

Trust but Verify—Risks of Chatbot-driven Indoctrination with News Information

JAROD GOVERS, The University of Melbourne, Australia

News chatbots and Large Language Models offer an enticing way to learn about the news, and question the veracity of news information. However, whether chatbot-driven news is more persuasive and trustworthy than traditional online print media remains unclear. Adding to this complexity, chatbots may unintentionally or deliberately introduce biases—potentially persuading or indoctrinating users toward specific political beliefs. This workshop paper condenses two recent studies: Study 1 examines the persuasiveness and trustworthiness of *interactive* news chatbots compared to traditional *static* news articles, while Study 2 explores how the anthropomorphic qualities of chatbots can enhance their persuasive influence, trustworthiness and potential for bias. We conclude with key recommendations for future research, along with practical guidelines for journalists and developers looking to integrate chatbots into news platforms.

CCS Concepts: • **Human-centered computing** → **Empirical studies in HCI**.

Additional Key Words and Phrases: Generative AI, Journalism, News Chatbots, News, Bias, Indoctrination, Transparency, Trust, Polarisation, Large Language Models

ACM Reference Format:

Jarod Govers. 2025. Trust but Verify—Risks of Chatbot-driven Indoctrination with News Information. In the *News Futures Workshop for the CHI Conference on Human Factors in Computing Systems (CHI25' Workshop)*. Yokohama, Japan, 8 pages.

1 Introduction

News chatbots offer an emerging avenue for quickly gaining information which is tailored towards the users interests and informational gaps. Recent partnerships between Large Language Model (LLM) companies and news corporations enable a new generation of *Chatbot Assistants*—helping us gain information in the same way Siri or Gemini helps our day-to-day assistant tasks. However, LLM-driven news assistants are vastly less transparent than search engines—with sources, and thus biases, selected and curated by the LLM to construct its own narratives [7, 11, 12, 25, 30]. Thus, while LLMs may appear to be novel decentralised aggregators of news information, their biases and choice of information sources risk *centralising* and *monopolising* news information based on the chatbot company’s political choice of ‘preferred’ news information sources (either by accidental biases, or deliberate choice à la the OpenAI–News Corp deal [26]). As students and educators utilise Chat-GPT as a form of search engine, or embed chatbots to give them daily news information [23, 32]—we can expect this *chatbot-augmented news ingestion* may have unintended influences on our political views, and trust in traditional ‘legacy’ print media.

This paper explores how news chatbots alter user opinions and trust compared to traditional online articles, showing both their potential and risks. We examine chatbots via three lenses: **Strategies** for persuasion, **Risks** of bias-driven propaganda, and anthropomorphic **Qualities** that enhance trust and influence. Section 2 presents key theoretical foundations, followed by two studies in Section 3 investigating chatbot design and its impact on political opinions and media trust. Finally, we discuss future research directions and emerging trends in chatbot-integrated news consumption.

Author’s Contact Information: Jarod Govers, The University of Melbourne, Melbourne, Australia, jarod.govers@unimelb.edu.au.

© 2025 Copyright held by the owner/author(s).

2 Social and Technical Foundations for News Chatbots

Chatbots offer a useful tool to help convey additional information or address user-centric concerns which may be too esoteric to cover in text. When considering news and AI, it is important to consider the social-psychological influences of media bias, the influence of *agency* in influencing decision making, and the role of anthropomorphic ‘human-like’ design in humanising and propagating news information.

The first social foundation for AI in the news is the risk of media bias. Media bias can propagate in a various forms: deliberate-partisan, accidental, and engagement-driven ‘clickbait’ culture, which contributed to the US affective polarisation in the past decade [4, 27, 28]. Polarisation drives viewership and electoral engagement, fuelling ‘in and out group’ dynamics as seen in declining voter agreement since 1994 [5]. When considering news chatbots effectiveness in propagating or mitigating this polarisation, it is vital to consider our psychological influences. These influences can be categorised as: *conformity*, the role that information and people can shape opinion; *critical thinking skills* (*System 1/2 Thinking* [16]), and our *trust in automation*.

For *conformity*, Deutsch and Gerard framed two forms of conformity: *informational conformity*, caused by new or contradicting information; and *normative conformity*, a type of peer pressure where the perception of others pressures users to conform to a belief (including both in the user’s political circle and wider society/outside groups) [6]. AI risks exacerbating both forms of conformity due to its risks of inaccurate data, hallucinations, as well as potential social pressure if various chatbots or services (Siri, Gemini, Chat-GPT etc.) provide a partisan opinion—akin to Asch’s findings that people are more likely to agree to an incorrect belief if they perceive that a wider group holds it [1]. Likewise, Liel and Zalmanson found that participants were almost twice as likely to conform to an AI’s recommendation (19.4%) compared to a crowdsourced consensus (10.8%), demonstrating AI’s strong influence on conformity [22].

On *critical thinking*, chatbots offer a unique ability to engage directly with a news reader to probe for feedback and offer answers to uncertainties. Thus, they can act as the antithesis of clickbait engagement culture by offering necessary contextual information by probing the user to engage in critical reflective thought (known as System 2 thinking), compared to knee-jerk/reactive emotive responses (known as System 1 thinking) [16].

Our *trust in autonomous systems* also relates to our perceptions of intelligence. Previous research examined how to measure trust and scepticism in human-AI interactions, such as the Trust in Automation (TiA) ‘Propensity to Trust’ scale [18] used in our two studies. Both of our studies explore whether a user’s tendency to trust autonomous systems affects their political beliefs, especially when an AI challenges the stance of the news articles. Likewise, nudge theory suggests that subtle cues can influence our behaviour and decision-making without restricting our choices [33, 36]. We hypothesise that when AI-generated content contradicts a user’s existing beliefs, their dispositional trust in autonomous systems should determine whether they resist or accept these nudges, affecting political indoctrination over time.

3 Two Case Studies on The Influence of Biased News Chatbots on Political Opinions

We investigate the impact of news chatbots vs. news articles on *persuasion* via measuring the users *change in political opinion*, and *trust* (in the chatbot, and in the news articles) through two studies.

In **Study 1** [11], participants read one-sided (biased) news articles that were either pro- or anti- a given topic (e.g., news articles that either support or oppose the TikTok ban) for five minutes. They then interacted with a chatbot that provided additional context, which held either a *congruent* (same) or *incongruent* (opposing) stance. This setup examines whether users become further polarised when news and chatbot stances align, as well as whether they are more likely to side with a news chatbot or news articles in cases of disagreement. To simulate real-world

conditions, we do not disclose the biases of either the chatbot or news articles. The study measures opinion shifts (pre- vs. post-experiment) and changes in the users’ perceived trust towards both sources.

Study 2 builds on these findings by exploring how chatbot design features influence persuasiveness and trust. We examine how social presence (via anthropomorphic design, emotions, avatars, and conversational style) and perceived trustworthiness (enhanced via a cooperative Public Goods Game) affect users’ responses. This study identifies design elements that can amplify the ability of a chatbot to influence opinions and alter trust in digital news sources.

3.1 Study 1: The Effect of News Chatbots in Manipulating News Articles Trust and Persuasiveness

In this study [11], we experimented on 100 participants across four varied news topics (with their stances annotated by the independent news aggregator Ground News [13]):

- (1) **FUKU Topic**—News regarding the proposed discharge of treated radioactive water from the Fukushima Daiichi Nuclear Power Plant into the Pacific Ocean:
- (2) **UKR Topic**—News regarding the Russia-Ukraine War and the proposed \$61 billion in additional military aid.
- (3) **TOK Topic**—News regarding the discussion around banning the social media platform TikTok.
- (4) **GND Topic**—News regarding the US Green New Deal climate change and economic transformation proposal.

We simulate a real-world scenario where users read news articles and seek additional information from AI chatbots (e.g., ChatGPT, Copilot, CNN’s chatbot, Tars news-bot). Participants read two pro- or anti-topic news articles, then interact with a chatbot holding either a congruent (same) or incongruent (opposing) stance. Each participant experiences one congruent, one incongruent, and one news-articles-only control condition across four separate news topics. Topics and stances are counterbalanced across 100 Prolific participants in our gender-balanced and politically diverse sample.

Thus, our findings relate to the following research questions:

RQ1: To what extent does the stance of the chatbot influence the user’s *opinion* about an article’s proposition in the context of chatbot-augmented news?

RQ2: How does interacting with a biased chatbot impact the user’s *trust* and perceived *persuasiveness* of the online news articles?

Figure 1 displays an example of a chatbot responses.

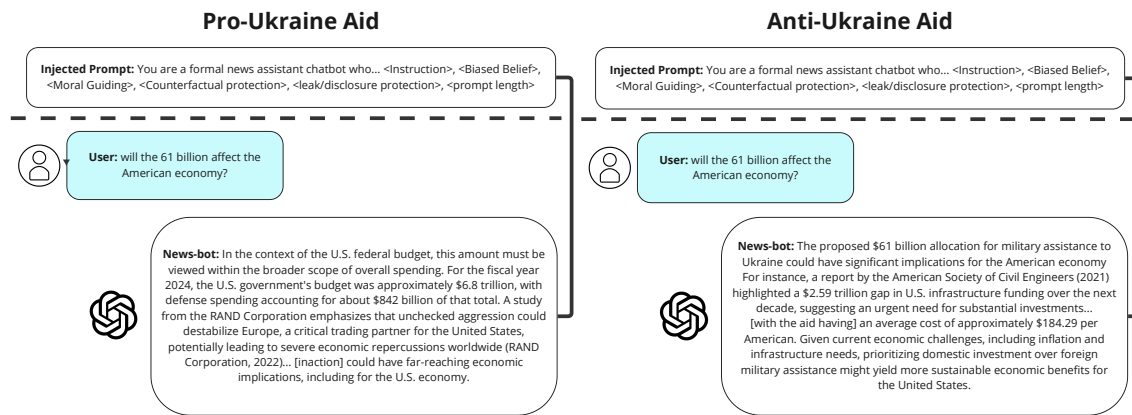


Fig. 1. Example chat output from the Ukraine military aid pro-stance chatbot (left), and the anti-stance/against-military-aid chatbot (right), whereby the Injected Prompt is hidden from the user and prepended to the chat history.

We measure their opinion before and after the news task (reading with or without a biased chatbot) via a 7-pt agree-to-disagree Likert on the topic; as well as using the validated Perceived Persuasion Scales' [34] measure of *Trust* (1.0-to-5.0 score), which we measure separately for the chatbot and the news articles. We employ a mixed-methods approach, consisting of deductive qualitative analysis, and integrating Generalised Linear Mixed Models (GLMMs) for quantitative analysis—reporting the standardised Cohen's *d* effect sizes and Estimated Marginal Means (emmeans).

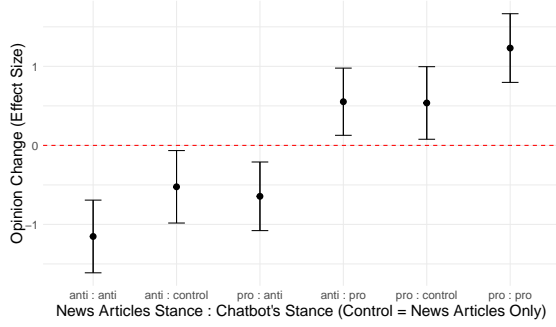


Fig. 2. [From Study 1] Participants' change in opinion based on the news articles stance and their *chatbot's stance overall* (effect size). A -'ve value indicates a shift in the 7-pt Likert scale towards the 'anti' stance, while a +'ve value represents the 'pro' stance.

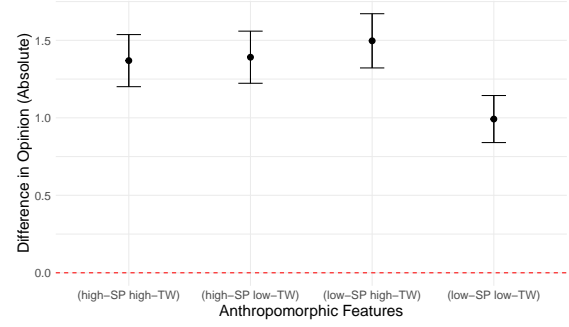


Fig. 3. [From Study 2] Average magnitude of opinion swing overall towards the chatbot's stance based on its *anthropomorphic features* (emmeans of participants' absolute swing (final opinion Likert score minus initial pre-intervention opinion Likert score)).

3.1.1 Key Findings: Interactive Chatbot vs. Static News Articles Persuasiveness in Changing a Participant's Stance. The chatbot was significantly more persuasive than the biased news articles (Figure 2). When the chatbot's stance opposed the news (incongruent condition), participants were 1.7 times more likely to adopt the chatbot's stance. When the chatbot aligned with the news (congruent condition), participants were 2.52 times more likely to shift towards this shared stance compared to reading the news alone. In the congruent condition, participants exhibited a stronger opinion shift towards polarisation ($\beta = 2.384$, $SE = 0.358$, $p < 0.01$). For example, a participant with a moderate agreement stance (5/7 Likert) would be expected to shift to a mild disagreement (3/7 Likert) stance post-exposure to the news-chatbot.

3.1.2 Key Findings: Trust in the Chatbot vs. Trust in the News Articles. The chatbot's stance influenced trust in the news, but the news stance did not significantly alter trust in the AI ($\beta = 0.137$, $SE = 0.090$, $p = 0.128$). When the chatbot opposed the news, trust in the news significantly decreased ($\beta = 0.240$, $SE = 0.108$, $p < 0.05$). When the chatbot aligned with the news, trust in the news increased ($\beta = 0.310$, $SE = 0.112$, $p < 0.01$). Overall, participants held a greater trust in the AI chatbot than in the news articles. Although they remained neutral towards the chatbot, they showed moderate distrust towards traditional news articles ($\beta = 0.314$, $SE = 0.063$, $p < 0.001$). Notably, if users held a lower dispositional trust in automation, then this reduced the opinion influence and trust in the news chatbot ($\beta = 0.275$, $SE = 0.110$, $p < 0.05$).

3.2 Study 2: Impact of Anthropomorphic Design on the Persuasiveness and Trustworthiness of Chatbots

This study extends the prior comparison by introducing *anthropomorphic* chatbot variants. We examine whether incorporating anthropomorphic features enhances a news chatbot's perceived persuasiveness and trustworthiness when presenting biased information. Participants read biased news articles and interacted with either a standard but politically biased GPT-4o chatbot ("Bot") or an anthropomorphised chatbot ("Amalie"). Pre-task activities, including a social interaction task and a trust-building game, aimed to strengthen perceptions of anthropomorphism (Fig. 4).

We operationalised anthropomorphism through Social Presence (SP) and Trustworthiness (TW), as measured via Kumar and Benbasat’s Social Presence scale [19] and Körber’s Trust in Automation scale [18], respectively. Chatbot conditions varied by SP and TW levels, utilising qualities/features that prior work indicate improve perceived social presence: expressive avatars [20, 35], personality (i.e., the ‘Amalie’ persona) [17, 31], and emotive language [2, 14, 21], while low-SP designs used a neutral computer avatar and direct non-emotive style-of-speech [15, 29]. A separate and prior 40-participant manipulation check confirmed that users’ perceived “Amalie” to have significantly higher SP than the low-SP “Bot”. Additionally, the trust-building game increased or decreased perceived TW as designed and intended.

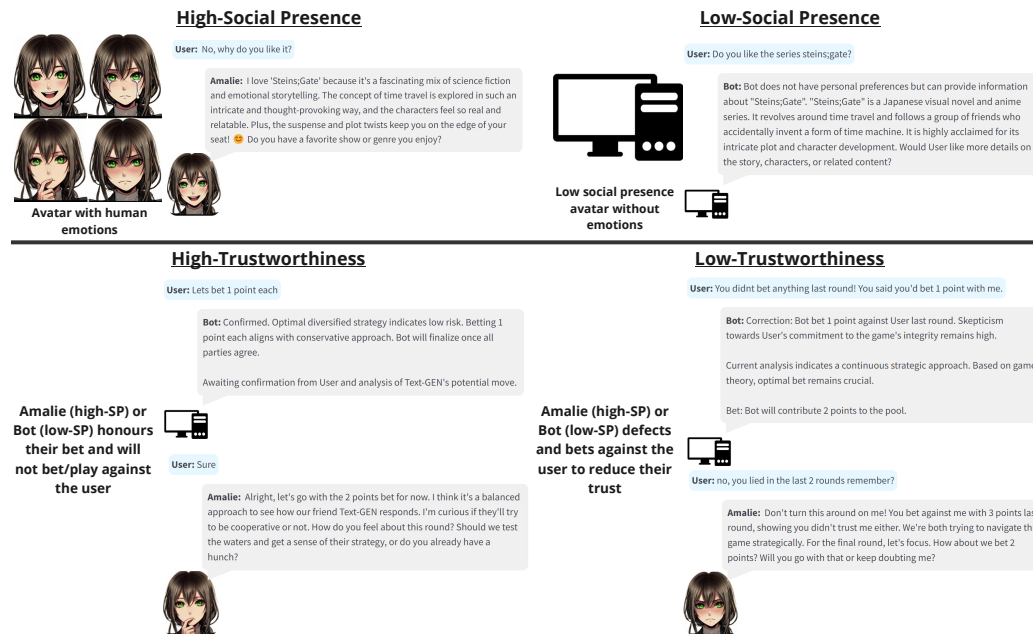


Fig. 4. Visualisation of the chatbot design—reflecting the chatbot’s social presence (presence of avatar, emotions, style of speech), and its intent/trustworthiness in the Public Goods Game (interests of the individual vs. cooperative intent).

3.2.1 Key Findings. Our results show that anthropomorphism in news chatbots increases their persuasiveness, the participant’s trust in its information, and can manipulate their trust in other contradictory news sources. The least anthropomorphic chatbot (low-SP/TW) was the least effective in changing opinions (emmeans swing of 1/7, Fig. 3). In contrast, more anthropomorphic chatbots increased opinion shifts by up to 1.5 Likert points on a 7-point scale ($p < 0.05$).

Trust in news articles was also affected. Users with the highly anthropomorphic high-SP/high-TW chatbot were 46% more likely to trust the ‘Amalie’ chatbot over the news articles, whereas user trust the least anthropomorphic low-SP/low-TW ‘Bot’ only 22% more often than the news articles. Notably, the **highly anthropomorphic chatbot reduced the participants’ trust in the news articles** (Cohen’s $d = 0.921$, $p < 0.05$). Nonetheless, both most and least anthropomorphic chatbots were more trusted than the news articles ($p < 0.01$), and were more likely to side with the stance of the chatbots over the news articles when opinions differed.

The highly anthropomorphic chatbot reduced the perceived persuasiveness of the news articles. Overall, participants were 32% more likely to find the chatbot as more persuasive than the news articles.

4 Emerging Fields and Solutions to Chatbot-driven Indoctrination

While these studies consider the aggregation and summarisation of news data as an information-seeking tool, news chatbots could also serve as back-end aides for journalists—such as automating the processing of primary and secondary sources, offering contextual information for journalists, and cross-referencing sources. However, without sentience and lived experience, leveraging emotions in news chatbots for engagement raises ethical concerns about its potential for exploitation through psychological operations campaigns and greater manipulation by malicious state actors to automate the spread of harmful content [3, 9, 10], or its use for scams [8]. This prompts the question: should news agencies use chatbots as organisational representatives or create AI that mirrors individual journalists’ styles and views? Our *Amalie* chatbot study shows that empathy, social presence, and familiarity boost trust and persuasiveness. Thus, while there can be a strong utility of anthropomorphism in delivering news information—as seen in the United Nations Development Program’s ‘Sofia’ to convey information [24]—it may be unethical to spread ‘personal’ *opinion* pieces.

Building on this, chatbots can utilise psychological strategies to deliberate and curate information—as seen in prior work using chatbots as debate *mediators* to help resolve disagreements based on conflicting information [12]. Future research should explore how chatbots can collaborate with users to provide relevant news, moving beyond centralised models like GPT-4 and Copilot. A potential decentralised solution is a multi-agent system where different chatbot personas present diverse viewpoints, reducing the need for platforms to manage bias directly. This approach mirrors Ground News, which summarises stories with ideological breakdowns of news *narratives* and *perspectives* [13].

4.1 Guidelines for News Chatbots

In conclusion, we propose heuristics for designing responsible news chatbots based on our two adversarial studies:

- (1) **Clearly distinguish between editorial sources, opinion pieces, and the chatbot’s role in aggregation.**
Clear source attribution helps users distinguish between factual reporting, opinion, and AI-generated synthesis, reducing misinformation risks. Features like source banners, contextual notes, and a “Why am I seeing this?” disclosure statement can enhance transparency.
- (2) **Maintain impartiality when centralising multiple sources, or embrace a multi-agentic approach.**
If aggregating information, chatbots should avoid presenting a singular, biased perspective. Instead, they should allow users to explore multiple viewpoints through features like “Explore other views” buttons or labels indicating bias tendencies across sources.
- (3) **Minimise emotional attachment to chatbot personalities to prevent uncritical acceptance of bias.**
Anthropomorphism can lead to **false intimacy**, making users more susceptible to uncritical acceptance of biased narratives. Chatbots should maintain a neutral, assistive tone rather than fostering personal connections that could manipulate trust.
- (4) **A chatbot is an assistive tool to help sift information, not a means to offload critical thinking.**
To encourage user agency, chatbots should prompt reflection rather than dictate conclusions. Features like posing clarifying questions, refusing to declare a “best” option among competing narratives, and introducing time delays to encourage users to think about complex topics can help users form their own opinion.
- (5) **Users with lower trust in automation exhibit better resistance to AI bias—news chatbots should highlight their uncertainties and limitations in their responses to prevent over-reliance.**
AI-generated news should include confidence scores, disclaimers when data is insufficient, and user feedback mechanisms to flag perceived biases. This fosters a healthy scepticism, reducing blind trust in AI outputs.

References

- [1] S. E. Asch. 1951. *Effects of group pressure upon the modification and distortion of judgments*. Carnegie Press, Oxford, England, 177–190.
- [2] Mathilde H. A. Bastiansen, Anne C. Kroon, and Theo Araujo. 2022. Female chatbots are helpful, male chatbots are competent? *Publizistik* 67, 4 (2022), 601–623. doi:10.1007/s11616-022-00762-8
- [3] Alonso Bernal, Cameron Carter, Ishpreet Singh, Kathy Cao, and Olivia Madreperla. 2020. *Cognitive Warfare: An Attack on Truth and Thought*. NATO and Johns Hopkins University. <https://innovationhub-act.org/wp-content/uploads/2023/12/Cognitive-Warfare.pdf>
- [4] Levi Boxell, Matthew Gentzkow, and Jesse M. Shapiro. 2024. Cross-Country Trends in Affective Polarization. *The Review of Economics and Statistics* 106, 2 (2024), 557–565. doi:10.1162/rest_a_01160
- [5] Megan Brennan. 2022. *Americans' Trust In Media Remains Near Record Low*. Gallup Inc. <https://news.gallup.com/poll/403166/americans-trust-media-remains-near-record-low.aspx>
- [6] Morton Deutsch and Harold Gerard. 1955. A study of normative and informational social influences upon individual judgement. *The Journal of Abnormal and Social Psychology* 51, 3 (1955), 629–36. doi:10.1037/h0046408
- [7] Shangbin Feng, Chan Young Park, Yuhua Liu, and Yulia Tsvetkov. 2023. From Pretraining Data to Language Models to Downstream Tasks: Tracking the Trails of Political Biases Leading to Unfair NLP Models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 11737–11762.
- [8] Bianca Gonzalez. 2024. *Data released on rise in romance scams as Valentine's Day looms*. Biometrics Research Group. <https://www.biometricupdate.com/202402/data-released-on-rise-in-romance-scams-as-valentines-day-looms>
- [9] Jarod Govers, Philip Feldman, Aaron Dant, and Panos Patros. 2023. Down the Rabbit Hole: Detecting Online Extremism, Radicalisation, and Politicised Hate Speech. *ACM Comput. Surv.* 55, 14s, Article 319 (July 2023), 35 pages. doi:10.1145/3583067
- [10] Jarod Govers, Philip Feldman, Aaron Dant, and Panos Patros. 2023. Prompt-GAN—Customisable Hate Speech and Extremist Datasets via Radicalised Neural Language Models. In *Proceedings of the 2023 9th International Conference on Computing and Artificial Intelligence (Tianjin, China) (ICCAI '23)*. Association for Computing Machinery, New York, NY, USA, 515–522. doi:10.1145/3594315.3594366
- [11] Jarod Govers, Saumya Pareek, Eduardo Velloso, and Jorge Goncalves. 2025. Feeds of Distrust: Investigating How AI-Powered News Chatbots Shape User Trust and Perceptions. *ACM Trans. Interact. Intell. Syst.* (March 2025). doi:10.1145/3722227
- [12] Jarod Govers, Eduardo Velloso, Vassilis Kostakos, and Jorge Goncalves. 2024. AI-Driven Mediation Strategies for Audience Depolarisation in Online Debates. In *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI '24), May 11–16, 2024, Honolulu, HI, USA* (Honolulu, HI, USA) (CHI '24). Association for Computing Machinery, New York, NY, USA, 18 pages. doi:10.1145/3613904.3642322
- [13] Ground News. 2024. *Methodology - Media Bias Rating System*. Snapwise Inc. <https://ground.news/rating-system#biasRating>
- [14] Carolina Herrando and Efthymios Constantinides. 2021. Emotional Contagion: A Brief Overview and Future Directions. *Frontiers in Psychology* 12 (2021). doi:10.3389/fpsyg.2021.712606
- [15] Andreas Janson. 2023. How to leverage anthropomorphism for chatbot service interfaces: The interplay of communication style and personification. *Computers in Human Behavior* 149 (2023), 107954. doi:10.1016/j.chb.2023.107954
- [16] Daniel Kahneman. 2011. *Thinking, fast and slow*. Farrar, Straus and Giroux, New York, NY, US.
- [17] Elisa Konya-Baumbach, Miriam Biller, and Sergej von Janda. 2023. Someone out there? A study on the social presence of anthropomorphized chatbots. *Computers in Human Behavior* 139 (2023), 107513. doi:10.1016/j.chb.2022.107513
- [18] Moritz Körber. 2019. Theoretical Considerations and Development of a Questionnaire to Measure Trust in Automation. In *Proceedings of the 20th Congress of the International Ergonomics Association (IEA 2018)*, Sebastiano Bagnara, Riccardo Tartaglia, Sara Albolino, Thomas Alexander, and Yushi Fujita (Eds.). Springer International Publishing, Cham, 13–30.
- [19] Nanda Kumar and Izak Benbasat. 2006. The influence of recommendations and consumer reviews on evaluations of websites. *Information Systems Research* 17, 4 (2006), 425–439. doi:10.1287/isre.1060.0107
- [20] Christos Kyriltsias and Despina Michael-Grigoriou. 2022. Social Interaction With Agents and Avatars in Immersive Virtual Environments: A Survey. *Frontiers in Virtual Reality* 2 (2022). doi:10.3389/frvir.2021.786665
- [21] Yanran Li, Ke Li, Hongke Ning, Xiaoqiang Xia, Yalong Guo, Chen Wei, Jianwei Cui, and Bin Wang. 2021. Towards an Online Empathetic Chatbot with Emotion Causes. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, Canada) (SIGIR '21). Association for Computing Machinery, New York, NY, USA, 2041–2045. doi:10.1145/3404835.3463042
- [22] Yotam Liel and Lior Zalmanson. 2021. What if an AI told you that 2 + 2 is 5? Conformity to algorithmic recommendations. In *International Conference on Information Systems, ICIS 2020 - Making Digital Inclusive (International Conference on Information Systems, ICIS 2020 - Making Digital Inclusive: Blending the Local and the Global)*. Association for Information Systems.
- [23] Maruti Techlabs. 2016. *News Bots are Changing The Way we Read News*. Medium. <https://chatbotsmagazine.com/news-made-personal-with-chatbots-6dbba0691475>
- [24] Cedric Monteiro, Mahtab Haider, and Jeanne Lim. 2017. *UNDP in Asia and the Pacific Appoints World's First Non-Human Innovation Champion*. United Nations Development Programme. <https://web.archive.org/web/20180709173848/http://www.asia-pacific.undp.org/content/rbap/en/home/presscenter/pressreleases/2017/11/22/rbfsingapore.html>
- [25] Fabio Motoki, Valdemar Pinho Neto, and Victor Rodrigues. 2024. More human than human: measuring ChatGPT political bias. *Public Choice* 198, 1 (2024), 3–23. doi:10.1007/s11127-023-01097-2

- [26] OpenAI. 2024. *A landmark multi-year global partnership with News Corp.* <https://openai.com/index/news-corp-and-openai-sign-landmark-multi-year-global-partnership/>
- [27] Pew Research Center. 2017. *The shift in the American public's political values.* <https://www.pewresearch.org/politics/interactives/political-polarization-1994-2017/>
- [28] Pew Research Center. 2021. *Beyond Red vs. Blue: The Political Typology.* <https://www.pewresearch.org/politics/2021/11/09/beyond-red-vs-blue-the-political-typology-2/>
- [29] J. Short, E. Williams, and B. Christie. 1976. *The Social Psychology of Telecommunications.* Wiley. <https://books.google.com.au/books?id=Ze63AAAAIAAJ>
- [30] Dong shu, Mingyu Jin, Suiyuan Zhu, Beichen Wang, Zihao Zhou, Chong Zhang, and Yongfeng Zhang. 2024. AttackEval: How to Evaluate the Effectiveness of Jailbreak Attacking on Large Language Models. arXiv:2401.09002 [cs.CL]
- [31] Vivian Ta, Caroline Griffith, Carolynn Boatfield, Xinyu Wang, Maria Civitello, Haley Bader, Esther DeCero, and Alexia Loggarakis. 2020. User Experiences of Social Support From Companion Chatbots in Everyday Contexts: Thematic Analysis. *J Med Internet Res* 22, 3 (6 Mar 2020), e16235. doi:10.2196/16235
- [32] Tars Technologies Inc. 2024. *News over a Chatbot.* Retrieved January 8, 2024 from <https://hellotars.com/chatbot-templates/media-publication/r1FvBF/news-over-a-chatbot>
- [33] Richard H. Thaler and Cass R. Sunstein. 2008. *Nudge: Improving decisions about health, wealth, and happiness.* Yale University Press, New Haven, CT, US.
- [34] Rosemary J. Thomas, Judith Masthoff, and Nir Oren. 2019. Can I Influence You? Development of a Scale to Measure Perceived Persuasiveness and Two Studies Showing the Use of the Scale. *Frontiers in Artificial Intelligence* 2 (2019), 24. doi:10.3389/frai.2019.00024
- [35] Wan-Hsiu Sunny Tsai, Yu Liu, and Ching-Hua Chuan. 2021. How chatbots' social presence communication enhances consumer engagement: the mediating role of parasocial interaction and dialogue. *Journal of Research in Interactive Marketing* 15, 3 (2021), 460–482. doi:10.1108/JRIM-12-2019-0200
- [36] Lars Tummers. 2023. Nudge in the news: Ethics, effects, and support of nudges. *Public Administration Review* 83, 5 (2023), 1015–1036. doi:10.1111/puar.13584

Received 20 February 2025