**RESEARCH**

# The devil is in the details: using machine-learning to scrutinize "state-of-the-art" language models' responses to public inquiries across 3 continents

Stany Nzobonimpa[1] · Jean-François Savard[1]

## Abstract

In this article, we seek to investigate bias and quality of so called "state-of-the art" Artificial Intelligence (AI) driven language models under the lens of geographical regions. We assess OpenAI's ChatGPT model and its responses to a diverse group of testers spanning North America, Europe, and Africa. In an era where AI-driven language models are increasingly influencing decision-making and shaping opinions, it is crucial to assess their advice and guidance, particularly on issues of public interest. This study employs a non-representative large sample of respondents, drawn from various walks of life across the three continents. While no representative sampling techniques were employed, the diverse geographical and cultural backgrounds of the participants provide a rich and varied dataset that is leveraged in a rigorously comparative analysis. Our research focuses on ChatGPT's responses to critical topics in public administration including paying taxes, Indigenous Peoples' rights, Indigenous Peoples' self-determination, and whistleblowing, among others. Using advanced machine-learning techniques, we assess the quality and topical orientation of ChatGPT's advice on these selected topics and conduct regional comparisons. This comparative approach helps flag regional variations in the model's responses and provides evidence ChatGPT's answers vary with geographical locations of users. Furthermore, we integrate conceptual insights from the field of public administration to situate our study in a discipline that has seen increased interest in AI adoption. With the diverse, thoroughly selected group of multiple respondents, this study contributes to understanding the implications of AI-driven language models and their limitations when looked at under regional perspectives. We believe our findings will contribute towards the development of region-fair AI-driven conversational agents and offer valuable insights for policymakers, educators, and users.

## 1 Introduction

This is a Generative Artificial Intelligence (GenAI) era. As AI technologies increasingly become integral part of today's society [1, 2], the scrutiny of AI-driven language models has never been more paramount. We designed a study to investigate bias, quality, and coherence of responses provided by the so-called "state-of-the-art" AI language models, with a focus on OpenAI's ChatGPT [3, 4], to a large group of participants. Our research spans respondents from diverse backgrounds across three continents: North America, Europe, and Africa.

Our motivation for this research is rooted in the growing influence of AI-driven language models on decision-making processes and public opinion. As these models gain prominence in providing advice and shaping perspectives on critical issues, it is crucial for us as public management researchers to evaluate their guidance on matters and topics of public interest. The speed of current GenAI language models' evolution has surfaced a gap between keeping pace with AI innovation and understanding its enormous impact on citizens' daily lives. The research objective of this article is to contribute to bridging that gap by attempting to understand potential regional discrepancies in AI-generated recommendations and advice on topics of public interest. To achieve this objective, we selected such public management contemporary topics as tax compliance, government

✉ Stany Nzobonimpa
Stany.nzobonimpa@enap.ca

Jean-François Savard
Jean-Francois.Savard@enap.ca

1 École Nationale d'Administration Publique, Quebec, Canada

directives, Indigenous Peoples' rights, whistleblowing, electoral participation, racism, and careers in public management [5] and prompted OpenAI's model on related questions from the three geographical zones mentioned above.

Using a large, non-representative, sample of respondents, we leverage machine-learning techniques to analyze the consistency and thematic orientation of ChatGPT's responses. By comparing the model's responses across the three different regions, our research aims to uncover regional variations and to assess the potential societal impact of AI-driven chatbots on decision-making processes. Our evaluation integrates theoretical insights from public administration concepts and literature to provide a nuanced understanding of how AI models navigate complex and sensitive topics. Through this comparative analysis, our study, we believe, provides valuable insights for policymakers and society at large, and contributes to informing developers on the perquisites for reliable, trustworthy, and equitable AI conversational agents.

This study contributes to the ongoing and timely discussion on evaluating GenAI models' performance and how these models might have an impact on the way in which citizens interact with public administrations. By looking at the quality of model recommendations on selected topics and across regions while testing these under the lens of public administration theories related to cop-production, we believe this work is at least twice contributive. On one hand, we test for potential bias that can be understood as model's regional preference. This surfaces an important debate on whether users can truly trust GenAI models regardless of where they query them from. On the other hand, we bring a novel point of view about co-production and how this cherished concept of public management might be heading toward a renewed understanding in the age of AI.

This article is organized as follows: we start by exploring relevant literature related to AI and GenAI as seen from a public administration point of view. We explore studies that cover the societal impact of AI, socio-economic biases, and (dis)information dissemination, among others. Our research is grounded in a framework that brings the public administration concept of co-production around the issue of generating information using AI. Our methodology section highlights how the study was designed and carried out with respondents from three continents. In the results and discussion sections, we show how variations across the selected GenAI model's responses shed light on its potential biases when the regional factor is accounted for. The study also presents some limitations and explores avenues for future research.

## 2 Literature review

Adoption of the recently deployed AI-powered language models such as OpenAI's ChatGPT has expanded significantly in various fields and for various uses. Their expansion, however, while generally appreciated and sometimes touted as game-changing tools and as the invention of the century [6], are not immune to well documented and inherent complexities and issues arising from AI [7]. Bias in AI systems has generally been a major concern, and recent studies have demonstrated that even the new language models could reflect and perpetuate existing societal biases present in their training data [8]. Some researchers have termed bias in GenAI as the new versions of the "old demons" of known and documented biases of "traditional" Machine Learning systems [9]. Notably, the quality and consistency of responses from these models play a crucial role in their effectiveness and reliability. For example, a study by Gehman et al. [10] on neural language models demonstrated how issues related to quality and consistency could impact the usability of AI systems in real-world applications.

The potential societal impact of AI-driven language models as we demonstrate in this study, is a growing area of concern and has attracted the attention of researchers in fields as diverse as public administration, computer science, health sciences and many more. While the generative AI domain is still a growing and burgeoning research field, recent studies have shown several potential societal impacts are likely to rise from the new technologies. Gupta et al. [11] explored the theoretical foundations and implications of this field. The authors proposed a comprehensive research agenda to understand the adoption and impacts of GenAI in various domains. They emphasized the potential transformative effects of AI on industries and societies and highlighted both opportunities and challenges. According to the authors, for example, "when ChatGPT is integrated with other technologies such as augmented reality and virtual reality, it would be feasible to develop even more immersive marketing experiences that blend the real world with the virtual world".

Capraro et al. [12] investigated how generative AI influences socioeconomic inequalities and the implications for policy making. Their study addresses both the exacerbation and improvement of existing inequalities by generative AI and provide insights into the nuanced effects of this technology on different socioeconomic strata. The authors discuss potential policy interventions to mitigate negative impacts. For example, studying the potential impact of GenAI on employment, they suggest that "AI will likely have outsized impacts on US workers with Bachelors' or Associates' degrees, compared to higher or lower levels of education"

and conclude that the advent of AI is likely to make the future strikingly different from what was expected:

> We have focused on the socioeconomic inequalities that are likely to be impacted-for better or worse-by the advent of generative AI. This technology has profound implications in the domain of information, where it has the potential to offer more tailored, efficient, and democratic ways to process information. Yet it also poses several challenges, including anticompetitive market advantages, data misuse and abuse, and misinformation. These mixed outcomes will certainly affect a very wide range of social organization and decision making [12, p. 14].

In a comprehensive literature review, Baldassarre [13] analyzed the societal impacts of ChatGPT and examined various dimensions of social influence, including ethical considerations, public perception, and the potential for GenAI to alter social interactions and communication patterns. They found that while the model has many potential benefits such as enhancing customer service, automating repetitive tasks and augmenting education, serious negative impacts also exist. The authors document the impacts related to bias, false information, negative effects on democratic processes, the increase of hate speech, privacy and job loss, among others. They also pointed out some emerging trends and concerns including the uncertainty surrounding copyright, possible inefficiency of regulatory frameworks, increased political biases and environmental impacts. The researchers also note that the concerns are likely to stay as they reflect the society and may not be easily removed:

> In light of this, it is essential to note that generative AI reflects the social context in which it was created. As these models are trained on data that captures various aspects of our reality, it becomes clear that addressing their flaws and biases requires a comprehensive understanding of the broader social context within which they operate. Rather than solely focusing on repairing the model, it is imperative to also engage in a critical examination of the social factors that contribute to these biases and limitations [13, p. 370]

Existing literature has also pined out the ability for AI systems to influence decision-making and shape public opinion. This was, indeed, one of the motivations to carry out this study.

AI systems are traditionally known for their ability to leverage vast amounts of data to provide insights and recommendations that can guide human decisions. For instance, in the healthcare sector, AI algorithms have been touted for their prowess in analyzing patient data to suggest optimal treatment plans, potentially improving patient outcomes [14]. Similarly, in finance, AI-driven models can predict market trends, and assist investors in making informed decisions. These capabilities are rooted in AI's ability to process and interpret complex datasets far more efficiently than humans. With the advent of GenAI, this power has only been increased.

Generative AI also has the potential to profoundly impact and shape public opinions through its applications in media and communications. AI-generated content is increasingly being used in social media to engage audiences and influence their perceptions [15, 16]. The capability of generative AI to produce realistic and compelling content raises ethical concerns about misinformation and manipulation, as users may find it difficult to distinguish between genuine and AI-generated material. Our assessment of ChatGPT's information dissemination to users across three continents contributes to this broader discourse on the role of AI in influencing and shaping public opinion. In this article, we highlight both its potential and challenges when looked at under the lenses of the discipline of public administration. While our research is innovative by comparatively testing ChatGPT' responses across several geographical regions, previous studies have tackled both challenges and opportunities of GenAI in public administration. Related past studies have focused on themes ranging from briefing notes analysis and policy analysis [17, 18], evidence-based policy [17, 19] and government operation [20].

## 3 Conceptual framework: generative AI for co-production

In public administration, co-production is a collaborative process where government services and policies are designed and delivered with active citizen participation [21, 22]. The concept of co-production emphasizes collaborative efforts between public administrators and citizens in service delivery and policymaking. It involves sharing power, resources, and responsibilities to create a more effective service delivery system and democratic governance structures [23–25]. Osborne and Strokosch [21] introduce an important aspect of co-production called "consumer co-production". The authors argue that in service theory, the role of the consumer includes the contribution to the production process and therefore argue that consumption itself is a cornerstone of co-production. Core elements of co-production include shared knowledge, active participation, and a focus on outcomes that meet the needs of the community. Co-production aims at improving service quality, enhance public trust, and foster a sense of community ownership over public policies

and services by engaging citizens directly. It is worth noting that while authors such as Brandsen and Honingh [26] distinguish between co-production and co-creation, our use of the concept does not make that distinction which, as the authors suggest, involves citizens' involvement not only in service co-production but also in strategic planning.

In this article, we examine co-production from a technological perspective and advance the claim that digital technologies can act within, and sometimes as, agents of the co-production process. Prior works show technology can both enable and constrain co-production. Lember [27] for example, studied the impact of digital technologies on co-production and co-creation and proposed a framework to analyze those impacts with various types of technologies. The authors suggest that technology may be a double-edged sword for co-production. On one hand, it may strengthen participation of citizens by making service channels and communication more accessible. On the other hand, technology may create less collaboration with citizens by bypassing the human interaction. Despite these tensions, however, technology-supported co-production can indeed enhance transparency, accountability, and responsiveness in governance and foster trust and engagement between public authorities and the community.

By bringing co-production as the theoretical foundation of our research, we focus on OpenAI's ChatGPT as an AI-mediated participant in co-production. Our empirical setting involves citizens querying a general-purpose chatbot about public matters such as tax obligations. We define two roles GenAI can play:

i.   *AI as channel* (tool): the chatbot primarily routes citizens to existing, authoritative government information.
ii.  *AI as agentic coproducer*: the chatbot performs value-adding micro-tasks, e.g., extracting rules, translating jargon, tailoring guidance to local context, or structuring user inputs-that shape the citizen's understanding and subsequent action.

Yuan [28] classified technology-enabled co-production into three categories namely citizen-sourcing, automatic co-production and government as an open platform. Following this author's suggestions, we treat the two roles outlined above as forms of automatic co-production when (a) the AI automates parts of the information exchange that would otherwise occur with public servants or government platforms and (b) its outputs are measurably aligned with jurisdiction-specific rules and standards.

In our case, the chatbot (OpenAI's ChatGPT) is not embedded in a public service workflow and does not draw directly on government datasets. Accordingly, we conceptualize our use case as AI-mediated informational co-production at the boundary: the citizen and AI jointly assemble task-relevant understanding that conditions later interactions with government (e.g., compliance, applications, service requests). This goes beyond technology-enabled self-service information seeking when the AI transforms and structures information in ways that are specific to the citizen's context and plausibly substitutive of first-line government advice. Consistent with Yuan's categorization of technology-enabled co-production, citizen-sourcing, automatic co-production, and government-as-platform [28], we therefore situate GenAI primarily within automatic co-production, with potential spillovers toward government-as-platform when governments publish machine-readable standards that GenAI models can reliably consume.

Nzobonimpa [14] proposed co-production as one of the ways "AI-enabled public services can be adapted to minimize bias by tailoring innovation to contextual particularities". The notion of citizen involvement is indeed key in co-production postulates and the broad usage and popularity of GenAI in general and of ChatGPT in particular [3, 17] means it can no longer be neglected as a source of information. Therefore, the question that arises is whether the models can be trusted to minimize bias and to adapt to different contextual realities, which, as shows Nzobonimpa [14], is key to ensuring co-production is translated into truly valuable outputs to citizens.

The advent of GenAI tools like ChatGPT has shown the tools have the potential to play a significant role in the co-production of public administration information. They can provide accessible information by processing and disseminating vast amounts of data quickly and have the capabilities to explore, scrape the internet and offer what can be regarded as data-informed insights. While academic evidence on the issue is still scarce, it has emerged that generative AI models can indeed be leveraged for value co-creation as shown by Demir and Demir [29] who studied value co-creation of ChatGPT in the travel industry and found that the model "significantly influences service individualization and service value co-creation, with information internalization and information value playing mediator roles. It is evident that AI technologies are inevitable and should be embraced as essential digital stakeholders for businesses".

Within the framework of the co-production postulates, we verify two hypotheses. First, jurisdiction-aware answers that cite authoritative sources and disclose uncertainty (localization + transparency) will increase citizen task success and perceived legitimacy while narrowing accuracy gaps across regions compared with generic, uncited responses. Second, when the chatbot performs substitutive advisory functions, classifying the case, localizing rules, and producing structured next steps, citizens' reliance on first-line human support will decline without reducing compliance,

and implementing assurance thresholds will further lower harmful error rates. Overall, theoretically, it is argued in this article that if GenAI tools like ChatGPT can meet the standards of accuracy, transparency, and fairness, they have the potential to be trusted participants in the co-production of public administration information to citizens.

## 4 Research design, data collection and methodology

Our research context stems from the growing evidence that GenAI models are being leveraged in various fields as a source of information on critical issues, including on those of public policy lens. As such, we postulate that such tools are increasingly influencing decision-making and shaping public opinions. It is therefore crucial to evaluate their performance and how they compare across contexts, languages, and topics. Our research question is as follows:

> RQ: To what extent do factors of regional and language variance impact the quality of public management inquiry responses of GenAI models and can those models be considered a trustworthy coproduction agent?

To answer the question, we designed a web data collection interface and disseminated it to respondents across Africa, North America, and Europe. Being Canadian researchers and as our institution uses French as the primary language of instruction, we designed our study in Canada's two official languages, French and English. We therefore report results for both languages and cross-compare findings. Recruitment eligibility was solely based on participants' self-reported geographical location. We did not use other methods to operationalize region tagging beyond this self-reported localization where participants had the three options to choose from (Africa, Europe and North America) when entering the data collection platform.

Existing studies show that when reviews are conducted over the Internet, users can profit from the inherent web anonymity to manipulate and falsify information such as their identities in order to falsify or manipulate reviews. In their work on "chasing spammers", Sáez-Ortuño et al. [30] show that techniques such as Internet Protocol (IP) address tracking can help mitigate the issue of location manipulation. While we recognize the robustness of this method, we did not perform such tracking techniques because of the voluntary nature of our study. In particular, in the above-mentioned work, the authors suggest four types of motivation (revenge, entertainment, profit and self-esteem), none of which was deemed to be a potentially major factor that could affect the quality of our experimental study which

involved no gain whatsoever for participants. Our research design also excluded sharing of any personal information including email addresses of respondents. We specified that participants willing to contact the researchers could do so directly using our address in a separate conversation. This ensured that the three spam actors pointed out by Stringhini et al. [31] as harvester, botmasters and spammers had no interest in participating in our study. Nonetheless, we ensured that our messaging around location, the voluntary and non-monetary nature of our study was clear when recruiting participants via the above-mentioned channels. When data collection was concluded, we received 14 emails where participants followed up individually asking to be sent a copy of our publication when the research is over.

Participants were recruited through social media sharing, institution communication, word-of-mouth and snowball sharing. All our participants were required to have a free account of ChatGPT and to use the GPT-3.5 version, which was the model freely available to all account holders at the time of designing our research. OpenAI released the GPT-3.5 model in November 2022. In March 2023, the company announced the release of its then newer version, GPT-4 which was only available to paid users with a ChatGPT Plus subscription. When recruiting our participants, we specified that multiple submissions were not allowed and therefore, no deduplication instances occurred. The design was exploratory in nature and, as such, no other models were included, although models like Google's Gemini, Anthropic's Claude and Mistral's Le Chat were also available at the time. Participants were asked to prompt the models using the following questions, which the researchers had tested and designed to be as simple and as straightforward as possible (Table 1):

It is worth noting that since this was a study involving human beings, we followed our institution's guidelines with regards to ethical conduct of research. As such, we ensured clear, active and well understood participant consent, voluntary participation and data protection. When announcing the study and recruiting participants, we ensured the voluntary nature of the study was clearly stated. With regards to consent, the platform we created (called GML Lab) for data collection was configured so that only participants who read the consent statement and selected "I consent" were able to complete the study. Consent statements included consenting on having to use a free ChatGPT account, that no personal information was to be shared and that all responses were to remain anonymous. Additionally, participants maintained the right to withdraw from the research at any given moment and they also had the right not to submit a ChatGPT response for personal reasons. All participants who did not consent were automatically taken to a different page where

**Table 1** Prompt questions into ChatGPT

| English prompts | French prompts |
|---|---|
| Should we pay taxes? | Devons-nous payer des impôts ? |
| How are my taxes being used? | Comment mes impôts sont-ils utilisés ? |
| Should public servants blow the whistle and denounce wrongdoing? | Les fonctionnaires devraient-ils dénoncer les actes répréhensibles ? |
| What are the most effective mechanisms for protecting whistleblowers from retaliation? | Quels sont les mécanismes les plus efficaces pour protéger les lanceurs d'alerte contre les représailles ? |
| Should public servants comply with government directives or leverage more their discretionary power? | Les fonctionnaires devraient-ils se conformer aux directives gouvernementales ou exploiter davantage leur pouvoir discrétionnaire ? |
| Should we trust the government? | Faut-il faire confiance au gouvernement ? |
| Should we participate in politics? | Devons-nous participer à la politique ? |
| What should motivate an individual to participate in politics? | Qu'est-ce qui devrait motiver un individu à participer à la politique ? |
| Who should pursue public administration careers? | Qui devrait poursuivre une carrière dans l'administration publique ? |
| How can public administration careers contribute to addressing pressing societal challenges? | Comment les carrières dans l'administration publique peuvent-elles contribuer à relever des défis sociétaux urgents ? |
| Should Indigenous people pay taxes? | Est-ce que les peuples autochtones devraient payer des taxes ? |
| Do Indigenous people have the right to self-determination? | Est-ce que les peuples autochtones ont le droit à l'autodétermination ? |
| Are there Indigenous governments? | Est-ce qu'il existe des gouvernements autochtones |
| Should governments keep transferring money to Indigenous communities? | Est-ce les gouvernements devraient continuer à transférer des fonds aux communautés autochtones ? |
| Is there racism against Indigenous people? | Est-ce qu'il y du racisme contre les peuples autochtones ? |

**Table 2** Total responses by region

| Region | Month–year | Total responses |
|---|---|---|
| Africa | Feb-24 | 292 |
| Africa | Mar-24 | 456 |
| Africa | Apr-24 | 218 |
| North America | Feb-24 | 127 |
| North America | Mar-24 | 45 |
| Europe | Feb-24 | 45 |
| Europe | Mar-24 | 15 |

they were encouraged to share the study to others who might be interested in participating.

Our analytical approach leverages Natural Language Processing and other machine learning techniques. As per our research question, we seek to analyse regional and language variances in the quality of returned responses to the prompts fed into the model. We propose a novel approach to topic extraction leveraging a machine learning model known as LatentDirichletAllocation (LDA) that we combine with Boolean search techniques to extract key dual topics representative of responses returned by the model across respondents. We also employ similarity and frequency analysis techniques as detailed in the following section on results. It is worth noting that to ensure replicability, we have shared all codes and input data in open-source repositories. These are available on the first author's GitHub page.

# 5 Results

In total, 1198 responses were received to our ChatGPT prompts across all questions in both French and English. Table 2 shows the breakdown of prompt responses by region and collection month-year.

As shown on the above table, the region "Africa" predominantly participated in our experiment compared to the other two regions. This dominance of African respondents can partly be explained by our reliance on social media and personal contacts in disseminating the experiment. The research was shared with African universities whose participants tended to respond more to the study. However, despite the predominantly African response rate, our results show ChatGPT responses to African respondents tended to be significantly shorter on average i.e., per prompt response than those from North America and European respondents. On average, a typical ChatGPT response to prompts from Africa-based respondents was 678 characters compared to over 1000 for North American respondents and over 2000 for European respondents, as shown in Fig. 1 below. This result held true regardless of the language of the prompt. It is worth noting that despite these evident variations, we do not make causal inferences from our observation mainly because we acknowledge that uncontrollable factors may have been in play during the data collection stage.
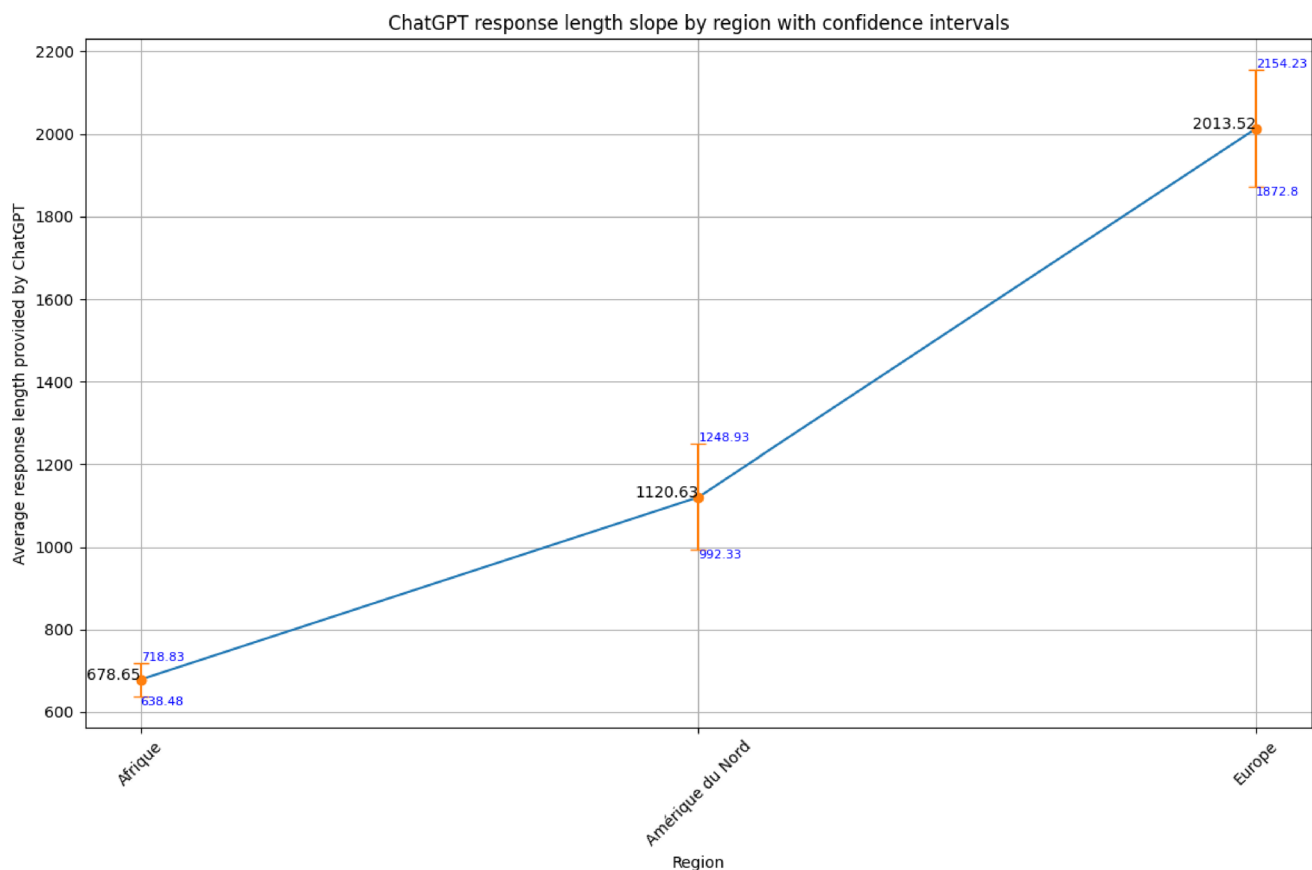
Fig. 1 ChatGPT response length slope by region with confidence intervals

**Table 3** Description statistics for response length across regions

| Region | Mean length | Median length | Std Dev | IQR (Q3-Q1) | Count |
|---|---|---|---|---|---|
| Africa | 678.65 | 396 | 637.09 | 412.5 | 966 |
| North America | 1120.63 | 884.5 | 858.49 | 1530.5 | 172 |
| Europe | 2013.52 | 2112 | 556.12 | 638.75 | 60 |

For example, for users who already had accounts and who may have activated memory functions, their responses may reflect their current use of ChatGPT and not be merely based on their geographical location. Nevertheless, since our platform automatically prompted users to start a new prompt in a new session, we believe that past conversations, history and memory had a minor impact on the model's response. Additionally, at the time of collecting data, OpenAI had not made model options available to users except for developers using the OpenAI API platform which was outside the scope of our study. This ensured that all our users queried the same GPT version.

As shown in Fig. 1, responses from North America and Europe were significantly longer than those from Africa and were, in return, more detailed. Since all users had the same questions and there was no change in wording, there are no other major factors that could explain the regional differences, expect for the region itself. We also looked at median and interquartile ranges to confirm this result. In the table below, we provide detailed results on response length across the region by quartiles, mean and medians.

In Table 3 , the regional analysis of response length shows distinct patterns across the three regions. Europe provides the longest responses by a significant margin, as shown by its highest mean (2,013.52) and median (2,112.0) character counts. The fact that the median is slightly higher than the mean shows the response length distribution in Europe is relatively symmetric and slightly negatively skewed which indicates a lack of extremely short responses for the region. This was the case for both languages. In comparison, North America showed the second-longest response lengths, with a mean of 1,120.63 and a median of 884.5. This region is characterized by a large Interquartile Range (IQR) of 1,530.50 which suggest most responses are clustered around the median length but there are still some long responses. In Table 3, Africa yields the shortest responses. Its mean (678.65) and median (396.0) are less than half those found in North America.

We also conducted a Kruskal–Wallis H-Test to verify robustness of the group difference in Table 3 since the data

may not be normally distributed as indicated by the large difference between mean and median for some groups. The test yielded a highly significant result ($H = 149.2$, $p < 0.0001$) which confirms our previous observations. Thus, we reject the null hypothesis that the median response lengths are the same across all three regions. This results is further clarified with the distribution plot in Fig. 2 below. To carry out our analyses, we began by checking word frequency. As we explained in the methodology section, we used relative frequency and inverse frequency analysis to avoid over-representing meaningless words (such as 'the', 'a', 'in', 'on', etc.) also known as stop words, to highlight the most significant words in our corpus. As Fig. 3 shows, the most significant words to emerge from our corpus relate to Indigenous Peoples' issues, public trust issues and tax-related themes. This is not surprising since it corresponds to the nature of our questions, but the interest for this first analysis lies in the fact that it confirms ChatGPT indeed answered our questions directly (Fig. 4).

Our thematic extraction analysis, a significant part of our research, further underscores the difference in ChatGPT responses across regions. The LDA analysis, based on two crucial keywords-government and communities, was able to extract the themes associated with these words. Techniques of topic filtering and extraction are especially useful in complex documents where texts may lack clarity or be diffuse [32–36]. Figures 3 and 4 clearly depict a major thematic difference between the African and North American regions, as indicated by the relative and inverted frequency analysis. This depth of analysis provides a comprehensive understanding of ChatGPT's regional variations. It is worth noting that when using LDA, we separated French and English to ease stop word exclusion. To achieve this, we defined a Python custom keywords variables specified using the bult-in ENGLISH_STOP_WORDS function. In our code, this was specified as custom_stopwords = list(text.ENGLISH_STOP_WORDS.union({'public', 'yes','including','ensure'})) for English. For French, we manually created a list of such stop words which contained words such as "le", "un", "je", "que", etc.[1]

The relative frequency and inverted frequency analysis were especially useful in building the list of words we needed to run our vectorization analysis. The idea was to determine the extent to which responses to the questions were similar from one region to another. Figures 4 and 5 show the positioning of each word across the responses for Europe, North America and Africa, according to their frequency. We were unable to take Europe into account for
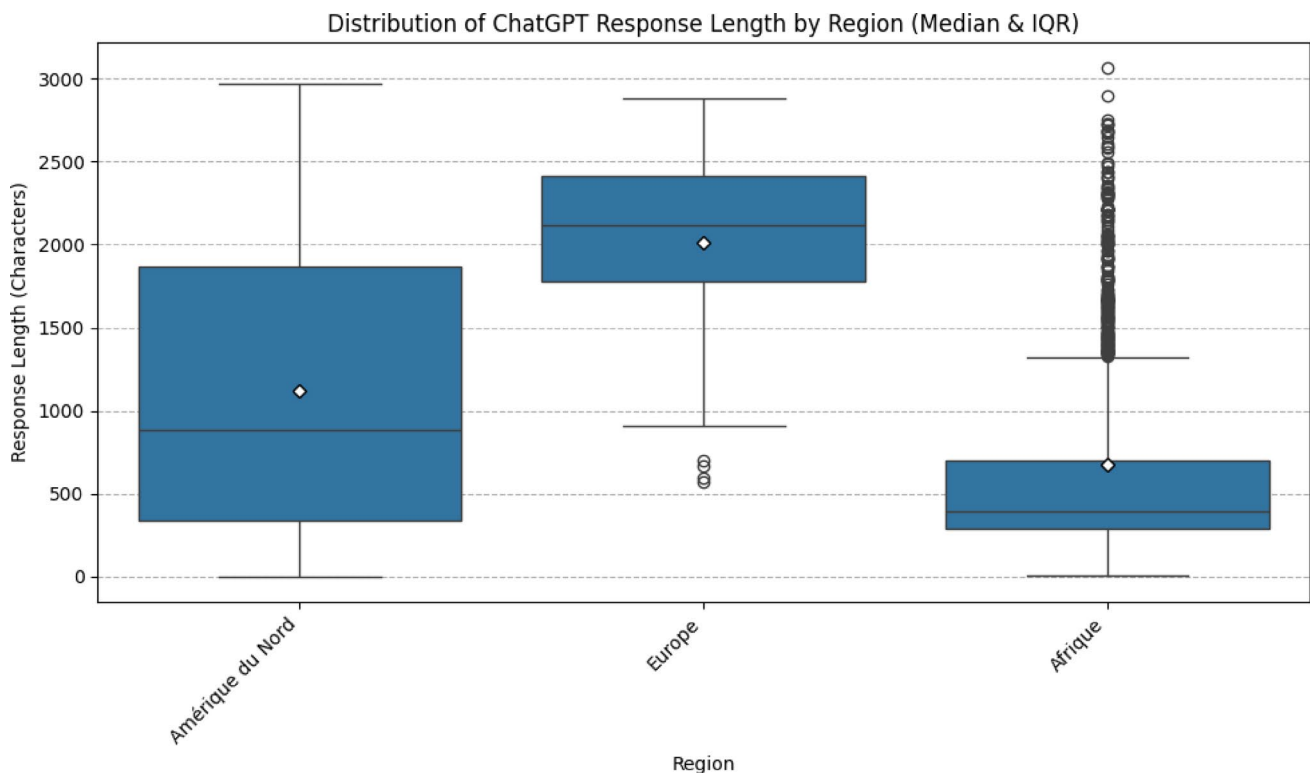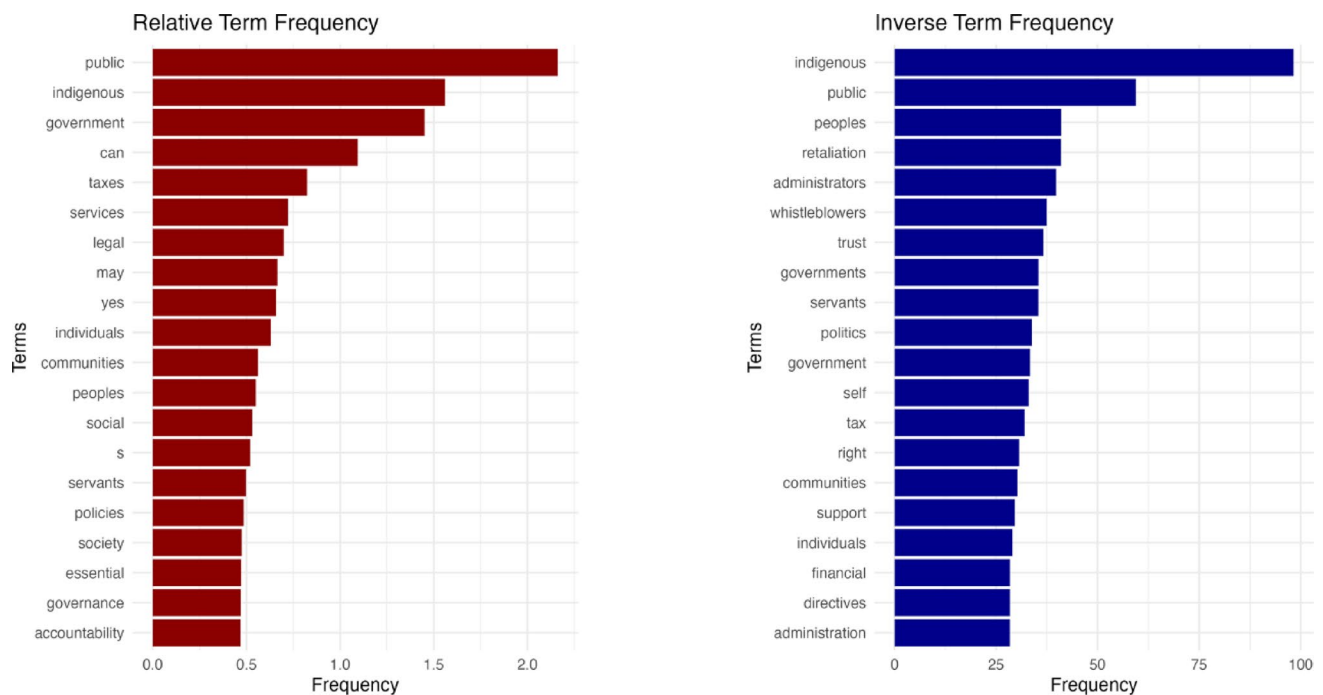


**Fig. 2** Distribution of responses

**Fig. 3** ChatGPT Responses' relative terms frequency and inverse term frequency
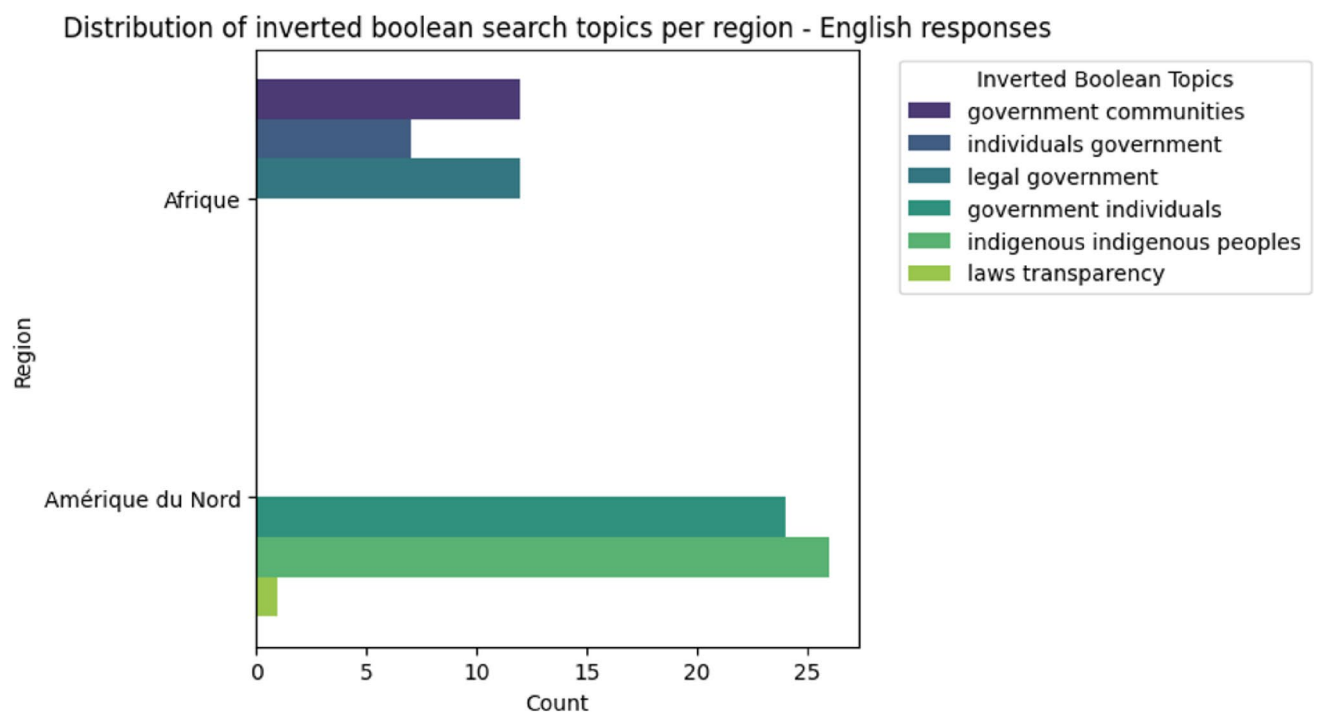


**Fig. 4** Distribution of inverted Boolean search topics per region—English responses

English at Fig. 4 as the number of responses was insufficient. As we can see from these graphs, the positioning of words between Europe, North America and Africa is different, which is an indicator that responses between these regions are very different. This led us to conclude that ChatGPT adapted its answers according to the region from which the question was asked.

To validate these initial results, we conducted a similarity analysis taking into account responses from Europe. To do this, we used a technique that calculates a similarity score for answers to the same question and then produces a

## Distribution of inverted boolean search topics per region - French responses
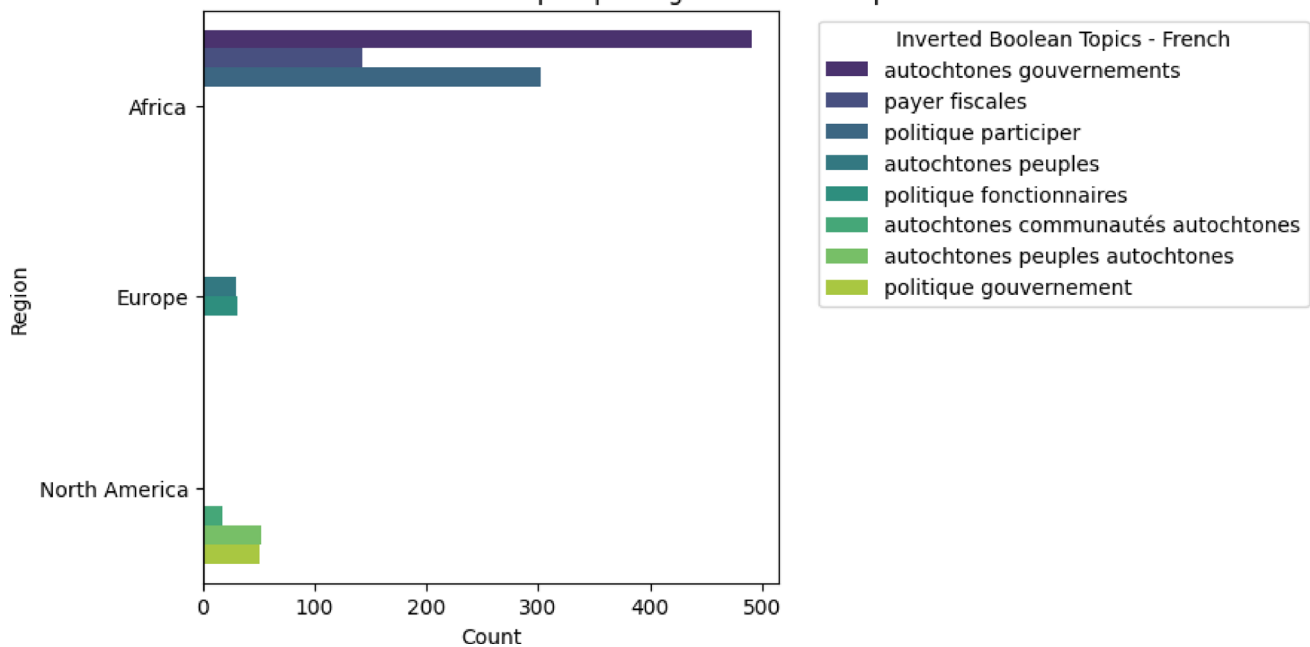


**Fig. 5** Distribution of inverted Boolean search topics per region for French responses

scatterplot. The higher the points on the graph, the higher the similarity of answers to a question, and conversely, the lower the points on the graph, the lower the similarity. Furthermore, the more concentrated the cloud, the greater the similarity between regions; the more dispersed the cloud, the lower the similarity between regions.

First, we notice that the similarity of responses is lower for the first questions on trust, taxes and the actions of public officials. We can, therefore, conclude that ChatGPT responses varied a great deal from region to region in terms of the values associated with public service, which would confirm a regional anchoring of ChatGPT responses since we can assume that the values associated with public service vary from one society to another. As the last questions are linked to Indigenous Peoples' realities, and as these issues are much more prevalent and present in North America than in Europe or Africa, it may well be that ChatGPT did not have enough information from the latter two regions to formulate answers more rooted in their realities, and produced answers based on data specific to North America (Figs. 6 and 7).

## 6 Discussion

To our question about whether factors of variance impact the quality of responses when using OpenAI's ChatGPT, our data show geographical location where prompts (or questions) originate impacts the responses produced by the model both in quantity and content. We observed that

responses in Africa were considerably shorter than in America or Europe and key themes differ accordingly. We also observed that there was a fair degree of dissimilarity between the responses produced in each region, and the nature of the themes in each region was also different.

Under our theoretical framework, the implications of this study are such that the use of GenAI technologies could be considered a potential agent of "automatic co-production". Given this type of technology is able to produce, as shown in our results, answers that could make it a complement to a civil servant or a government platform for citizens' queries, co-production in the age of GenAI should be regarded as a changing concept that is likely to put on a different hat in the near future. This could also imply significant time, money, and human resources savings for both the government and its citizens. For example, by having a GenAI that answers citizen's queries accurately on topics such as those we used in our experiment, governments can make significant savings and invest on model quality and accuracy that would benefit citizens while making their interaction with the public sector easy and handy. From this perspective, the use of GenAI can be regarded as a way of improving public services. However, despite this promise, it is worth noting that questions remain regarding the sense of community ownership over public policies, services and public trust. While governments can control the sources of information provided to citizens through such AI agents, it remains challenging to determine how that information is processed and rendered by the technology behind GenAI, because, as we

**Fig. 6** Token distribution between each answer for North America and Africa
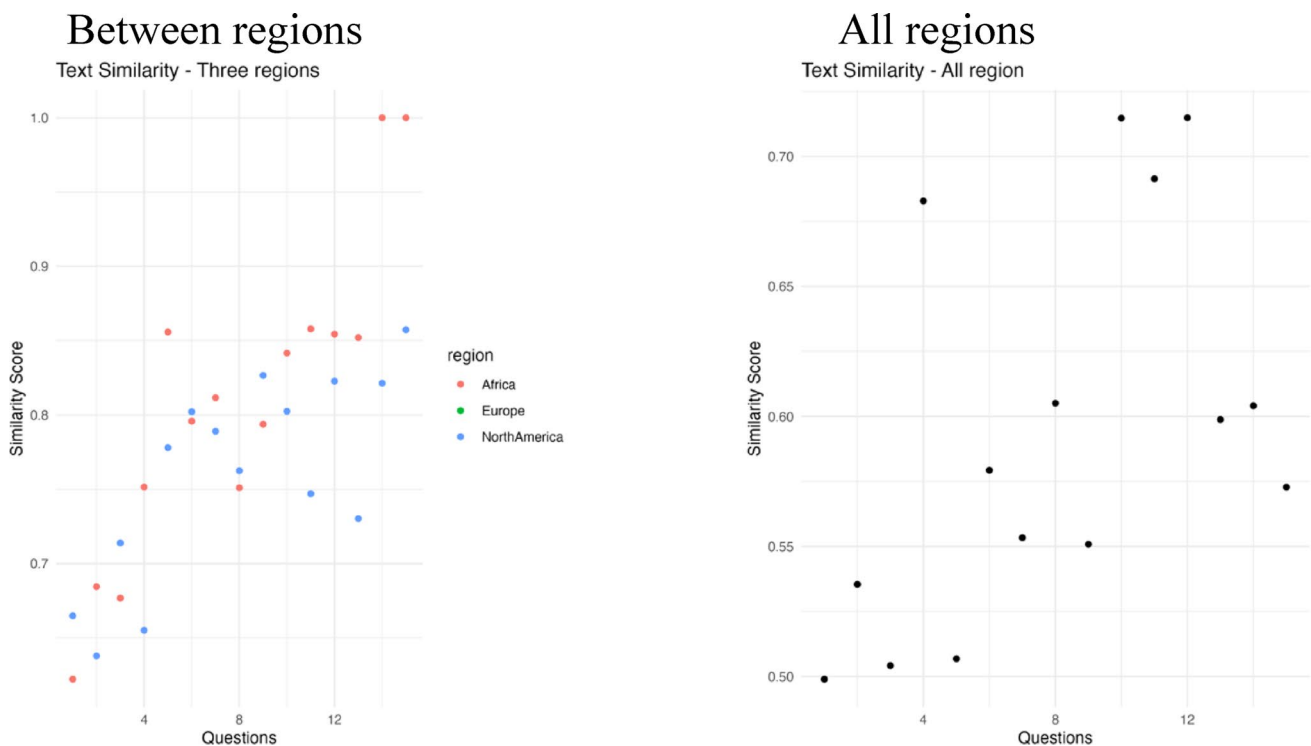


**Fig. 7** Similarity scores for each question—between and all regions

have shown, factors such as regions could play key roles in the nature of responses returned by the models.

When it comes to trust, we believe that the issue remains complex. We have pointed out that GenAI can be a two-edged sword in co-production. We have also shown that

there are potentially positive effects of this technology that are likely to encourage public trust, mainly if, in the form of an echo chamber phenomenon, the responses obtained are consistent with the values, experience and knowledge of the public using them. However, as pointed out by many

authors such as Noble [37], Collins [38] and Esposito [39], AI systems are known for the reproduction of social biases. Governments should therefore be mindful of possible short-comings even as they integrate GenAI as a truly agent of co-production. Thus, we believe that the public's use of GenAI to better understand public interest issues or answer questions about policies and programs can make it an effective tool for co-production. But the danger would be in letting the governments' guardrails down. As Collins [38] so aptly points out, despite these benefits, we must remain aware of AI's limits and ensure there is no "humanity's surrender" to technology.

# 7 Conclusion

Our research sought to determine the extent to which responses to questions prompted to OpenAI's ChatGPT about public interest topics varied depending on regional location of inquirers. To this end, we designed a web-based data collection application that allowed to query the model from three regions. The tool allowed participants we recruited in North America, Europe and Africa to prompt ChatGPT using fifteen questions and copy the answers into the application. These answers were then added to a database that we later used to conduct the analyses discussed in this article. In total, we gathered 1198 answers from participants across the three regions. These answerers were unevenly split between the regions with an over-representation of Africa (966) compared to North America (172) and Europe (60). This over-representation of Africa can be explained by a more effective recruitment campaign in this region than in others.

At our analytical phase, we carried out several types of textual analyses to compare answers across regions. The techniques we used include relative and inverse frequency analysis, thematic analysis and similarity analysis. Our results show that ChatGPT responses to our inquiries did indeed vary according to the regional location of the participants. Firstly, in Africa, the responses obtained by our participants were much shorter than those offered in North America or Europe. Secondly, vectorization and similarity analysis showed significant differences in response content between these regions. Furthermore, our thematic analysis also indicated a difference in the themes addressed in the responses produced by OpenAI's model. We have therefore concluded from these analyses that there is a real semantic difference in ChatGPT responses when analyzed under these regional lenses.

Based on our preliminary observations, the analyzed model could be seen as a type of automating co-production as defined in our theoretical framework. However, the

issues of regional variations we uncovered mean that these models may still be far from being truly "state-of-the-art" agents of co-production in the public administration sense. Users should be mindful of documented problems, including quality and length of responses as these may hide more problematic issues. Additionally, as we have shown in our literature review section, users still need to ensure that the responses offered by these tools do not produce hallucinations or reproduce other known social biases. In the words of Collins [38] it is crucial that both citizens and governments do not heavily rely or "surrender" to these models as reliable helpers of co-production.

## 7.1 Limitations of the study

This study's results and conclusions should be understood within the framework of certain limitations mainly related to its experimental nature. For example, as mentioned previously, our study was limited to using OpenAI's freely available model to avoid incurring costs to our participants and to ensure comparability of results. At the time of our data collection stage, the freely available model was GPT-3.5. Thus, it is possible that paid models such as the GPT-4 or GPT-5 families could have yielded different outputs. However, we believe that using freely available models reflected what we expect most users would have access to and use, making this limitation less significant. Similarly, we did not attempt to design the study in a way that would have allowed model comparison as we aimed at studying regional, rather than model, variations. This explains why other models such as Google's Gemini, Anthropic's Claude and Mistral's Le Chat were excluded from our study. Another limitation relates to the data collection tool that was employed. It came to our attention, during the analytical stage, that the data collection questionnaire might have been too long as some participants did not answer all the questions. Because of this limitation, we avoided making conclusive remarks on such questions. We also face the challenge of working with two different languages (French and English). To overcome this limitation, we decided to process the input data as well as the analyses in parallel which allowed us to conduct similar analysis for each language and provided us an opportunity for comparison.

## Declarations

**Conflict of interest**  The authors declare no competing interests.

## References

1. Thomas, H.D., Rajeev, R.: Artificial Intelligence for the Real World, https://www.hbsp.harvard.edu/product/R1801H-PDF-ENG

2. Topol, E.J.: High-performance medicine: the convergence of human and artificial intelligence. Nat. Med. **25**, 44–56 (2019). https://doi.org/10.1038/s41591-018-0300-7

3. Aydin, Ö., Karaarslan, E.: Is ChatGPT leading generative AI? What is beyond expectations? Acad. Platf. J. Eng. Smart Syst. **11**, 118–134 (2023). https://doi.org/10.21541/apjess.1293702

4. Rane, N., Choudhary, S., Rane, J.: Gemini versus ChatGPT: applications, performance, architecture, capabilities, and implementation. SSRN J. (2024). https://doi.org/10.2139/ssrn.4723687

5. Côté, L., Savard, J.-F.: Dictionnaire encyclopédique de l'administration publique: la référence pour comprendre l'action publique. École nationale d'adminstration publique, Québec (2012)

6. Leslie, D.: Does the sun rise for ChatGPT? Scientific discovery in the age of generative AI. AI Ethics (2023). https://doi.org/10.1007/s43681-023-00315-3

7. Kapoor, S., Henderson, P., Narayanan, A.: Promises and pitfalls of artificial intelligence for legal applications, https://arxiv.org/abs/2402.01656 (2024)

8. Bender, E.M., Gebru, T., McMillan-Major, A., Shmitchell, S.: On the dangers of stochastic parrots: can language models be too big?. In: Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency. pp. 610–623. ACM, Virtual Event Canada (2021)

9. Nzobonimpa, S.: new tech, old demons: generative artificial intelligence could inherit the flaws of machine learning. SSRN (2025). https://doi.org/10.2139/ssrn.5118289

10. Gehman, S., Gururangan, S., Sap, M., Choi, Y., Smith, N.A.: RealToxicityPrompts: Evaluating Neural Toxic Degeneration in Language Models, https://arxiv.org/abs/2009.11462 (2020)

11. Gupta, R., Nair, K., Mishra, M., Ibrahim, B., Bhardwaj, S.: Adoption and impacts of generative artificial intelligence: theoretical underpinnings and research agenda. Int. J. Inf. Manag. Data Insights **4**, 100232 (2024). https://doi.org/10.1016/j.jjimei.2024.100232

12. Capraro, V., Lentsch, A., Acemoglu, D., Akgun, S., Akhmedova, A., Bilancini, E., Bonnefon, J.-F., Brañas-Garza, P., Butera, L., Douglas, K.M., Everett, J.A.C., Gigerenzer, G., Greenhow, C., Hashimoto, D.A., Holt-Lunstad, J., Jetten, J., Johnson, S., Kunz, W.H., Longoni, C., Lunn, P., Natale, S., Paluch, S., Rahwan, I., Selwyn, N., Singh, V., Suri, S., Sutcliffe, J., Tomlinson, J., Van Der Linden, S., Van Lange, P.A.M., Wall, F., Van Bavel, J.J., Viale, R.: The impact of generative artificial intelligence on socioeconomic inequalities and policy making. PNAS Nexus **3**, pgae191 (2024). https://doi.org/10.1093/pnasnexus/pgae191

13. Baldassarre, M.T., Caivano, D., Nieto, B.F., Gigante, D., Ragone, A.: The Social Impact of Generative AI: An Analysis on ChatGPT. (2024). https://doi.org/10.48550/ARXIV.2403.04667

14. Nzobonimpa, S.: Artificial intelligence, task complexity and uncertainty: analyzing the advantages and disadvantages of using algorithms in public service delivery under public administration theories. Digit. Transform. Soc. **2**, 219–234 (2023). https://doi.org/10.1108/DTS-03-2023-0018

15. Dunn, A.G., Shih, I., Ayre, J., Spallek, H.: What generative AI means for trust in health communications. J. Commun. Healthc. **16**, 385–388 (2023). https://doi.org/10.1080/17538068.2023.2277489

16. Pavlik, J.V.: Collaborating with ChatGPT: considering the implications of generative artificial intelligence for journalism and media education. J. Mass. Commun. Educ. **78**, 84–93 (2023). https://doi.org/10.1177/10776958221149577

17. Nzobonimpa, S., Savard, J.-F., Caron, I., Lawarée, J.: Automating public policy: a comparative study of conversational artificial intelligence models and human expertise in crafting briefing notes. AI Soc. **40**, 3627–3639 (2025). https://doi.org/10.1007/s00146-024-02103-x

18. Safaei, M., Longo, J.: The end of the policy analyst? Testing the capability of artificial intelligence to generate plausible, persuasive, and useful policy analysis. Digit. Gov.: Res. Pract. **5**, 1–35 (2024). https://doi.org/10.1145/3604570

19. Newman, J., Mintrom, M.: Mapping the discourse on evidence-based policy, artificial intelligence, and the ethical practice of policy analysis. J. Eur. Public Policy **30**, 1839–1859 (2023). https://doi.org/10.1080/13501763.2023.2193223

20. Mamalis, M.E., Kalampokis, E., Karamanou, A., Brimos, P., Tarabanis, K.: Can large language models revolutionize open government data portals? A case of using ChatGPT in statistics. gov.scot. In: Proceedings of the 27th Pan-Hellenic Conference on Progress in Computing and Informatics. pp. 53–59. ACM, Lamia Greece (2023)

21. Osborne, S.P., Strokosch, K.: It takes two to tango? Understanding the Co-production of public services by integrating the services management and public administration perspectives. Br. J. Manag. 24 (2013). https://doi.org/10.1111/1467-8551.12010

22. Bandola-Gill, J., Arthur, M., Leng, R.I.: What is co-production? Conceptualising and understanding co-production of knowledge and policy across different theoretical perspectives. Evid. Policy. **19**, 275–298 (2023). https://doi.org/10.1332/174426421X16420955772641

23. Farr, M.: Power dynamics and collaborative mechanisms in co-production and co-design processes. Crit. Soc. Policy **38**, 623–644 (2018). https://doi.org/10.1177/0261018317747444

24. Turnhout, E., Metze, T., Wyborn, C., Klenk, N., Louder, E.: The politics of co-production: participation, power, and transformation. Curr. Opin. Environ. Sustain. **42**, 15–21 (2020). https://doi.org/10.1016/j.cosust.2019.11.009

25. Verschuere, B., Brandsen, T., Pestoff, V.: Co-production: the state of the art in research and the future agenda. Voluntas **23**, 1083–1101 (2012). https://doi.org/10.1007/s11266-012-9307-8

26. Brandsen, T., Honingh, M.: Definitions of co-production and co-creation. In: Co-Production and Co-Creation. Routledge (2018)

27. Lember, V., Brandsen, T., Tõnurist, P.: The potential impacts of digital technologies on co-production and co-creation. Public Manag. Rev. **21**, 1665–1686 (2019). https://doi.org/10.1080/14719037.2019.1619807

28. Yuan, Q.: Co-production of public service and information technology: a literature review. In: Proceedings of the 20th Annual International Conference on Digital Government Research. pp. 123–132. ACM, Dubai United Arab Emirates (2019)

29. Demir, M., Demir, ŞŞ: Is ChatGPT the right technology for service individualization and value co-creation? Evidence from the travel industry. J. Travel Tourism Mark. **40**, 383–398 (2023). https://doi.org/10.1080/10548408.2023.2255884

30. Sáez-Ortuño, L., Forgas-Coll, S., Huertas-Garcia, R., Puertas-Prats, E.: Chasing spammers: using the Internet protocol address for detection. Psychol. Mark. **41**, 1363–1382 (2024). https://doi.org/10.1002/mar.21985

31. Stringhini, G., Hohlfeld, O., Kruegel, C., Vigna, G.: The harvester, the botmaster, and the spammer: on the relations between

the different actors in the spam landscape. In: Proceedings of the 9th ACM symposium on Information, computer and communications security. pp. 353–364. ACM, Kyoto Japan (2014)

32. Damarell, R.A., May, N., Hammond, S., Sladek, R.M., Tieman, J.J.: Topic search filters: a systematic scoping review. Health Info. Libr. J. **36**, 4–40 (2019). https://doi.org/10.1111/hir.12244

33. Kim, Y., Seo, J., Croft, W.B.: Automatic boolean query suggestion for professional search. In: Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval. pp. 825–834. ACM, Beijing China (2011)

34. Lashkari, A.H., Mahdavi, F., Ghomi, V.: A boolean model in information retrieval for search engines. In: 2009 International Conference on Information Management and Engineering. pp. 385–389. IEEE, Kuala Lumpur, Malaysia (2009)

35. Murphy, K.A., Bassett, D.S.: Information decomposition in complex systems via machine learning. Proc. Natl. Acad. Sci. U. S. A. **121**, e2312988121 (2024). https://doi.org/10.1073/pnas.2312988121

36. Wang, S., Scells, H., Koopman, B., Zuccon, G.: Can ChatGPT write a good boolean query for systematic review literature search? In: Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 1426–1436. ACM, Taipei Taiwan (2023)

37. Noble, S.U.: Algorithms of oppression: how search engines reinforce racism. New York University Press, New York (2018)

38. Collins, H.M.: Artifictional intelligence: against humanity's surrender to computers. Polity Press, Medford (2018)

39. Esposito, E.: Artificial communication: how algorithms produce social intelligence. The MIT Press, Cambridge (2022)