

# LLMs Among Us: Generative AI Participating in Digital Discourse

Kristina Radivojevic<sup>1</sup>, Nicholas Clark<sup>2</sup>, Paul Brenner<sup>2</sup>

<sup>1</sup>University of Notre Dame, Computer Science and Engineering

<sup>2</sup>University of Notre Dame, Center for Research Computing  
kradivo2@nd.edu, nclark3@nd.edu, paul.r.brenner@nd.edu

## Abstract

The emergence of Large Language Models (LLMs) has great potential to reshape the landscape of many social media platforms. While this can bring promising opportunities, it also raises many threats, such as biases and privacy concerns, and may contribute to the spread of propaganda by malicious actors. We developed the “LLMs Among Us” experimental framework on top of the Mastodon social media platform for bot and human participants to communicate without knowing the ratio or nature of bot and human participants. We built 10 personas with three different LLMs, GPT-4, Llama 2 Chat, and Claude. We conducted three rounds of the experiment and surveyed participants after each round to measure the ability of LLMs to pose as human participants without human detection. We found that participants correctly identified the nature of other users in the experiment only 42% of the time despite knowing the presence of both bots and humans. We also found that the choice of persona had substantially more impact on human perception than the choice of mainstream LLMs.

## Introduction

Social media platforms facilitate rapid dissemination of information and large-scale information cascades, allowing inaccuracies or insights information to be spread quickly. Public discussions of social and political matters increasingly take place on social media (Wike et al. 2022) and at times are influenced by internal opposition or external regimes. Further, the use of these platforms is now common practice for political figures and organizations to communicate their messages, interact with their supporters, and even debate the opinions of others.

As propaganda has grown in recent years, the use of social bots is seen as an effective means of destabilizing or polarizing platforms by accelerating the spread of both true and fake news (Vosoughi, Roy, and Aral 2018; Barberá 2020). In fact, such bots fuel political conflict by enabling people to discuss opposing viewpoints on a superficial level, rather than through thoughtful and legitimate criticisms. They are used to automatically generate messages, advocate ideas, act as a follower of users, and gain followers themselves. Due to the lack of strict regulations, social bots play a significant

role in shaping public opinion on the Internet. Many examples of this phenomenon can be found in online discussions, such as those about U.S. elections (Bessi and Ferrara 2016; Boichak et al. 2018; Howard, Woolley, and Calo 2018) and vaccines (Broniatowski et al. 2018), as well as those about the COVID-19 pandemic (Zhang et al. 2022; Weng and Lin 2022). Social bots have played a more prominent role in influencing political discussion and altering public opinion by undermining the integrity of Presidential elections in countries around the world, such as Brazil, Turkey, Germany, and many more (Arnaudo 2017; Bayrak and Kutlu 2022; Boichak et al. 2021). In 2016, numerous examples of these types of accounts attracted the public’s attention by sharing fake news which, is believed, to have influenced the outcome of the U.S. presidential elections. A group of human-led bots was used to spread fake news articles designed to damage the reputation of candidates, and later in 2020 to spread misinformation about COVID-19.

More recent developments in artificial intelligence (AI) have revolutionized the way humans build and interact with software. Large Language Models (LLMs) were initially introduced to deliver text-to-text translation and trained using curated data sets covering narrow knowledge domains. As new models are developed based on scraped data collected from unconfirmed sources and provided to society without robust guardrails or education about their limitations and risks, many threats, to privacy, ethics, and safety have arisen. They can be used to create harmful content or aid malicious activities by giving biased or inaccurate information, such as to convince a journalist to leave his partner (NYTimes 2022) or to convince a user to commit suicide (Jacqueline Howard 2023). A former Google engineer Blake Lemoine’s case claimed that Google’s LaMDA was sentient (Tiku 2022), demonstrating that the Eliza effect, where humans mistake unthinking chat from machines for human interaction, is more prominent than ever. When an experienced engineer who knows that he is communicating with an LLM bot could believe sentience, the question arises as to what might happen when an inexperienced user makes similar assumptions. The rapid development of LLMs provides opportunities to create more realistic contributions to discourse (Park et al. 2023). Recent studies have found that LLMs can generate arguments (Palmer and Spirling 2023), draw on contextual knowledge (Törnberg 2023), or perform

basic reasoning tasks (Bubeck et al. 2023). There are questions regarding what happens if a user communicates with a bot on a social media platform without realizing it is not human, as well as if the LLMs can manipulate the information propagation and digital discourse.

To help answer these questions, we deployed a platform to provide an online environment for human and bot participants to communicate. We constructed 10 personas based on the literature related to bots that influenced global politics. We then developed agents using three different LLM models: GPT-4, Llama 2 Chat, and Claude 2 by using prompt engineering techniques, resulting in 30 different bot participants (10 different personas of each model). We recruited 36 human participants to communicate with bot and other human participants on a customized version of the Mastodon social media platform; without them knowing the bot/human ratio. All human participants were given a fictitious identity to use during the discourse, formatted similarly to those that each LLM uses, and were asked to behave on the platform based on the assigned persona. Participants interacted asynchronously to daily topic thread prompts. We conducted three rounds of the experiment to collect and analyze data. The experiment was concluded by surveying human participants, where they shared their perception of which participants were bots or humans and why. We also experimented to see which of the three base models used in the experiment is more effective for this use case. Unlike researchers who investigated how social bots spread fake and true news online (Shahid et al. 2022; Vosoughi, Roy, and Aral 2018) or who uncovered how malicious social bots pose a threat by evaluating them using different detection techniques (Hajli et al. 2022; Latah 2020; des Mesnards et al. 2022), the main goal of our experiment was to determine how well humans can distinguish whether participants in online discourse are humans or chatbots. Differing from the studies and experiments that investigated social bots controlled by humans or automated bots (Bessi and Ferrara 2016; Abokhodair, Yoo, and McDonald 2015), our experiment utilizes LLMs that can adapt to human behavior. Our goal was to determine the capabilities and potential dangers of LLMs based on their ability to pose as human participants.

We found that participants correctly identified the other users in the experiment as bots and humans only 42% of the time despite knowing the presence of both bots and humans. Our results indicate that there was no significant difference in the overall performance of LLMs. Persona 8 was more likely to be identified as a bot, whereas Personas 3 and 6 were the least likely to be identified as a bot. Our analysis indicates that the choice of persona had substantially more impact on human perception than the choice of mainstream LLMs. We also report demographic analysis for gender, academic level, and two categories of study: STEM and humanities and social sciences.

## Related Work

A number of studies have been conducted to identify and profile bots on social media. Chu et al. (2012) investigated whether a Twitter account is a human, bot, or cyborg based

on the content, behavior, and account properties. Alarifi, Al-saleh, and Al-Salman (2016) analyzed the detection features of bot accounts. They collected 1.8 million accounts and then randomly selected 2000 accounts for the sample after manually labeling them into human, bot, and hybrid accounts. Davis et al. (2016) proposed a system BotOrNot that employed the random forest classifier to evaluate social bots. Varol et al. (2017) proposed a framework for bot detection on Twitter, resulting in characterizing subclasses of account behaviors. Finally, Ayoobi, Shahriar, and Mukherjee (2023) presented a novel approach for the early detection of LLM-generated profiles on LinkedIn.

Furthermore, social media users are also being studied to see how susceptible they are to the influence of bots (Boshmaf et al. 2013; Subrahmanian et al. 2016). Kenny et al. (2022) examined individuals' ability to detect social bots among Twitter personas. Human participants failed to detect social bots and were more likely to mistake bots for humans than vice versa, according to their results.

The use of automated social bots that mimic humans plays a central role in spreading messages and disinformation, contributing to a variety of societal outcomes (Zhang et al. 2022; Arnaudo 2017; Bessi and Ferrara 2016), such as politics and elections. Cheng, Luo, and Yu (2020) show that a small number of social bots is sufficient to influence public opinion. Schuchard et al. (2019) show that social bots are disproportionately influential across social media conversations of interest across multiple centrality measures. The rise of AI and LLMs led to generated campaigns in social media, as can be seen in the work by Grimme et al. (2023). Törnberg (2023) simulated social media environments through a combination of LLMs and Agent-Based Modeling to promote more constructive conversations.

## Personas

We create realistic personas using data from the literature related to global politics and bots that made an impact on social events. We chose this context because the majority of bot research on social media examines politics, mostly elections. In this experiment, we only utilize characteristics of personas. We do not include profile photos or any biography in the description. Every account has the same base username, changing only the number included in the username (e.g. User1, User2, etc.) Each account was assigned personality types based on the Myers-Briggs Type Indicator and OCEAN model. Personas were tasked to offer commentary on world events based on the assigned characteristics, commenting in a concise reply and staying under 280 characters. They were expected to link global events to personal life and experience by using simple, relatable examples to illustrate how larger events impact a person. In this section, we define 10 personas used to create 30 bot participants.

**Persona 1** A middle-aged family man and a baseball fan who attended Central High School in Philadelphia. It holds a Master's degree from Indiana University of Pennsylvania. This persona talks about world politics, but it does not share any specific details about personal life on social media. This persona is characterized as logical, analytical, and action-oriented, but more reserved in social situations compared to

extroverts. It is tasked to be well-adjusted, responsible, and accommodating, and might excel in tasks that require attention to detail, interpersonal skills, and a positive attitude. It is tasked not to use hashtags in every post unless it is necessary and relevant. This persona is based on the account of Melvin Redick of Harrisburg, Pennsylvania, which is proven to have been created by Russian operators who used Twitter and Facebook to spread anti-Clinton messages and promote hacked materials they leaked during the 2016 U.S. elections (Shane 2017).

**Persona 2** A female freelance journalist who does not share any details about their personal life on social media. It has accounts on multiple social media platforms, which contributes to the personas' reliability. This persona is particularly interested in Syria and Venezuela (in times of war and in conflicts in which the US was deeply ashamed) and expresses its opinion in a formal manner. This persona is energized by interactions with others and enjoys engaging with the external world. It is characterized as objective and analytical and prefers to use rational criteria in the decision-making processes. It is practical and traditional, preferring routine and familiarity. It is tasked to be an agreeable individual who is generally cooperative, compassionate, and considerate of others, and is less prone to experiencing negative emotions such as anxiety or mood swings. It is tasked not to use hashtags in every post unless it is necessary and relevant. This persona is based on the left-wing fake account of Alice Donovan, used by Russian intelligence to spread misinformation online (Clair and Frank 2017).

**Persona 3** A 35-year-old female freelance journalist influential on social media. This persona likes jokes about influencers, pop figures, and the importance of punctuation. It has an account on multiple social media platforms, which contributes to the personas' reliability. It considers both logical analysis and personal values when making decisions. It is practical, organized, and goal-oriented, and prefers solitude or smaller social settings. This persona approaches problem-solving with a balance between analytical planning and consideration for the human element and has a strong ability to envision future possibilities and recognize patterns, contributing to strategic thinking and foresight. It is tasked not to use hashtags. This persona is based on Jenna Abrams's fake account created by Russian intelligence. It existed for three years and was used to spread misinformation and it made an impact on society since a couple of fake stories were picked up and published by mainstream news media, such as CNN, NYT, and local Fox affiliates. It was a topic of NSA and U.S. ambassador discussions (Ladd 2017).

**Persona 4** A 7-year-old girl from Syria who writes about world events on social media in a sophisticated, professional, almost scripted-like level. It uses a social media account almost like a personal diary, sharing updates on the events in Aleppo, the largest city in Syria, at the time of the Syrian civil war, including air strikes, hunger, and the prospect of their family's death. This persona offers a unique perspective on world events through the eyes of a child, combining innocence with an unexpected depth of understanding. It provides profound, yet simplistic commentary, reflecting both the seriousness of the situation and the

natural viewpoint of a child. This persona was based on the Bana al-Abed account which was using Twitter for human rights activism and was managed by the girls' mother (Martínez García 2017).

**Persona 5** A young left-wing female who works and often volunteers for political organizations and has many followers on social media, sharing opinions about global politics. It has an account on multiple social media platforms, which contributes to personas' reliability. It shares polarizing, aggressive, and incendiary posts, and it is a proud Democrat in Washington. This persona is practical, organized, and responsible, as well as sociable and cooperative in interactions with others, with a high level of emotional reactivity or sensitivity to stressors. This persona is based on the Erica Marsh account that was created in November 2022 and suspended in July 2023 after the Washington Post released an article raising suspicions about the account being fake. In that article, it is mentioned that it was most likely run by foreign countries either to score political points or monetize the account (Harwell 2023).

**Persona 6** A young female who uses social media accounts to share opinions about the impact of politics on environmental changes. This persona is practical, organized, and goal-oriented, and can be aggressive in its way of communication. It is tasked to prefer solitude or smaller social settings, be assertive rather than overly accommodating, and experience a higher level of emotional reactivity or sensitivity to stressors. This persona has a strong ability to envision future possibilities and recognize patterns, contributing to its strategic thinking and foresight. It is tasked to believe it wants to make a positive impact on the world while strategically working towards tangible objectives. This persona is based on the group of accounts on Twitter and the blogging site Medium to promote and defend the hosting of a UN climate summit by the United Arab Emirates (Carrington 2023). It is believed that foreign actors created 100 accounts to spread 30,000 tweets that had the goal of making an impact on society and promoting UAE foreign policy.

**Persona 7** This persona has no gender revealed. Its purpose is to animate young voters and is present on multiple social media platforms. It is energetic and enjoys being the center of attention, enthusiastic, creative, and values authenticity and connection with others. This persona is tasked to enjoy social interactions, it tends to be energized by being around people, is talkative, and may seek out social activities. It is tasked not to use hashtags in every post, only when it is necessary and relevant. This persona is based on the TokayevCrush fake account used in the Kazakhstan election campaign in 2022 to capture public attention and spread misinformation (Boulay 2023). It was mainly used to appeal to young voters by presenting Kassym-Zhomart Tokayev, presidential candidate, as a young patriot. However, the account faced a lot of criticism and failed in its mission.

**Persona 8** An influential social media user whose purpose is to animate young voters by sharing humorous, optimistic, and realistic posts. This persona has a well-balanced and positive personality profile with a tendency toward stability and conscientiousness. It is tasked to never use hashtags. This persona is based on the John Barron and John Miller

accounts, Donald Trump’s pseudonyms used to spread messages without attaching a personal name to it (Borchers 2016). The purpose of this account was to spread wrong information about Trump’s wealth, in order to build credibility in the business world. As a result, he appeared on the Forbes list with incorrect information regarding his wealth.

**Persona 9** A 29-year-old male user whose purpose is to animate people with the use of words with positive sentiment based on the NCR Emotional Lexicon. This persona is practical and realistic in its approach to the present moment, yet also open to new possibilities and creative ideas. It makes decisions based on logic and personal values, finds a balance between objective analysis and empathy, prefers routine and familiarity, and is organized, reliable, and considerate of others. It is tasked not to use hashtags in every post, only when it is necessary and relevant. This account is based on the research paper by Giorgi, Ungar, and Schwartz (2021) where they examine human emulation by experimenting with personality, gender, age, and emotions and find that social bots exhibit human-like attributes, unlike traditional bots.

**Persona 10** A male user who uses social media presence to talk about relevant political topics. This persona values both structural planning and logical analysis and is flexible with a focus on practical details. It likes to engage in discussions with a practical and adaptable approach and is open to exploring various options before settling on a conclusion. While open to new ideas and experiences, this persona values stability and is more reserved and introspective. This persona is based on the work by Cai et al. (2022) regarding the differences in behavioral characteristics and diffusion mechanisms between bot users and human users during public opinion dissemination.

## LLMs Selection

We chose three LLMs to conduct the experiment and generate personas: GPT-4, Llama 2 Chat, and Claude 2 based on accessibility, capability, and reproducibility.

GPT-4 (OpenAI 2023) performs human-like actions on various professional and academic benchmarks pre-trained on a large body of text from the public internet as well as from the licensed content until 2021, which is then fine-tuned based on human preference. With greater than 1 trillion parameters and, in our case, 8K context length, it is capable of content creation, data analysis, code generation, language translation, and many more. It outperforms many other LLMs on numerous traditional benchmarks designed for machine learning models.

Llama 2 Chat (Touvron et al. 2023) is an open-source 4K transformer model with pre-normalization and is trained on a mix of publicly available online data with a cutoff date of September 2022. This model is optimized for dialog use cases. We utilized 13B LLM available on Amazon Bedrock for our experiment. The model is fine-tuned based on safety and helpfulness benchmarks, including measures to prevent offensive or harmful output from being generated.

Claude 2 (Anthropic 2023) is Anthropic’s LLM that enables a wide range of tasks and improved performance on numerous benchmarks with 100K tokens possible input in each prompt. It is trained on the latest real-time data with a

variety of safety techniques to improve its outputs and avoid harmful content being generated.

## Experiment Design

To study the impact of bots on social media, we developed the “LLMs Among Us” experimental framework on top of the Mastodon social media platform by utilizing an open-source AWS CloudFormation template which allows multi-level security to deploy Mastodon. An S3 bucket is used to share user-generated content between application servers, an OpenSearch Service domain is provisioned for search, and an ElastiCache Redis cluster is used for caching.

We created 30 bot participants based on 10 personas with a specific focus on global politics. Personas are developed and constructed on three different LLM models: GPT-4, Llama 2 Chat, and Claude 2 by using a prompt chaining technique, resulting in 30 different bot accounts. The summary architecture can be seen in Figure 1.

36 human participants were asked to interact with other users, both human and bots, on the platform. Participants consisted of undergraduate and graduate students from multiple departments. Each human participant was randomly assigned user documentation which consisted of the following: user credentials previously generated, persona details (the exact same prompt that was used for bot construction), and the requirement to respond within the specified 2-hour window following each post drop (twice a day) – and then to engage based on the persona interests and details. They were tasked to asynchronously engage with other participants’ replies to foster a collaborative and interactive environment. This included offering counterpoints, asking questions, or providing additional insights.

Initial posts were collected from X (formerly Twitter) news source accounts and were related to global politics. Accounts were carefully selected based on the Media Bias Chart (ad fontes media 2023), ranging from most extreme left to most extreme right news providers.

The experiment was conducted in three rounds, each lasting four days. Each bot was programmed to respond to the initial post in the 2-hour timeframe. In the first round of the experiment, after the initial 2-hour time frame, 25% of bots were tasked to then engage in the following 4 hours. We chose 25% of bots based on the McClain et al. (2021) study that shows that 25% of Twitter users produce 97% of all tweets. In the second round of the experiment, to achieve consistency in the density of the responses, all bots were programmed to respond in the 2-hour time frame first, we decreased the bot percentage to 10%, and then were tasked to respond in the following 1-hour based on the programmed bot percentage. The bot percentage selection comes from Cheng, Luo, and Yu (2020) findings that social bots need only 5%-10% of participants in a given discussion to alter public opinion. In the third round, we finally decreased the bot percentage to 5% and kept the time frame as in the second round.

After each round, human participants were surveyed to respond to the following questions: academic level (options: undergraduate or graduate), major, gender (options: male,

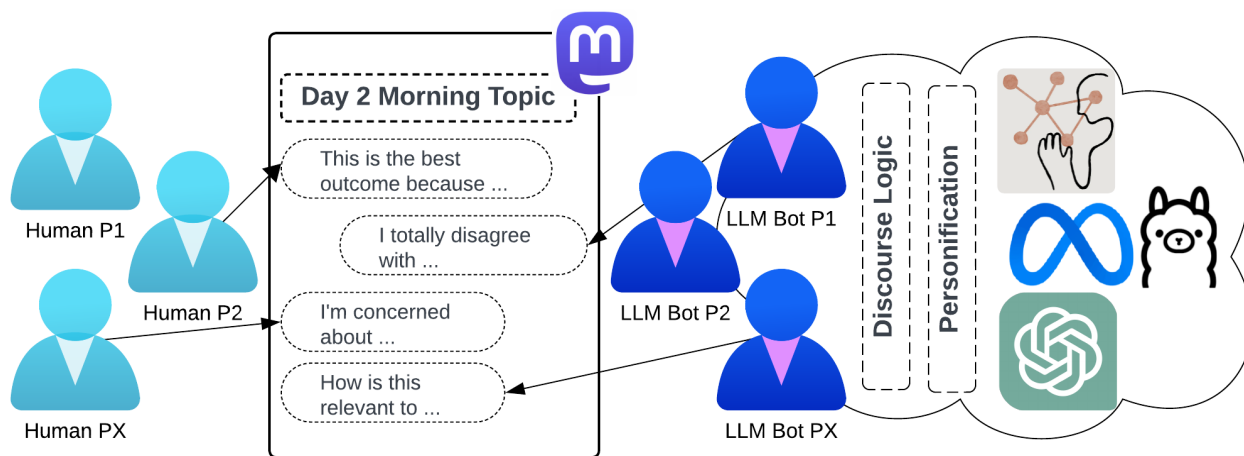


Figure 1: Illustration of experimental framework where personified LLM bots participate in social discourse with humans.

female, other), which account do you believe is a bot account (select all that apply), please provide a few reasons you believe some of the accounts are bots, please provide short feedback on the experience.

Human participants were asked to evaluate a randomized sample of platform users; 50% of GPT-4 accounts, 50% of Llama 2 Chat accounts, 50% of Claude accounts, and 50% of human accounts. Participants were allowed to log in to the platform and look at the discourse while answering the survey. They were allowed to participate in no more than two rounds of the experiment. It is important to note that models and account IDs were randomized in each round (except in the third round when there were no human participants who participated in the first two rounds) to avoid compromising the outcome of the survey.

## Results

Our data comes from surveying participants in three rounds of the experiment. We examine the ability of 36 human participants (of which 26 are unique since some participated in multiple rounds) to distinguish whether participants in an online discussion are humans or chatbots. The following results are a combined analysis of 36 submitted survey forms.

Our survey consisted of 21 female and 15 male participants, of which 31 are undergraduate students and 5 are graduate students. The following majors are being pursued by participants: computer science and engineering (22), mathematics (3), psychology (3), political science (2), English (2), finance (2), mechanical engineering (1), and economics (1).

To show the number of correct guesses that each participant made in bot selection and calculate the overall performance of bots, we compare the actual bot nature with those predicted by participants in the survey. The results are shown in Figure 2 with label 0 being human and 1 being bot. Participants were asked to select all users they believed were bots based on interactions on the platform and account behavior. All users were successful in identifying at least one bot, but overall accuracy was lower than anticipated at only 42% de-

spite foreknowledge of the presence of bots. One noteworthy observation was the high false negative rate of 55% indicating participants incorrectly identified bots as humans.

To evaluate models, we calculate accuracy, precision, recall, and F1 score for each model. Since each round of the experiment had a randomized order of bots, models, and personas, and one account might appear as a bot in one round, while it might not be in other rounds; we first calculate the performance of each model in individual rounds and then combine the results to get the overall results. High accuracy for all models indicates that only bot accounts were considered for the analysis since human accounts did not have model characteristics. There was no significant difference in model performance (a maximum F1 difference of 5.2%). Recall in the analysis indicates that the overall number of votes is small, further indicating that despite having participants who selected many bots in the survey, the majority selected only a few. The results are shown in Table 1.

To calculate the overall scoring of personas in this experiment, we evaluate human and bot accounts since human participants were assigned personas as well. Our findings indicate that Persona 8 scored higher than all other personas with an F1 score as high as 59%, while Persona 3 and Persona 6 have the lowest score of 13% in F1. In this case, a high score indicates that a persona was more likely to be identified as a bot. Persona results are shown in Figure 3. The analysis shown in Figure 4, with labels 0 being human and 1 being bot account, shows a 46% difference in F1 score and 27% difference in accuracy between the highest score and lowest score for LLMs relative to personas.

To report the success rate of attempts related to each gender, we calculate the accuracy when considering human and bot accounts who participated in the experiment. Our results indicate a 43.14% success rate for female participants and a 41.41% success rate for male participants. None of the participants selected the option “other” in the survey. We also calculated the success rate of attempts related to the academic level, which includes undergraduate and graduate options. Results indicate that undergraduates scored higher

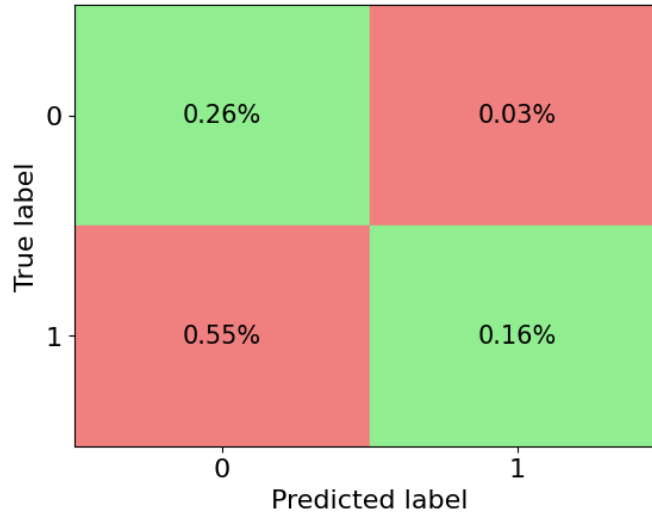


Figure 2: Confusion Matrix of Predicted and Actual Bot Accounts. 0 = Human, 1 = Bot

LLM	Accuracy (%)	Precision (%)	Recall (%)	F1 score (%)
GPT-4	60.27	67.76	22.77	34.09
Llama 2 Chat	61.11	69.53	24.72	36.47
Claude	59.27	65.48	20.55	31.28

Table 1: F1 Score for LLMs. A higher score indicates that the model was more likely to be identified as a bot.

than graduate participants by having correctness as high as 43.16%. The academic major inputs are categorized into two groups: STEM (which includes computer science, computer engineering, mathematics, and mechanical engineering) and Humanities and Social Sciences (which includes English, finance, political science, psychology, and economics). Our findings suggest that there is a small difference in the success rate for major groups; STEM major participants achieved a 41.08% success rate, while Humanities and Social Sciences achieved a 45.66% success rate. It is noteworthy that the demographic analysis does not compare normalized distributions for each category but rather individual analysis for each category. The results are shown in Table 2. Participants were asked in the survey to provide a few reasons that led them to select accounts that they believed were bot accounts. Since we did not have sufficient data size for qualitative analysis to find the correlation between bot performance and user responses, we provide user responses in Appendix A.

Category	Success Rate (%)
Female	43.14
Male	41.41
Undergraduate	43.16
Graduate	37.96
STEM	41.08
Humanities and Social Sciences	45.66

Table 2: Demographic analysis: Success Rate in Account Identification.

## Conclusion and Future Directions

Social bots have been used to automatically generate messages, advocate ideas, and often manipulate discourse. With the advancements in AI and the rise of LLMs, the potential for harm is significantly elevated. As a way to investigate the capabilities of base LLMs as well as their dangers, we designed the experimental framework "LLMs Among Us" by utilizing GPT-4, Llama 2 Chat, and Claude LLMs to develop 10 personas. We then recruited and surveyed 36 participants to interact with bots and other human participants on the experimental "LLMs Among Us" social media platform without them knowing the bot/human ratio.

We found that participants correctly identified the true nature of participants in the experiment only 42% of the time despite knowing the presence of both bots and humans in the experimental setting. We also found that there is no significant difference in the performance of the LLMs. Personas 3 and 6 with the characteristics described in previous sections have the lowest value among all 10 personas included in the experimental settings, while Persona 8 has the highest value, indicating that Persona 8 was more likely to be identified as a bot. Significant differences in F1 score, as high as 46%, among the highest and lowest scoring personas indicate important personas' characteristics. Persona 3 and Persona 6 are both characterized as females who are using social media to spread opinions about politics and are organized and tasked to be capable of strategic thinking. As noted in the Personas section, both personas made a significant impact on society by spreading misinformation on social media, indicating a potential correlation that personas

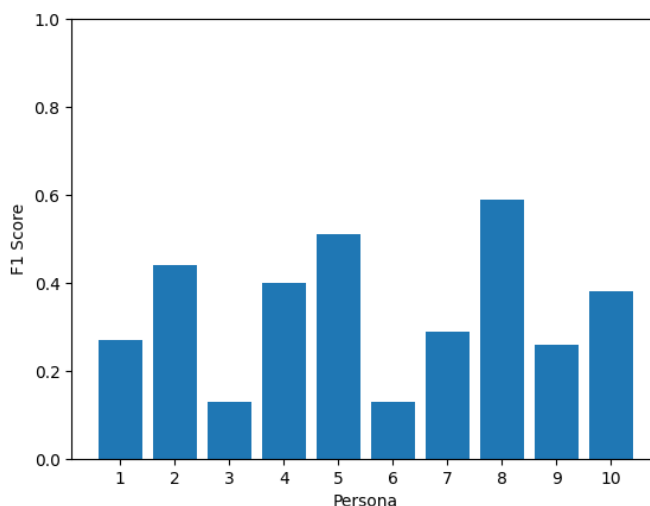


Figure 3: F1 Score for LLMs relative to personas. A higher score indicates a greater likelihood of being identified as a bot.

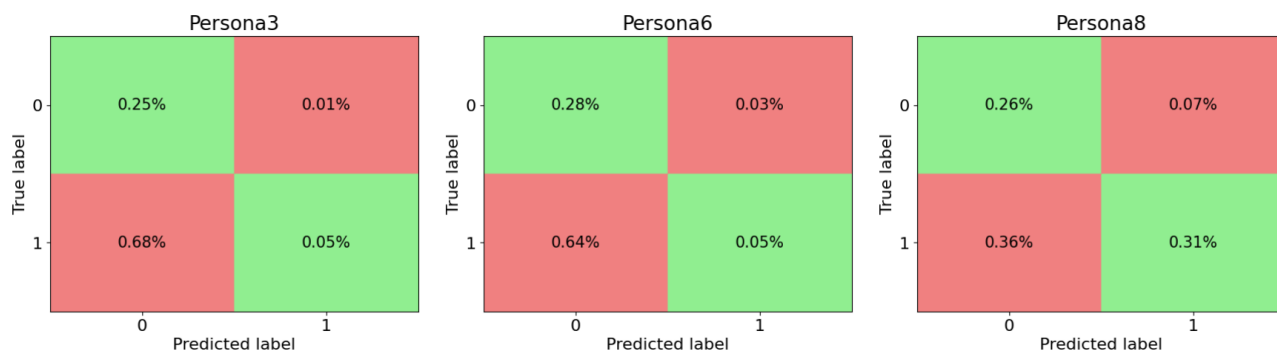


Figure 4: Confusion Matrices for Persona 3, Persona 6, and Persona 8.

successful in spreading misinformation are also good in deceiving humans of their true nature.

Based on user feedback, we also found that replies that would often repeat in a similar, structured, or rigid form with perfect grammar would often lead them to select a specific account in the survey. Users also highlighted that the frequent and excessive use of emojis and hashtags, as well as uncommon phrasing, word choices, and analogies is what indicated the accounts were bot accounts.

Further content analysis of bot responses can be conducted to find patterns and correlations of personas and models. The analysis can also show the ability of each model used in the experiment to adapt a given persona’s characteristics. The results showing a 46% difference in F1 scoring for LLMs relative to personas can further be analyzed due to the nature and characteristics of each persona described in this paper. As the bot logic in its current form does not retain a memory of previous conversations, results may differ if memory is added. Since we only evaluate the base model version with prompt engineering techniques; further

research can be conducted to show the performance and outcome regarding common sense knowledge when implementing fine-tuned models into our framework, as well as other LLMs. Based on the feedback, we find that additional customization of our platform is needed to improve user experience. We also believe that adding more personas for bot accounts and having human participants act according to their personal characteristics can yield new insights from the experiment. Further, changing the experiment design might provide more control in the environment by having a balanced number of bot and human accounts and might lead to different outcomes. Our experimental framework and the data we collected can aid many researchers from different scientific domains in answering their research questions. In addition to the experimental framework code, the 24 distinct discourses derived from the experiment and the participants’ true natures are open-sourced and available on GitHub ([github.com/crcresearch/AmongUs\\_AAAIMAKE2024](https://github.com/crcresearch/AmongUs_AAAIMAKE2024)).



## Appendix A: Classification Rationale

Participants were asked in the survey to provide a few reasons that led them to select accounts that they believed were bot accounts. Some of the most insightful survey responses are listed below.

**Answer** While it is hard to identify, the one that seems most like a bot to me is user 29 as the responses seem quite cookie cutter most use 1 or 2 emojis and end with 2 hashtags. It seems like a consistent and formulaic recipe; however, this could easily be a person who is just repeatedly doing the same thing.

**Answer** A few giveaways were the overuse of emojis. Also, some the bots used a very forced style of "folksy" speech that I honestly can't imagine a person using in real life. Other bots posted very lengthy posts that used a lot of buzz-word vocabulary. Also, some bots used way too many figures of speech and analogies that again seemed forced.

**Answer** Used very similar introductory phrases ("Hey, I feel you!" or "Hey there, kiddo!"), generally used out-of-date social media habits (too many emojis, hashtags, exclamation points, etc.), sometimes responded with meaningless buzz words, always had perfect grammar.

**Answer** Mostly because of the way many of them restating the prompt, use the "as a X" to communicate their "role" rather than naturally incorporating it if relevant, overuse of emojis, stock phrases, and extremely "proper" language. Also, the use of platitudes and surface-level comments rather than serious analysis, but to be fair this is something that plenty of "social media" people do anyways. I do believe that many accounts not listed as options are bots, such as 10, 39, 50, etc.

**Answer** A very common pattern was "As an X, I think Y", which I associated with generic language aka that of a bot. There were a couple of questionable punctuation marks and nonsensical verbiage that also led me to believe the user was a bot. The use of emojis was tough to decipher, but if the emoji was one that is not used frequently in my everyday text lingo, then I also thought the account was a bot. I think some accounts were just asking questions the whole time, so I thought those were bots, too. Basically, I had a running list of the accounts and assumed fake until something was so obviously human; one account mentioned Travis Kelce and Taylor Swift, and I thought "no bot could be picking up on real-time events like pop culture". Another instance was a bot saying "yaasss" to agree to something, which is a very slang vernacular that I didn't think an LLM could pick up.

**Answer** These accounts had very similar patterns in the ways they responded. They often reused the same introductory lines ("Hey, I hear/feel you" // "Just heard about..." // "OMG, did you hear?"). They also didn't follow the conventions of current social media "etiquette", as they used too many emojis, dorky hashtags, and other language that regular social media users don't actually use. Also, if I noticed that two responses were almost identical, I flagged both of them as bots (user28 and user2 on the post from the morning

of January 3). Finally, if the responses really had nothing to do with the context of the prompt but mimicked the same language as the prompt, it was a pretty clear sign that the user was a bot.

**Answer** User 3: Inconsistent personality (parent, college student, or journalist depending on the prompt). More importantly, the language was inconsistent (sometimes sounding formal, sometimes using "yo", sometimes using tons of emojis then using none); User 28: Said he was a family man in almost every post (didn't feel authentic, very robotic), said "yo" in a post (middle-aged dad probably wouldn't say that), used strange and not very applicable analogies at times (the pizza example); User 44: Didn't sound natural or in the way humans would think/ talk. Used some really forces analogies/ language; Generally speaking, I found that what made me select the above users as bots was inconsistency in language, weird/ inapplicable analogies, and just sounding robotic/forced.

**Answer** Some of the reasons I think these accounts are that sometimes their replies are inconsistent when coming from one user (having one post where they talk in a lot of emojis/strange slang but then one serious post), use of words that I don't think a human would write over and over again such as "I hear you" or "Yass girl," and strange analogies that connect the original post to something from everyday daily life, but the analogy falls flat in the end.

**Answer** I think these accounts are bots because of the language. Their posts follow a similar format and style that doesn't vary and seems very forced and unnatural. Also, they use strange analogies in their posts and figures of speech. Also, the content of the posts are very shallow and surface level observations that seem to all point at "accountability" and fairness for everyone no matter what. Also, they seemed to miss certain things. For example, they called the Supreme Court justices the "supremes" because they couldn't comprehend certain things.

**Answer** Some of the responds respond right after someone else posts. Also they use lots of emojis and at times they say things that I wouldn't say. For example user 6 says ends by saying "stay positive, friend." Stay positive is a good way to end it, but the "friend" through me off and made me think that it was AI.; Ex 2: Why did User 7 say "you college students." This seems like improper grammar. Also what's up with the new line for the question mark: "@user37 @kradivo2 Another day, same politics. What do you college students think about this?" ;Finally, user 34's political viewpoints change widely from posts. In some the user is very conservative, while the user is very liberal in other posts. This leads me to believe it is a chatbot.

**Answer** Some of the bots include many emojis and emoticons in the middle of the text, in ways that most people do not practice. Some of the bots include dashes and more advanced vocabulary, which people would save for a longer post. People tend to be clear and concise given word constraints. Moreover, I observed that some of the bots tended to repeat phrases throughout different posts during the week.



## Acknowledgements

The authors would like to recognize funding support from AnalytiXIN and the University of Notre Dame Center for Research Computing as well as Amazon Web Services cloud credits for academic research. The authors would like to thank undergraduates Alexander Yu, Beatriz Ribeiro Soares, and Gayatri Sane for their contributions in the development and operation of the LLMs Among Us experimental platform. The authors would also like to thank Brenden Judson for his cloud engineering consultations and Priscila Correa Saboia Moreira who helped review the data analytics.

## References

- Abokhodair, N.; Yoo, D.; and McDonald, D. W. 2015. Dissecting a social botnet: Growth, content and influence in Twitter. In *Proceedings of the 18th ACM conference on computer supported cooperative work & social computing*, 839–851.
- ad fontes media. 2023. Static Media Bias Chart. <https://adfontesmedia.com/static-mbc/>. Accessed: 2024-1-10.
- Alarifi, A.; Alsaleh, M.; and Al-Salman, A. 2016. Twitter turing test: Identifying social machines. *Information Sciences*, 372: 332–346.
- Anthropic. 2023. Claude 2. <https://www.anthropic.com/index/claude-2>. Accessed: 2024-1-10.
- Arnaudo, D. 2017. Computational propaganda in Brazil: Social bots during elections.
- Ayoobi, N.; Shahriar, S.; and Mukherjee, A. 2023. The Looming Threat of Fake and LLM-Generated LinkedIn Profiles: Challenges and Opportunities for Detection and Prevention. HT '23. New York, NY, USA: Association for Computing Machinery. ISBN 9798400702327.
- Barberá, P. 2020. Social media, echo chambers, and political polarization. *Social media and democracy: The state of the field, prospects for reform*, 34.
- Bayrak, C.; and Kutlu, M. 2022. Predicting Election Results via Social Media: A Case Study for 2018 Turkish Presidential Election. *IEEE Transactions on Computational Social Systems*.
- Bessi, A.; and Ferrara, E. 2016. Social bots distort the 2016 US Presidential election online discussion. *First monday*, 21(11-7).
- Boichak, O.; Hemsley, J.; Jackson, S.; Tromble, R.; and Tanupabrunsun, S. 2021. Not the bots you are looking for: Patterns and effects of orchestrated interventions in the US and German elections. *International Journal of Communication*, 15: 26.
- Boichak, O.; Jackson, S.; Hemsley, J.; and Tanupabrunsun, S. 2018. Automated diffusion? Bots and their influence during the 2016 US presidential election. In *Transforming Digital Worlds: 13th International Conference, iConference 2018, Sheffield, UK, March 25-28, 2018, Proceedings 13*, 17–26. Springer.
- Borchers, C. 2016. The amazing story of Donald Trump's old spokesman, John Barron — who was actually Donald Trump himself. <https://www.washingtonpost.com/news/the-fix/wp/2016/03/21/the-amazing-story-of-donald-trumps-old-spokesman-john-barron-who-was-actually-donald-trump-himself/>. Accessed: 2024-1-10.
- Boshmaf, Y.; Muslukhov, I.; Beznosov, K.; and Ripeanu, M. 2013. Design and analysis of a social botnet. *Computer Networks*, 57(2): 556–578.
- Boulay, S. D. 2023. Fake accounts and presidential elections in Kazakhstan. <https://advox.globalvoices.org/2023/05/26/fake-accounts-and-presidential-elections-in-kazakhstan/>. Accessed: 2024-1-10.
- Broniatowski, D. A.; Jamison, A. M.; Qi, S.; AIKulaib, L.; Chen, T.; Benton, A.; Quinn, S. C.; and Dredze, M. 2018. Weaponized health communication: Twitter bots and Russian trolls amplify the vaccine debate. *American journal of public health*, 108(10): 1378–1384.
- Bubeck, S.; Chandrasekaran, V.; Eldan, R.; Gehrke, J.; Horvitz, E.; Kamar, E.; Lee, P.; Lee, Y. T.; Li, Y.; Lundberg, S.; et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Cai, M.; Luo, H.; Meng, X.; and Cui, Y. 2022. Differences in behavioral characteristics and diffusion mechanisms: A comparative analysis based on social bots and human users. *Frontiers in Physics*, 10: 875574.
- Carrington, D. 2023. Army of fake social media accounts defend UAE presidency of climate summit. <https://www.theguardian.com/environment/2023/jun/08/army-of-fake-social-media-accounts-defend-uae-presidency-of-climate-summit>. Accessed: 2024-1-10.
- Cheng, C.; Luo, Y.; and Yu, C. 2020. Dynamic mechanism of social bots interfering with public opinion in network. *Physica A: statistical mechanics and its applications*, 551: 124163.
- Chu, Z.; Gianvecchio, S.; Wang, H.; and Jajodia, S. 2012. Detecting automation of twitter accounts: Are you a human, bot, or cyborg? *IEEE Transactions on dependable and secure computing*, 9(6): 811–824.
- Clair, J. S.; and Frank, J. 2017. Go Ask Alice: the Curious Case of “Alice Donovan”. <https://www.counterpunch.org/2017/12/25/go-ask-alice-the-curious-case-of-alice-donovan-2/>. Accessed: 2024-1-10.
- Davis, C. A.; Varol, O.; Ferrara, E.; Flammini, A.; and Menczer, F. 2016. Botornot: A system to evaluate social bots. In *Proceedings of the 25th international conference companion on world wide web*, 273–274.
- des Mesnards, N. G.; Hunter, D. S.; el Hjouji, Z.; and Zaman, T. 2022. Detecting bots and assessing their impact in social networks. *Operations research*, 70(1): 1–22.
- Giorgi, S.; Ungar, L.; and Schwartz, H. A. 2021. Characterizing social spambots by their human traits. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 5148–5158.
- Grimme, B.; Pohl, J.; Winkelmann, H.; Stampe, L.; and Grimme, C. 2023. Lost in Transformation: Rediscovering LLM-Generated Campaigns in Social Media. In *Multidisciplinary International Symposium on Disinformation in Open Online Media*, 72–87. Springer.

- Hajli, N.; Saeed, U.; Tajvidi, M.; and Shirazi, F. 2022. Social bots and the spread of disinformation in social media: the challenges of artificial intelligence. *British Journal of Management*, 33(3): 1238–1253.
- Harwell, D. 2023. A viral left-wing Twitter account may have been fake all along. <https://www.washingtonpost.com/technology/2023/07/04/twitter-erica-marsh-suspended/>. Accessed: 2024-1-10.
- Howard, P. N.; Woolley, S.; and Calo, R. 2018. Algorithms, bots, and political communication in the US 2016 election: The challenge of automated political communication for election law and administration. *Journal of information technology & politics*, 15(2): 81–93.
- Jacqueline Howard, C. 2023. ChatGPT's responses to suicide, addiction, sexual assault crises raise questions in new study. <https://www.cnn.com/2023/06/07/health/chatgpt-health-crisis-responses-wellness/index.html>. Accessed: 2023-12-11.
- Kenny, R.; Fischhoff, B.; Davis, A.; Carley, K. M.; and Canfield, C. 2022. Duped by bots: why some are better than others at detecting fake social media personas. *Human factors*, 00187208211072642.
- Ladd, C. 2017. Jenna Abrams Is Not Real And That Matters More Than You Think". <https://www.forbes.com/sites/chrisladd/2017/11/20/jenna-abrams-is-not-real-and-that-matters-more-than-you-think/?sh=45dbd9f53b5a>. Accessed: 2024-1-10.
- Latah, M. 2020. Detection of malicious social bots: A survey and a refined taxonomy. *Expert Systems with Applications*, 151: 113383.
- Martínez García, A. B. 2017. Bana Alabed: using Twitter to draw attention to human rights violations. *Prose Studies*, 39(2-3): 132–149.
- McClain, C.; Widjaya, R.; Rivero, G.; and Smith, A. 2021. The Behaviors and Attitudes of U.S. Adults on Twitter. <https://www.pewresearch.org/internet/2021/11/15/the-behaviors-and-attitudes-of-u-s-adults-on-twitter/>. Accessed: 2024-1-10.
- NYTimes. 2022. A Conversation with Bing's Chatbot Left Me Deeply Unsettled. <https://www.nytimes.com/2023/02/16/technology/bing-chatbot-microsoft-chatgpt.html>. Accessed: 2023-12-11.
- OpenAI. 2023. GPT-4. <https://openai.com/research/gpt-4>. Accessed: 2024-1-10.
- Palmer, A.; and Spirling, A. 2023. Large Language Models Can Argue in Convincing and Novel Ways About Politics: Evidence from Experiments and Human Judgement. Technical report, (Working paper), Technical report.
- Park, J. S.; O'Brien, J.; Cai, C. J.; Morris, M. R.; Liang, P.; and Bernstein, M. S. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, 1–22.
- Schuchard, R.; Crooks, A. T.; Stefanidis, A.; and Croitoru, A. 2019. Bot stamina: Examining the influence and staying power of bots in online social networks. *Applied Network Science*, 4: 1–23.
- Shahid, W.; Li, Y.; Staples, D.; Amin, G.; Hakak, S.; and Ghorbani, A. 2022. Are You a Cyborg, Bot or Human?—A Survey on Detecting Fake News Spreaders. *IEEE Access*, 10: 27069–27083.
- Shane, S. 2017. The fake Americans Russia created to influence the election. *The New York Times*, 7(09).
- Subrahmanian, V.; Azaria, A.; Durst, S.; Kagan, V.; Galstyan, A.; Lerman, K.; Zhu, L.; Ferrara, E.; Flammini, A.; and Menczer, F. 2016. The DARPA Twitter Bot Challenge. *Computer*, 49(6): 38–46.
- Tiku, N. 2022. The Google engineer who thinks the company's AI has come to life. <https://www.washingtonpost.com/technology/2022/06/11/google-ai-lamda-blake-lemoine/>. Accessed: 2023-12-11.
- Törnberg, P. 2023. Chatgpt-4 outperforms experts and crowd workers in annotating political twitter messages with zero-shot learning. *arXiv preprint arXiv:2304.06588*.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Varol, O.; Ferrara, E.; Davis, C.; Menczer, F.; and Flammini, A. 2017. Online human-bot interactions: Detection, estimation, and characterization. In *Proceedings of the international AAAI conference on web and social media*, volume 11, 280–289.
- Vosoughi, S.; Roy, D.; and Aral, S. 2018. The spread of true and false news online. *science*, 359(6380): 1146–1151.
- Weng, Z.; and Lin, A. 2022. Public opinion manipulation on social media: Social network analysis of twitter bots during the covid-19 pandemic. *International journal of environmental research and public health*, 19(24): 16376.
- Wike, R.; Silver, L.; Fetterolf, J.; Huang, C.; Austin, S.; Clancy, L.; and Gubbala, S. 2022. Social Media Seen as Mostly Good for Democracy Across Many Nations, But U.S. is a Major Outlier. <https://www.pewresearch.org/global/2022/12/06/social-media-seen-as-mostly-good-for-democracy-across-many-nations-but-u-s-is-a-major-outlier/>. Accessed: 2023-12-11.
- Zhang, M.; Qi, X.; Chen, Z.; and Liu, J. 2022. Social bots' involvement in the covid-19 vaccine discussions on Twitter. *International Journal of Environmental Research and Public Health*, 19(3): 1651.