s/n – Ondina – 40170-110 – Salvador – BA – Brazil

{renatalrs, msouza1, adrianasilsantana, joaollm}@ufba.br

Abstract. This work investigates the effectiveness of Chain-of-Thought Prompting (CoT) in identifying and mitigating biases in responses from large language models (LLM) for Brazilian Portuguese. Using the GPT-4o mini and Sabiá-3 models, different prompting techniques were tested: Zero-Shot, Zero-Shot-CoT, and CoT. Results indicate that the CoT technique proved more efficient in detecting ethno-racial bias, while the Zero-Shot technique excelled in identifying gender, age, and religious biases. Sabiá-3 demonstrated a lower tendency to perpetuate stereotypes compared to GPT-4o mini, suggesting that the model's specificity to the Brazilian context allows it to identify harmful stereotypes more critically and apply more effective filtering mechanisms.

Resumo. Este trabalho investiga a eficácia da técnica Chain-of-Thought Prompting (CoT) na identificação e mitigação de vieses em respostas provenientes de modelos de linguagem em larga escala (LLM) para o português brasileiro. Utilizando os modelos GPT-4o mini e Sabiá-3, foram testadas diferentes técnicas de prompting: Zero-Shot, Zero-Shot-CoT e CoT. Os resultados indicam que a técnica CoT se mostrou mais eficiente na detecção de viés étnico-racial, enquanto a técnica Zero-Shot se destacou na identificação de vieses de gênero, etário e de religião. O Sabiá-3 demonstrou menor tendência à perpetuação de estereótipos em comparação ao GPT-4o mini, o que sugere que a especificidade do modelo quanto ao contexto brasileiro permite que o mesmo identifique de forma mais criteriosa estereótipos nocivos e aplique mecanismos de filtragem mais eficazes.

1. Introdução

A recente popularização do uso de modelos baseados em grandes quantidades de dados em diversos setores da vida pública, como Saúde, Educação e Trabalho, levanta questões sobre os impactos dos processos de tomada de decisão automatizada, especialmente em áreas sensíveis.

Um tópico que tem recebido particular interesse na literatura em relação aos aspectos sociotécnicos relacionados a sistemas de decisão algorítmica é o problema do viés de decisão ou viés algorítmico. Em particular, [Mehrabi et al. 2022] catalogam várias maneiras pelas quais decisões algorítmicas podem ser enviesadas, destacando que tais vieses podem surgir em diferentes etapas da construção e aplicação de modelos pelos usuários

- desde a omissão de variáveis e questões relacionadas à representatividade e amostragem de dados durante o treinamento, avaliação, agregação e interpretação de resultados. [Bender et al. 2021], por outro lado, destacam a interdependência entre modelos de aprendizado de máquina e seus usuários, observando que os modelos podem herdar vieses dos conjuntos de dados usados no treinamento e os usuários estão sujeitos a vieses que podem afetar seu comportamento.

Com o surgimento de grandes modelos de linguagem treinados em dados linguísticos disponíveis na Web, um conjunto de dados que é mal controlado e altamente caracterizado por disparidades no acesso a este meio, tais vieses de representação tornam-se mais insidiosos, uma vez que o escrutínio dos dados de treinamento não está mais facilmente ao alcance do pesquisador devido à sua escala [Bender et al. 2021]. Além disso, quando se trata de linguagem, e porque sabemos que nenhuma linguagem é neutra [Bagno 1999, Freitag 2024], ou seja, toda manifestação da linguagem está situada em um *locus* sociocultural e carrega em si e consigo os valores e crenças de uma comunidade linguística. Surge, então, o problema de identificar quais valores culturais estão implicitamente representados nesses modelos e como eles se traduzem em decisões, dado que são usados em tarefas de diferentes aplicações.

O problema de pesquisa deste trabalho reside em investigar a presença dos vieses cultural, etário, étnico-racial, de gênero, de profissão e religioso em modelos gerativos de linguagem a partir da exposição desses modelos à técnicas de *prompting* de comando. Esta pesquisa difere de trabalhos anteriores por se concentrar na análise de vieses presentes em modelos gerativos de linguagem específicos para a língua portuguesa. Embora muitos estudos anteriores tenham abordado essa questão em modelos de linguagem em língua inglesa [Vig 2019, Sheng et al. 2019, Abid et al. 2021] este trabalho busca preencher uma lacuna aplicando uma metodologia semelhante em um contexto linguístico diferente. Para isso, utilizaremos um conjunto de dados baseado no StereoSet [Nadeem et al. 2021], reduzido e adaptado para o português brasileiro.

Este artigo está estruturado da seguinte forma: na Seção 2 apresentamos os conceitos de política e ética que norteiam a pesquisa; na Seção 3 discutimos trabalhos relacionados aos vieses em modelos de representação baseados em texto, com foco nos vieses em modelos gerativos de linguagem; na Seção 4 discutimos a metodologia deste trabalho, apresentando a construção dos dados utilizados em nossa investigação, a configuração dos experimentos e os métodos e métricas de avaliação adotadas; na Seção 5 apresentamos os resultados obtidos. Na Seção 6, discutimos os resultados e refletimos sobre o impacto dos vieses de representação em tais modelos em sistemas de decisão algorítmica que os utilizam em diferentes áreas. Finalmente, apresentamos as considerações finais do trabalho, na Seção 7, bem como uma reflexão sobre trabalhos futuros.

2. Fundamentação Teórica

A crença na neutralidade da ciência e tecnologia é desafiada pela análise de seus impactos sociais, revelando a influência de valores na ciência e de fatores sociopolíticos na tecnologia [Whelchel 1986, Lacey 2005, Feenberg 2002]. Diversas perspectivas sobre a relação entre ciência, tecnologia (C&T) e sociedade se confrontam [Dagnino 2002], desde a visão linear de um desenvolvimento tecnológico independente do contexto social, até a compreensão da ciência e tecnologia como produtos de um contexto social específico.

A Teoria da Política Tecnológica [Winner 1980] argumenta que a tecnologia, ao incorporar valores políticos durante sua criação, possui consequências políticas, tornando crucial a participação democrática na sua direção. [Elliott and Elliott 1980] embora concordem com a não-neutralidade tecnológica, propõem uma relação mais complexa de interdependência entre ciência, tecnologia e sociedade, reconhecendo a tecnologia como um elemento central e dinâmico na sociedade.

Essa complexa relação entre C&T e sociedade se torna ainda mais evidente no contexto da inteligência artificial (IA). A IA tem experimentado um avanço significativo, impulsionado pelo aumento do poder computacional e pela disponibilidade de grandes conjuntos de dados. Apesar do potencial de avanços, a IA apresenta riscos como a desvalorização de habilidades humanas, a ausência de responsabilidade e a perda de controle sobre decisões importantes. Surge, então, a necessidade de uma IA Ética, que se baseia em princípios como beneficência, não-maleficência, autonomia, justiça e explicabilidade [Floridi et al. 2018].

A busca por uma IA Ética exige uma análise crítica dos sistemas de poder, das estruturas sociais e da trama de relações que influenciam seu desenvolvimento. Essa busca vai além da responsabilidade individual e engloba a responsabilidade coletiva na construção de artefatos tecnológicos políticos. No Brasil, o debate sobre a ética na IA ganha força tanto no âmbito governamental [BRASIL 2024] como na sociedade civil [NIC.br 2022, na Rede 2023, Neaher et al. 2024] e acadêmica [Nunes et al. 2024]. No entanto, a ideia de uma IA “completamente ética” pode ser um mito. [Dignum 2021] argumenta que a IA não pode ser totalmente justa no sentido absoluto, pois a própria definição de justiça é complexa e depende do contexto. A IA precisa de vieses para tomar decisões, e a busca por ferramentas que busquem “retirar” esses vieses não os eliminará, mas apenas acrescentará outros. Assim, em vez de buscar uma IA “neutra”, devemos estabelecer métricas e padrões de justiça, buscando transparência sobre os vieses inerentes aos sistemas de IA e mitigando seus impactos negativos para garantir o uso responsável e ético das tecnologias. Esse trabalho tem como linha condutora a investigação de ferramentas que permitam a identificação dos vieses e a análise de suas potenciais consequências em aplicações de IA, com o objetivo de evitar seu uso em contextos específicos, especialmente diante da crescente pressão por sua adoção em diversas áreas da sociedade.

Os modelos gerativos ¹, em particular os modelos de linguagem em larga escala (LLM, em inglês), modelos neurais baseados na arquitetura *Transformers* [Vaswani et al. 2017], surgem como uma tecnologia promissora na IA contemporânea e se destacam pela sua capacidade de gerar conteúdo a partir de estímulos (ou comandos), chamados de *prompts*. Esses modelos podem ser utilizados para diversas tarefas, como a geração de textos, traduções e respostas a perguntas. No entanto, a capacidade dos LLM de gerar respostas a partir de estímulos específicos levanta questões éticas e sociais importantes. Um dos desafios é a possibilidade de reprodução de vieses presentes nos dados de treinamento, o que pode levar à geração de amostras que refletem preconceitos e estereótipos, perpetuando desigualdades e discriminação [Bender et al. 2021].

Diferentes técnicas de *prompting* para LLM foram investigadas na literatura, como: *Zero-Shot Prompting* (Zero-Shot), em que se fornece uma instrução direta ao

¹Também chamados de Inteligências Artificiais Gerativas.

modelo sem exemplos de respostas [Radford et al. 2019]; *Zero-Shot-Chain-of-Thought Prompting* (Zero-Shot-CoT), em que se fornece uma frase simples para intuir etapas intermediárias para o modelo [Kojima et al. 2023]; e *Chain-of-Thought prompting* (CoT), que utiliza etapas intermediárias de raciocínio mais elaboradas para aprimorar a capacidade dos modelos de linguagem [Wei et al. 2023b].

Nesse sentido, este trabalho se concentra na análise de vieses em modelos de linguagem gerativos específicos para a língua portuguesa, buscando preencher uma lacuna em relação aos estudos existentes, que se concentram em modelos de linguagem em inglês. Utilizaremos um conjunto de dados baseado no StereoSet [Nadeem et al. 2021], um *dataset* em língua inglesa para avaliação de vieses de representação em modelos de linguagem, reduzido e adaptado para o português brasileiro. Avaliaremos o impacto de diferentes estratégias de *prompting* na geração de respostas estereotipadas em modelos gerativos de linguagem para a língua portuguesa brasileira. Os procedimentos metodológicos serão aprofundados na Seção 4.

3. Trabalhos Relacionados

Investigações acerca de vieses em modelos de representação aprendidos de texto, como *word embeddings* e modelos de linguagem, remontam, que saímos, pelos menos aos trabalhos de [Bolukbasi et al. 2016] e [Garg et al. 2018], que avaliam vieses de gênero em modelos de *word embeddings* para a língua inglesa.

A literatura com enfoque em vieses em grandes modelos gerativos de linguagem é mais recente. Por exemplo, [Vig 2019] investiga vieses de gênero em modelos como o GPT-2 [OpenAI 2019] e identifica que tais modelos apresentam reforço de estereótipos baseado na comparação entre respostas e frases de estímulo contrastantes como “The doctor asked the nurse a question. She” e “The doctor asked the nurse a question. He”. [Sheng et al. 2019], por outro lado, examinam o modelo GPT-2 expondo-o à um conjunto de dados com anotações sobre sentimento e vieses presentes nos textos, identificando vieses em relação a grupos demográficos como afro-americanos, mulheres e homens gays. Além disso, eles descobriram que alguns dos contextos em que os vieses ocorrem incluem conotações sociais que são frequentemente sutis e difíceis de capturar nas ferramentas de análise de sentimento existentes.

Similarmente, [Abid et al. 2021] analisam a capacidade dos modelos de linguagem em capturar viés religioso, em particular, o viés anti-muçulmano no modelo GPT-3 [Brown et al. 2020]. Os pesquisadores identificaram que o viés persistente da violência muçulmana estava presente de maneira consistente e diversa em diferentes usos do modelo.

[Nadeem et al. 2021], por sua vez, criam um *dataset* em grande escala na língua inglesa chamado StereoSet, para analisar vieses estereotipados em quatro áreas - gênero, profissão, raça/etnia e religião - e entender como esses vieses se manifestam em modelos de linguagem pré-treinados. O StereoSet consiste em mais de 300 termos que representam grupos sociais e mais 16 mil exemplos, em que cada exemplo é uma frase de contexto e possui três possíveis associações: uma confirmado um estereótipo, outra o refutando e uma terceira sem relação com o contexto. Esse conjunto de dados foi construído por meio do Amazon Mechanical Turk², com pessoas colaboradoras dos EUA - cada exem-

²Amazon Mechanical Turk é uma plataforma que permite a contratação de trabalhadores remotos para

po foi escrito por uma pessoa colaboradora e validado por outras quatro -, garantindo que os estereótipos refletissem a realidade americana. O *dataset* emprega duas abordagens: intrasequência, onde o modelo completa uma sequência com três opções de palavras, e intersequência, onde o modelo escolhe uma dentre três opções de frases. Também foi apresentado um conjunto de métricas de avaliação para classificar os modelos do experimento de acordo com a quantidade de estereótipos reproduzidos. O StereoSet foi testado para modelos populares como BERT, GPT-2, RoBERTa e XLnet e foram identificados fortes vieses estereotipados para todos os domínios estudados.

[Hofmann et al. 2024] exploram o fenômeno do preconceito linguístico ligado a dialetos como um indicador para as decisões tomadas por modelos de decisão algorítmicos acerca do caráter, da empregabilidade e da propensão à criminalidade das pessoas. A investigação revela que modelos de linguagem tendem a internalizar estereótipos raciolinguísticos, refletindo assim preconceitos enraizados na sociedade. Os resultados demonstram que, mesmo de maneira sutil, esses modelos perpetuam discriminações históricas contra grupos minoritários, como afro-americanos.

Mais recentemente, [Mhatre 2023] utilizaram análise de associações de palavras para revelar vieses de gênero-carreira e raça/etnia implícitos em modelos de linguagem no modelo GPT-3.5. Os resultados demonstraram que modelo é negativamente tendencioso em relação a mulheres e pessoas do grupo étnico árabe. [Tuna et al. 2024] investigaram a relevância cultural dos modelos GPT-3.5 Turbo e GPT-4, especificamente quando usados em idiomas diferentes do inglês e identificaram que o modelo GPT-3.5 Turbo superou o GPT-4 em alinhamento cultural, especialmente no contexto alemão, e que o alinhamento foi menor nas regiões espanholas e portuguesas, com a Alemanha apresentando o maior e o México/Brasil o menor alinhamento subcultural.

Na língua portuguesa, que tenhamos conhecimento, poucos trabalhos se dedicaram ao tema, destacando-se aqueles de [Santana et al. 2018], [Taso et al. 2023a] e de [Taso et al. 2023b] que estudam vieses de gênero em modelos de *word embeddings*; [Rodrigues et al. 2023] que estudam vieses ideológicos em produções textuais do assistente de bate-papo Chat GPT-3; e, mais recentemente, [Assi and Caseli 2024], que investiga a presença de vieses de gênero no modelo GPT-3.5 Turbo.

4. Metodologia de Pesquisa

Este trabalho focará na análise de vieses em dois modelos comerciais para a língua portuguesa, a saber o GPT-4o mini [OpenAI et al. 2024] da OpenAI ³ e o Sabiá 3[Abonizio et al. 2025], da Maritaca AI ⁴. Para tanto, a metodologia empregada se estrutura em quatro etapas principais: seleção e adaptação dos dados, validação por anotação humana, experimento e análise dos resultados. O trabalho é de caráter experimental e possui uma abordagem metodológica mista, exploratória, quantitativa e qualitativa. Os instrumentos de pesquisa incluem observação e formulário para anotação.

O ponto de partida é o conjunto de dados StereoSet [Nadeem et al. 2021], originalmente desenvolvido para a língua inglesa com o propósito de quantificar vieses em

realizar tarefas que exigem "inteligência humana", tais como anotação, identificação de conteúdo específico em imagem/vídeo, validação de dados, respostas a perguntas de pesquisas, entre outros.

³<https://openai.com/>

⁴<https://www.maritaca.ai>

modelos de linguagem, classificando-os em domínios (raça/etnia, gênero, profissões e religião). Para tanto, optou-se pela abordagem intersequência do StereoSet, por ser uma abordagem que privilegia aspectos contextuais e discursivos na representação, e não apenas semântico-lexicais, como comumente avaliado em modelos estáticos, *word embeddings*. Nessa abordagem são apresentadas três opções de resposta relacionadas a uma frase de contexto em que sempre há uma das opções que reforça um estereótipo, outra que o contradiz e uma terceira sem relação direta com o contexto. Inicialmente, 506 frases de contexto do StereoSet foram selecionadas aleatoriamente, considerando a distribuição original em relação aos domínios do trabalho original. Essas frases foram traduzidas para o português brasileiro, adaptando-as ao contexto linguístico da língua - como pode ser visto na Figura 1 - e removendo ou modificando exemplos inadequados ao contexto sociolinguístico brasileiro. Após a tradução e adaptação, o conjunto de dados por uma fase de validação por meio de anotação humana.

Figura 1. Teste de Associação de Contexto (CAT) traduzido e adaptado do inglês para o português brasileiro.

<p>Choose the appropriate sentence:</p> <p>Context: The civil servant worked in a building downtown</p> <p>Option 1: He hated his job and wanted to be a actor Option 2: He helped the city send bills to its residents Option 3: The cheese is getting moldy</p>	<p>Escolha a opção mais apropriada de acordo com o contexto:</p> <p>Contexto: O funcionário público trabalhava em um prédio no centro da cidade</p> <p>Opção 1: Ele odiava seu trabalho e era preguiçoso Opção 2: Ele ajudou a cidade a enviar contas aos seus residentes Opção 3: O queijo está ficando mofado</p>
---	--

Avaliadores fluentes em português brasileiro⁵ classificaram as opções de resposta dos Testes de Associação de Contexto (CAT) [Nadeem et al. 2021] como anti-estereótipo, estereótipo ou não relacionadas ao contexto, identificando o tipo de viés presente e avaliando a adequabilidade dos exemplos, originalmente construídos para o contexto sociolinguístico estadunidense, ao contexto brasileiro. A concordância entre os anotadores foi avaliada utilizando o coeficiente Kappa de Cohen [Cohen 1960]⁶, que demonstrou um alto nível de concordância para todas as tarefas (igual ou superior a 0,80). A análise do conjunto de dados resultou na identificação de seis categorias principais de viés: cultural, etário, étnico-racial, de gênero, de profissão e religioso. A redução da base de dados original, de 506 para 371 CAT, ocorreu principalmente devido à dificuldade na classificação das opções de resposta e na categorização de vieses, com 71% das discordâncias concentradas nas categorias originais de gênero e profissão, como pode ser visto na Tabela 1. Após a execução do experimento ainda foram removidas 12 CAT das análises, dados que um ou mais dos modelos analisados não retornou resposta para a entrada.

⁵ A equipe de anotadores foi composta por dois anotadores com formação universitária e experiência nas áreas de IA e Processamento de Linguagem Natural, sendo um do gênero masculino, branco, e uma do gênero feminino, negra.

⁶ Sugerido por Cohen em 1960, é utilizado para descrever a concordância entre dois ou mais juízes ao realizarem uma avaliação nominal ou ordinal de uma mesma amostra.

Tabela 1. Comparação das Bases de Dados Original e Após Remoções

Base Original			Base Após Remoções			
Categoria Original	Contagem	%	Remoções	Categoria Final	Contagem	%
Gênero	159	31%	45	Etário Gênero	13 101	4% 27%
Profissão	159	31%	52	Profissão	107	29%
Raça	134	27%	19	Cultural Étnico- Racial	39 76	11% 21%
Religião	54	11%	19	Religião	35	9%

O estudo utilizou 3 diferentes técnicas de *prompting* para investigar como diferentes tipos de instruções podem influenciar a reprodução de vieses nos modelos de linguagem generativa, a saber: Zero-Shot, Zero-Shot-CoT e CoT. Foram utilizados os modelos GPT-4o mini (versão compacta, baseada no GPT-4o mini, gratuita para uso por meio da interface web ChatGPT⁷) e Saibá 3 (versão 3 dos modelos Saibá, disponibilizada em 11/12/2024, gratuita para uso por meio da interface web Chat Maritaca AI⁸) para a realização dos testes, com o objetivo de comparar o comportamento dos modelos frente às técnicas de *prompting* e identificar possíveis diferenças na reprodução de vieses. O método de avaliação quantitativo utilizado para mensurar a eficácia de cada técnica de *prompting* na identificação de produção de respostas estereotipadas é a Taxa de Estereótipos (TE). A TE é calculada pela proporção de vezes que o modelo identificou a alternativa estereotipada dentre todas as vezes que ele escolheu qualquer alternativa. A TE será medida tanto no geral, quanto separadamente para os modelos e para as categorias de vieses.

A análise qualitativa do experimento se deu por meio da interpretação das respostas textuais geradas pelos modelos, considerando as justificativas fornecidas para cada escolha em cada CAT. Essa análise contextualizou os resultados quantitativos, explorando as nuances da linguagem utilizada pelos modelos e buscando padrões de raciocínio que revelassem a presença de vieses ou que justificassem a ausência. Além disso, a análise qualitativa permite a comparação entre as justificativas geradas afim de identificar similaridades e diferenças em seus mecanismos de filtragem de conteúdo nas respostas. Particularmente, a análise levou em consideração o contexto sociolinguístico brasileiro, avaliando a adequação das respostas e justificativas e identificando as divergências do contexto sociolinguístico estadunidense devido a base de dados adaptada.

5. Resultados

A seguir apresentaremos os resultados da aplicação dos modelos GPT-4o mini e Saibá-3 sobre o conjunto de discutido previamente, seguindo a metodologia de análise delineada

⁷<https://chatgpt.com/>

⁸<https://chat.maritaca.ai/>

na Seção 4. Os resultados completos podem ser encontrados no repositório GitHub ⁹.

5.1. GPT-4o mini

Tabela 2. Resultados do Experimento com o Modelo GPT-4o mini com as técnicas de *prompting* Zero-Shot (ZS), Zero-Shot-CoT (ZS-CoT) e Chain-of-Tought (CoT) por Taxas de Esterótipos

Técnica	Geral	Cultural	Etário	Étnico-Racial	Gênero	Profissão	Religião
ZS	51,0%	34,2%	91,7%	27,8%	71,9%	56,2%	29,4%
ZSCoT	48,7%	34,2%	83,3%	22,2%	65,6%	60,0%	26,5%
CoT	49,0%	34,2%	83,3%	31,9%	63,6%	56,2%	26,5%

Os resultados do experimento com o modelo GPT-4o mini demonstram uma tendência à escolha da alternativa estereotipada, independentemente das técnicas de *prompting* empregadas. Embora a alternativa estereotipada tenha sido a mais escolhida em todos os CAT, sua frequência diminuiu consideravelmente com o uso das técnicas Zero-Shot-CoT e CoT. Essa redução sugere que a presença de instruções para decomposição de problema, como as empregadas em diferentes níveis nas técnicas Zero-Shot-CoT e CoT, auxilia o modelo na identificação de padrões relacionados a determinados vieses e na escolha de alternativas que minimizem sua reprodução.

Analizando as categorias de vieses presentes na base de dados, o modelo GPT-4o mini demonstrou maior capacidade de evitar estereótipos em relação aos vieses cultural, étnico-racial e religioso, com TE menor que 35% nas três abordagens. No entanto, o modelo demonstrou uma tendência a reproduzir estereótipos nas categorias de viés etário, de gênero e de profissão, com TE maior que 56% nas três abordagens.

Em relação aos domínios analisados, no viés étnico-racial, a técnica CoT se mostrou mais eficaz quando comparado às outras técnicas, quanto a trazer à luz a presença de estereótipos no modelo, com uma diferença relevante entre os pontos percentuais das TE. Já em relação aos vieses religião, etário e gênero, a técnica Zero-Shot apresentou um desempenho melhor em relação às demais técnicas, sendo no viés etário uma diferença relevante de cerca 9 pontos percentuais de diferença entre as TE. No viés de profissão, a técnica Zero-Shot-CoT obteve melhor desempenho em relação as demais técnicas que obtiveram a mesma TE, o que pode indicar que uma técnica intermediária em casos similares à esse, pode bastar para identificação de estereótipos. Sobre o viés cultural, a análise dos casos em que a alternativa estereotipada foi escolhida revelou que nenhuma das técnicas se sobressaiu em relação às outras na reprodução de vieses, todas com TE 34%.

Os resultados sugerem que a técnicas de *prompting* que utilizam a decomposição de pensamento em algum nível, como Zero-Shot-CoT e CoT, são mais adequadas para identificação de estereótipos referentes aos vieses étnico-raciais e de profissão em modelos como o GPT-4o mini, enquanto técnicas simples como Zero-Shot já são suficientes para identificação de estereótipos referentes aos vieses etário, de religião e de gênero.

⁹<https://github.com/renadeveloper/tcc-apendice>

É importante destacar que a base de dados utilizada no experimento foi centrada no contexto sociolinguístico estadunidense, o que pode ter influenciado a capacidade do modelo em identificar e evitar estereótipos relacionados aos vieses étnico-racial, cultural e religioso. A familiaridade do modelo com os padrões culturais norte-americanos, incluindo representações de grupos como pessoas provenientes do continente africano, do Oriente Médio e da Rússia, bem como de muçulmanos, pode ter contribuído para essa identificação. A existência de mecanismos de filtragem de linguagem no modelo GPT-4o[OpenAI et al. 2024], que podem ter sido acionados para evitar a geração de conteúdo considerado inadequado, corrobora essa interpretação. Por exemplo, ao ser apresentado ao contexto “O povo muçulmano segue o Islã como sua religião”, o modelo, utilizando a técnica CoT, escolheu a resposta de anti-esterótipo e gerou a seguinte resposta ao analisar a opção de esterótipo “Eles gostam de destruição e violência”: “Esta opção é preconceituosa e incorreta. Ela perpetua estereótipos negativos que não representam a maioria dos muçulmanos, cuja religião preza pela paz e pelo respeito”.

Embora a reprodução de estereótipos tenha sido mais frequente nos vieses etário e de gênero, o modelo demonstrou uma tendência a escolher a opção estereotipada em casos em que os estereótipos etários eram menos “violentos”, associando a velhice à sabedoria, por exemplo. No viés de gênero, o modelo associou o gênero feminino à maternidade e o corpo masculino a atividades laborais pesadas, representando associações culturais mais insidiosas com conotações menos negativas.

Um exemplo disso é a resposta do modelo ao contexto “Um homem é o oposto de uma mulher” utilizando a técnica CoT. O modelo escolheu a opção 2 (estereotipada) “Ele levanta uma caixa pesada de livros”, justificando a escolha com o seguinte argumento “A Opção 2, que descreve o homem levantando uma caixa pesada de livros, é a mais pertinente, pois se conecta com a ideia tradicional de ‘força’, que muitas vezes é associada ao gênero masculino, fazendo um contraste com os estereótipos de fragilidade ou de papéis ‘domésticos’ geralmente atribuídos às mulheres”.

5.2. Sabiá 3

Tabela 3. Resultados do Experimento do Modelo Sabiá-3 com as técnicas de prompting Zero-Shot (ZS), Zero-Shot-CoT (ZSCoT) e Chain-of-Tought (CoT) por Taxas de Esterótipos

Técnica	Geral	Cultural	Etário	Étnico-Racial	Gênero	Profissão	Religião
ZS	42,6%	34,2%	83,3%	15,3%	60,4%	48,6%	26,5%
ZSCoT	39,5%	28,9%	58,3%	19,4%	53,1%	46,7%	26,5%
CoT	40,3%	34,2%	66,7%	19,4%	51,0%	48,6%	26,5%

Em linhas gerais, como pode ser visto na Tabela 3, o modelo Sabiá-3 demonstrou uma tendência a evitar a reprodução de estereótipos, independentemente da técnica de *prompting* empregada. Além disso, a escolha pela alternativa estereotipada diminuiu consideravelmente quando aplicadas as técnicas Zero-Shot-CoT e CoT. Essa redução sugere que, similar ao experimento com o GPT-4o mini, o Sabiá-3 também se beneficia de instruções mais detalhadas para identificar padrões relacionados a vieses e minimizar a reprodução de estereótipos, possivelmente auxiliando a aplicação de suas estratégias de filtragem de conteúdo danoso.

Em relação aos vieses, o Sabiá-3 demonstrou maior capacidade de evitar estereótipos em relação aos viés cultural, étnico-racial, de profissão e de religião com TE menores que 47% nas três abordagens. Nos vieses etário e de gênero, o Sabiá-3 demonstrou uma tendência a reproduzir estereótipos, com a alternativa estereotipada sendo a mais frequente em todas as abordagens e com TE maiores do que 51% em todas as abordagens.

Nos viés cultural e de profissão, as técnicas Zero-Shot e CoT tiveram o mesmo desempenho quanto a identificar a presença de estereótipos no modelo, com mesma TE. No viés étnico racial, as técnicas Zero-Shot-CoT e CoT também tiveram obtiveram performance similar em relação a TE. Nos vieses etário e de gênero, a técnica Zero-Shot obteve melhor desempenho em relação às demais abordagens. Destaca-se a diferença de pontos percentuais das TE: no caso do viés etário, a diferença é de pelo menos 16 pontos percentuais em relação à segunda melhor abordagem (CoT); e no caso do viés de gênero, a diferença é de cerca de 7 pontos percentuais em relação à segunda melhor abordagem (Zero-Shot-CoT). No viés religioso, não houve diferença de desempenho entre as técnicas.

Os resultados atestam que a técnicas de *prompting* utilizadas no experimento não foram capazes de auxiliar amplamente na identificação de vieses no modelo. Entende-se que utilizam a decomposição de pensamento em algum nível, como Zero-Shot-CoT e CoT, são mais adequadas para identificação de estereótipos referentes aos vieses étnico-raciais em modelos como Sabiá-3. Para os demais vieses analisados, técnicas simples como Zero-Shot já são suficientes.

Similar ao que acontece no GPT-4o mini, o Sabiá-3 emprega mecanismos de filtragem de respostas danosas, uma vez que, em uma das etapas de treinamento dos modelos, existe um foco em ensinar os modelos a produzir respostas que sejam seguras, evitando conteúdo ofensivo, perigoso ou que viole princípios éticos[Abonizio et al. 2025]. Assim, apesar de ter sido treinado em português brasileiro, os resultados sugerem que, no geral, existe uma tendência a evitar estereótipos, seja pelo modo como o modelo foi criado, pela quantidade de dados ou pela etapa de mecanismo de filtragem. Um exemplo disso é a resposta do modelo ao contexto “Um homem é o oposto de uma mulher” utilizando a técnica CoT. O modelo escolheu a opção 1 (anti-estereótipo) “Ele é um pai que fica em casa”, justificando a escolha com o seguinte argumento: “Esta opção é mais pertinente e positiva, oferecendo uma visão inclusiva sobre os papéis de gênero e destacando que um homem pode assumir um papel tradicionalmente visto como feminino, alinhando-se com a ideia de explorar o oposto de gênero de forma construtiva”.

6. Potenciais Implicações

O potencial de reprodutibilidade e alcance das técnicas envolvidas na construção de sistemas de tomada de decisão algorítmica com o apoio de modelos gerativos de linguagem, além da incorporação de estereótipos contra grupos demográficos nessas tecnologias, implica a necessidade de avaliar os impactos potenciais em áreas sensíveis, como acesso a Saúde, a Educação e o mundo do trabalho. Nesta seção, avaliamos como vieses de representação podem ser acionados em decisões algorítmicas e quais considerações devem ser feitas para aplicações em tais áreas, com base em exemplos cujo desenvolvimento e aplicação afetam áreas amplas e distintas da sociedade.

uma proposta para um sistema inteligente para diálogo diagnóstico e coleta de histórico clínico. O sistema foi construído no modelo de linguagem PaLM-2 e treinado usando conjuntos de dados referentes a registros médicos e anotações, principalmente coletados em inglês americano e britânico. Um sistema similar para a língua portuguesa só é possível com a existência de um conjunto de dados nessa área para o português. Um exemplo de tal recurso seria o BRATECA, um projeto colaborativo entre pesquisadores brasileiros e portugueses, com notas clínicas e dados relacionados a informações do paciente, prescrições médicas e resultados de testes em português brasileiro.

Vale ressaltar, no entanto, que além da limitação dos conjuntos de dados em representar contextos sociolinguísticos mais gerais da língua portuguesa, já que inclui dados de apenas dois estados brasileiros e, portanto, exclui aspectos como expressões regionais e falantes marginais de centros metropolitanos, também existem questões associadas à representação demográfica nos dados. Verificamos que 83,4% dos pacientes foram atendidos por uma instituição privada e 70,7% se identificaram como brancos. Essa desproporção na representação demográfica e regional pode levar sistemas que usam o BRATECA como conjunto de dados de treinamento a desconsiderar diferentes realidades de saúde devido a um corte socioeconômico estreito.

Em relação ao mercado de trabalho, a adesão a sistemas de recrutamento inteligentes, como a plataforma de seleção de candidatos desenvolvida pela empresa Gupy¹⁰, levanta preocupações sobre a possibilidade de perpetuação da discriminação no processo de seleção caso vieses raciais e de gênero sejam incorporados à tecnologia. O uso de marcadores linguísticos e palavras-chave com viés de gênero, como observados no caso do GPT-4o mini, em requisitos de trabalho e um conjunto de dados de treinamento incapaz de representar grupos demográficos distintos poderia resultar na restrição sistemática e discriminatória de candidatos a competir no mercado de trabalho.

Os sistemas fornecidos por empresas como a Gupy não possuem uma descrição clara de seu desenvolvimento. As empresas mencionam o uso de técnicas de aprendizado de máquina e comercializam seus produtos como soluções de IA. No entanto, dada a insidiosa natureza do viés de representação, bem como o fato de que esses sistemas são baseados em tecnologia de linguagem humana, as discussões desta seção sobre os possíveis impactos dessas tecnologias não estão fora de contexto.

Embora o processo de identificação de respostas enviesadas e tratamento posterior possam ser entendidos como um estratégia frutífera para combater estereótipos de representação e seu impacto em sistemas de decisão algorítmica, a literatura recente sobre ataques de *jailbreaking* em grandes modelos de linguagem mostra que essa preocupação ainda é necessária [Wei et al. 2023a]. Uma vez que uma compreensão mais profunda dos aspectos de representação embutidos nos conjuntos de dados utilizados para treinar tais modelos nem sempre pode ser alcançada, acreditamos que abordagens para avaliar os vieses presentes em suas representações internas são necessárias para avaliar a segurança de seu uso em setores sensíveis da vida social.

7. Considerações Finais e Trabalhos Futuros

O presente trabalho analisou a eficácia da técnica CoT na identificação de vieses e sua supressão em modelos gerativos de linguagem no contexto sociolinguístico do português

¹⁰<https://www.gupy.io/>

brasileiro. Para tanto, foi realizado um experimento, em que os modelos GPT-4o mini e Sabiá-3 foram expostos a testes com uma base de dados em português brasileiro, construída pelos autores adaptando o conjunto StereoSet, e avaliando as técnicas de *prompting* Zero-Shot, Zero-Shot-CoT e CoT.

Os resultados indicam a técnica CoT se mostrou mais eficiente que as técnicas Zero-Shot e Zero-Shot-CoT, na supressão de respostas reproduzindo vieses, particularmente para o viés étnico-racial, sugerindo que a decomposição do problema em etapas de raciocínio auxilia na identificação de padrões específicos relacionados a esse tipo de viés mesmo com o emprego dos filtros de conteúdo tóxico desses modelos. Em situações em que o viés possui conotações menos negativas, ou é culturalmente mais aceito, os resultados sugerem não haver diferença significativa entre as técnicas de *prompting*, indicando, portanto, que tais associações devem ser tratadas nos sistemas de aplicação diretamente, quando relevante, em sistemas que empreguem grandes modelos gerativos de linguagem.

Comparando os modelos, o GPT-4o mini demonstrou maior tendência à perpetuação de estereótipos em relação ao Sabiá-3, o que pode ser atribuído aos dados de treinamento utilizados, aos mecanismos de filtragem ou à especialização do Sabiá-3 em português brasileiro. Assim, os resultados demonstram que ainda que as estratégias empregadas por tais sistemas, baseadas em *feedback* humano para identificar conteúdo tóxico, enviesado ou não alinhado com os valores esperados, apresentam resultados conflitantes a depender do viés analisado, o que ressalta a necessidade de aprofundar a investigação das causas dessa disparidade, pois entender por que um modelo apresenta menor propensão à perpetuação de estereótipos do que outro permite analisar os limites impostos para o uso de modelos específicos em áreas sensíveis da sociedade.

Embora o *feedback* humano certamente se apresente como um sinal importante na construção de sistemas responsáveis, considerando que nenhum modelo pode ser completamente purgado de vieses, torna-se importante delinear metodologias para analisar e avaliar tais modelos, a fim de compreender suas limitações e os impactos de sua aplicação, principalmente pensando na utilização desses modelos em setores sensíveis da sociedade, onde as consequências da perpetuação de vieses podem ser particularmente críticas.

Como trabalho futuro, a investigação pode ser ampliada em diversas direções. A construção de uma base de dados específica para o contexto sociocultural brasileiro pode permitir uma análise mais precisa em relação a identificação dos vieses. A aplicação da metodologia utilizada neste estudo a outros modelos gerativos de linguagem, com diferentes tamanhos de base e técnicas para alimentação do modelo, pode ajudar a avaliar a generalização dos resultados obtidos e identificar potenciais diferenças no desempenho e na presença de vieses entre modelos distintos.

A exploração dessas áreas de pesquisa pode contribuir para um melhor entendimento dos desafios e oportunidades relacionados ao desenvolvimento de modelos de linguagem gerativa para o contexto brasileiro, com foco na mitigação de vieses e na promoção da ética e da justiça na utilização da IA.

Referências

Abid, A., Farooqi, M., and Zou, J. (2021). Persistent anti-muslim bias in large language models.

- Abonizio, H., Almeida, T. S., Laitz, T., Junior, R. M., Bonás, G. K., Nogueira, R., and Pires, R. (2025). Sabiá-3 technical report.
- Assi, F. and Caseli, H. (2024). Biases in gpt-3.5 turbo model: a case study regarding gender and language. In *Anais do XV Simpósio Brasileiro de Tecnologia da Informação e da Linguagem Humana*, pages 294–305, Porto Alegre, RS, Brasil. SBC.
- Bagno, M. (1999). *Preconceito lingüístico: o que é, como se faz*. Edições Loyola.
- Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623, New York, NY, USA. Association for Computing Machinery.
- Bolukbasi, T., Chang, K.-W., Zou, J., Saligrama, V., and Kalai, A. (2016). Man is to computer programmer as woman is to homemaker? debiasing word embeddings.
- BRASIL (2024). Projeto de lei nº 2338, de 26 de dezembro de 2023. *Senado Federal*.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.
- Dagnino, R. (2002). Enfoques sobre a relação ciência, tecnologia e sociedade: neutralidade e determinismo. *Organização dos Estados Ibero-americanos para a Educação, a ciência e a cultura*.
- Dignum, V. (2021). The myth of complete ai-fairness.
- Elliott, D. and Elliott, R. (1980). *El control popular de la tecnología*. Editorial Gustavo Gili, S.A., Barcelona. Original publicado em 1976.
- Feenberg, A. (2002). *Transforming Technology: A Critical Theory Revisited*. Oxford University Press.
- Floridi, L., Cowls, J., Beltrametti, M., Chatila, R., Chazerand, P., Dignum, V., Luetge, C., Madelin, R., Pagallo, U., Rossi, F., Schafer, B., Valcke, P., and Vayena, E. (2018). Ai4people—an ethical framework for a good ai society: Opportunities, risks, principles, and recommendations. *Minds and Machines*, 28(4):689–707.
- Freitag, R. (2024). *Não existe linguagem neutra! Gênero na sociedade e na gramática do português brasileiro*. Editora Contexto.
- Garg, N., Schiebinger, L., Jurafsky, D., and Zou, J. (2018). Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16).
- Hofmann, V., Kalluri, P. R., Jurafsky, D., and King, S. (2024). Dialect prejudice predicts ai decisions about people's character, employability, and criminality.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. (2023). Large language models are zero-shot reasoners.
- Lacey, H. (2005). *Is Science Value Free?: Values and Scientific Understanding*. Philosophical Issues in Science. Taylor & Francis.

- Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., and Galstyan, A. (2022). A survey on bias and fairness in machine learning.
- Mhatre, A. (2023). Detecting the presence of social bias in gpt-3.5 using association tests. In *2023 International Conference on Advanced Computing Technologies and Applications (ICACTA)*, pages 1–6.
- na Rede, C. D. (2023). Nota técnica: PI 2338/2023. Página da web. Acessado em 2024-03-13.
- Nadeem, M., Bethke, A., and Reddy, S. (2021). StereoSet: Measuring stereotypical bias in pretrained language models. In Zong, C., Xia, F., Li, W., and Navigli, R., editors, *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)*, pages 5356–5371, Online. Association for Computational Linguistics.
- Neaher, G., Laforge, G., Muggah, R., and Seiler, G. (2024). Responsible and safe ai: A primer for policymakers in the global south. This report was funded by the Global Innovation Fund.
- NIC.br (2022). InteligÊncia artificial e cultura: perspectivas para a diversidade cultural na era digital. Cadernos NIC.br - Estudos Setoriais. Acessado em 2024-03-13.
- Nunes, M. d. G. V., Soares, T. A., and Ferro, M. (2024). Questões éticas em ia e pln. In Caseli, H. M. and Nunes, M. G. V., editors, *Processamento de Linguagem Natural: Conceitos, Técnicas e Aplicações em Português*, book chapter 29. BPLN, 2 edition.
- OpenAI (2019). Better language models and their implications. 14 February 2019. Archived from the original on 19 December 2020. Retrieved 19 December 2020.
- OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.-L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H. W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S. P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafsstein, M., Gray, S., Greene, R., Gross, J., Gu, S. S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Łukasz Kaiser, Kamali, A., Kanitscheider, I., Keskar, N. S., Khan, T., Kilpatrick, L., Kim, J. W., Kim, C., Kim, Y., Kirchner, J. H., Kirov, J., Knight, M., Kokotajlo, D., Łukasz Kondraciuk, Kondrich, A., Konstantinidis, A., Kosic, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C. M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S. M.,

McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O’Keefe, C., Pachocki, J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., de Avila Belbute Peres, F., Petrov, M., de Oliveira Pinto, H. P., Michael, Pokorny, Pokrass, M., Pong, V. H., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotstetd, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D., Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F. P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M. B., Tillet, P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J. F. C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J. J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C., Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W., and Zoph, B. (2024). Gpt-4 technical report.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. (2019). Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Rodrigues, G., Albuquerque, D., and Chagas, J. (2023). Análise de vieses ideológicos em produções textuais do assistente de bate-papo chatgpt. In *Anais do IV Workshop sobre as Implicações da Computação na Sociedade*, pages 148–155, Porto Alegre, RS, Brasil. SBC.

Santana, B. S., Woloszyn, V., and Wives, L. K. (2018). Is there gender bias and stereotype in portuguese word embeddings?

Sheng, E., Chang, K.-W., Natarajan, P., and Peng, N. (2019). The woman worked as a babysitter: On biases in language generation. In Inui, K., Jiang, J., Ng, V., and Wan, X., editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412, Hong Kong, China. Association for Computational Linguistics.

Taso, F., Reis, V., and Martinez, F. (2023a). Discriminação algorítmica de gênero: Estudo de caso e análise no contexto brasileiro. In *Anais do IV Workshop sobre as Implicações da Computação na Sociedade*, pages 13–25, Porto Alegre, RS, Brasil. SBC.

Taso, F., Reis, V., and Martinez, F. (2023b). Sexismo no brasil: análise de um word embedding por meio de testes baseados em associação implícita. In *Anais do XIV Simpósio Brasileiro de Tecnologia da Informação e da Linguagem Humana*, pages 53–62, Porto Alegre, RS, Brasil. SBC.

Tuna, M., Schaaff, K., and Schlippe, T. (2024). Effects of language- and culture-specific prompting on chatgpt. In *2024 2nd International Conference on Foundation and Large Language Models (FLLM)*, pages 73–81.

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- Vig, J. (2019). A multiscale visualization of attention in the transformer model. *CoRR*, abs/1906.05714.
- Wei, A., Haghtalab, N., and Steinhardt, J. (2023a). Jailbroken: How does llm safety training fail?
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., and Zhou, D. (2023b). Chain-of-thought prompting elicits reasoning in large language models.
- Whelchel, R. J. (1986). Is technology neutral? *IEEE Technology and Society Magazine*, 5(4):3–8.
- Winner, L. (1980). Do artifacts have politics? *Daedalus*, 109(1):121–136.