# Opportunities and Risks of LLMs for Scalable Deliberation with Polis

Christopher T. Small[*,1], Ivan Vendrov[2], Esin Durmus[2], Hadjar Homaei[1],
Elizabeth Barry[1], Julien Cornebise[1], Ted Suzman[1], Deep Ganguli[2], and Colin
Megill[1]

[1]The Computational Democracy Project
[2]Anthropic

June 2023

## Abstract

Polis is a platform that leverages machine intelligence to scale up deliberative processes. In this paper, we explore the opportunities and risks associated with applying Large Language Models (LLMs) towards challenges with facilitating, moderating and summarizing the results of Polis engagements. In particular, we demonstrate with pilot experiments using Anthropic's Claude that LLMs can indeed augment human intelligence to help more efficiently run Polis conversations. In particular, we find that summarization capabilities enable categorically new methods with immense promise to empower the public in collective meaning-making exercises. And notably, LLM context limitations have a significant impact on insight and quality of these results.

However, these opportunities come with risks. We discuss some of these risks, as well as principles and techniques for characterizing and mitigating them, and the implications for other deliberative or political systems that may employ LLMs. Finally, we conclude with several open future research directions for augmenting tools like Polis with LLMs.

---

[*]Corresponding author: chris@compdemocracy.org

# Contents

# 1 Introduction

## 1.1 Applying Machine Intelligence to Deliberation

Polis is a powerful platform for collective self-representation, enabling the construction of an interactive model of the public opinion landscape (Small et al., 2021). Since initial development began in 2012, it has been used to find common ground in thousands of conversations tied to policy making processes across dozens of countries and five continents, and been reported on widely in outlets like the New York Times (Coy, 2023), BBC (Miller and Anubi), and The Atlantic (Applebaum and Pomerantsev, 2021). It has been adopted nationally by Taiwan, Singapore and Finland, as well as by numerous municipalities. Its methodology was also built on as a basis for Twitter's Community Notes (Wojcik et al., 2022).

The goal of Polis is to enable open and constructive dialogue between people with diverse viewpoints. Using a combination of machine learning and human interaction, Polis helps groups understand each other and identify points of common ground. Polis can be used for a variety of purposes, including group decision-making, public engagement, and public opinion research.

Polis works by allowing participants to submit comments in response to an open ended prompt, and votes in response to other people's comments. Figure 1 shows the participation interface. By arranging the votes in a *vote matrix* and applying dimension reduction and clustering techniques—specifically Principle Components Analysis (PCA) (Pearson, 1901) and K-means clustering (MacQueen, 1967)—Polis is able to learn a 2-dimensional *opinion space* and *opinion groups*, which are used as the basis for syntheses of the deliberations. Specifically, its possible to surface comments which best distinguish opinion groups, as well as points of consensus between groups (*group-informed consensus*). This manifests in a real-time visualization within the participation interface, as well as in an automated report for deeper analysis. Figure 2 shows the full workflow of a Polis conversation.

To date, Polis as a system for large scale deliberation has not leveraged Natural Language Processing (NLP) techniques. Circa 2012, when Polis was developed, NLP techniques were considered not robust enough for capturing nuance in short text documents. Building the platform on comparatively simple and well-trod machine learning methods meant that results were more robust, as well as easier to interpret and explain.

Figure 1: Polis participation interface. The topic of the conversation is at the top. A box displays another participant's statement, along with vote options: Agree, Disagree, and Pass/Unsure. The participant has the option to add a new statement of their own, on which other participants will then vote. Below is a visualization of the interactive opinion landscape as summarized by 2D-PCA and K-means. It shows the two opinion groups that emerged, their relative size, and some example statements. The user can also see some key statements and how representative they are of each group.

Nevertheless, we continued to monitor advances in NLP, and recognized a critical threshold in capacity had been crossed with the advent of Large Language Models (LLMs), e.g. Brown et al. (2020). We now believe that they have significant potential to augment human intelligence in deliberative settings, but recognize that this comes with risks, and find it imperative that systems be designed to account for and mitigate these risks.

## 1.2    Challenges with Polis

Polis has a proven track record of hosting open and constructive dialogue between people with diverse viewpoints (see e.g. Horton, 2018; CPI, 2019). However, the scale of data produced, upwards of millions of votes per conversation, comes with the challenge of synthesizing the results to inform policy makers, researchers, participants, and the public at large.

Polis' machine learning techniques and automated report provide scaffolding for this task, and indeed make it much more feasible than it would otherwise be, but successful applications of the platform nevertheless require dozens of hours of (human) analysis per conversation. This typically manifests in the production of a written report which builds on the automated results to tell a story about the data in terms of the learned opinion groups. This role requires significant expertise in data analysis, writing skills to summarize for a general audience, and contextual knowledge about Polis and the topic of discussion. These background requirements, training required to build up familiarity with Polis, and ultimately time spent analyzing and reporting on the data, pose a significant cost to successful adoption of the platform.

Polis engagements also require a significant commitment of human time and energy devoted to facilitating conversations. To ensure that the initial group of participants have statements to consider as soon as they join the discussion, preliminary comments, known as *seed statements*, are typically provided. These statements help establish the conversation's tone and offer a starting point for participants to delve into the discussion. Participant submitted statements are also typically moderated to ensure they are on-topic, not abusive, and sufficiently distinct from existing statements. More active facilitators will sometimes add clarifying statements based on content they see emerging from the participants themselves. All of these tasks require significant training in best practices and commitment of time as the conversation unfolds, and run the risk of silencing voices or biasing results.

There are a wide range of methods used by groups to process information together and form collectively considered positions. Citizen assemblies facilitate intensive deliberation, at a considerable cost of time and money (Landemore, 2022). Focus groups can produce robust summaries more quickly, but remain expensive and do not scale. Meanwhile, traditional polling sacrifices robustness, nuance and emergence in favor of inclusiveness and quantitative interpretability. Polis bridges qualitative and quantitative methodologies to produce a kind of rough ethnography at greater speed and lower cost relative to citizen assemblies or focus groups, trading some robustness for the speed and inclusiveness of a traditional poll. Nevertheless, there remains significant potential for further improving outcomes *and* affordability.

## 1.3   Addressing Polis Challenges with LLMs

LLMs are artificial neural networks trained on massive corpora of training data (typically text) in order to predict new textual completions (OpenAI, 2023; Touvron et al., 2023; Bai et al., 2022; Schulman et al., 2022). Importantly, this ability can be used to generate the most contextually relevant and plausible text as a response to or continuation of a provided prompt. With careful prompt engineering, LLMs can be used to analyze, synthesize, categorize and in other ways process and produce textual content. LLMs can also be used to create *textual embeddings*, which map bodies of text into high dimensional "semantic spaces" which numerically capture the semantic essence of the text in question. These embeddings can be used to measure semantic similarity between bodies of text, among other uses.

LLMs offer a number of opportunities for improving Polis as a deliberative platform, which we investigate in detail in Section 2, ordered from most feasible and lowest risk to speculative and high risk:

**Topic Modelling:** Identifying topics among comments to assist with analysis and reporting.

**Summarization:** More automated production of nuanced, digestible reports and summaries of the outputs of engagements.

**Moderation:** Reduce the burden of moderating comments.

**Comment routing:** Make better use of participant and moderator time by more effectively "routing" comments to participants for voting, i.e. deciding which comments should be voted on first.

**Identifying consensus:** Discovery of statements with broad consensus between identified opinion groups (*group informed consensus*; Small et al., 2021), as well as those which describe the worldview of specific groups.

**Vote Prediction:** Based on how a participant voted on N statements, can the system predict agreement or disagreement on a random unseen statement?

However, while LLM technology has massively broadened the scope of what's possible to achieve with Natural Language Processing, there remain significant limitations and risks associated with their use (Ganguli et al., 2022b,a; Weidinger et al., 2021; Bommasani et al., 2022; Bender et al., 2021b). These risks are especially critical to address if LLMs are to find a place in deliberative technology, where the impact, by the virtue of their purpose, is both societal and political. Despite significant efforts on the part of researchers developing LLMs, there remain issues around bias (Santurkar et al., 2023; Hutchinson et al., 2020; Kurita et al., 2019; Basta et al., 2019; Abid et al., 2021b; Sap et al., 2020), prompt injection attacks (Greshake et al., 2023; Perez and Ribeiro, 2022), so called "hallucinations" (LLMs making things up when they don't know something; Ji et al., 2023b), and basic transparency. Therefore, it is our position that the utmost care be taken in considering how LLMs are employed in deliberative settings.

Because Polis is designed to facilitate participant feedback on text, it is particularly amenable to addressing some of the risks associated with the application of LLMs described above. By integrating human evaluation of LLM-generated content as a part of the participation process, it may be possible to both mitigate risk, and further empower

participants. Where necessary, dedicated human moderators and facilitators can also serve as a safety check.

**We believe that maintaining human feedback is vital as a starting point for safe application of these immensely promising technologies to the task of deliberative engagement.**

## 1.4  Related Work

Recently, there have been several efforts to utilize LLMs in political or sociological contexts. For example, Romania's government has created Ion, an Artificial Intelligence Agent that scours social media for content, solicits feedback from citizens, and reports back its findings to the nation's Prime Minister, who can also interact with it directly, e.g. posing questions about public sentiment (France-Presse, 2023). DeepMind has fine tuned an LLM specifically designed to generate statements that resonate with a broad audience, based on a given set of input statements (Bakker et al., 2022). The Remesh platform (Bilich et al., 2019) has pioneered the use of LLMs for inferring opinion data in the absence of votes across and within conversations by integrating semantic similarity via text embeddings into its Latent Factor Model (Konya et al., 2022).

This diverse set of examples reflects the broad design space of LLM application to deliberative technology. In some of these cases, such as Romania's Ion, we see an LLM being used in a very unconstrained manner where there is very little risk mitigation. By contrast, DeepMind's fine-tuned model is modular, and thus presents opportunities for using in a responsible manner, in which it is continually supervised.

Our motivation for this paper is to consider the broad design space of LLMs as applied to Polis, and identify strategies for taking advantage of the many opportunities presented by LLMs, while mitigating the risks they might pose.
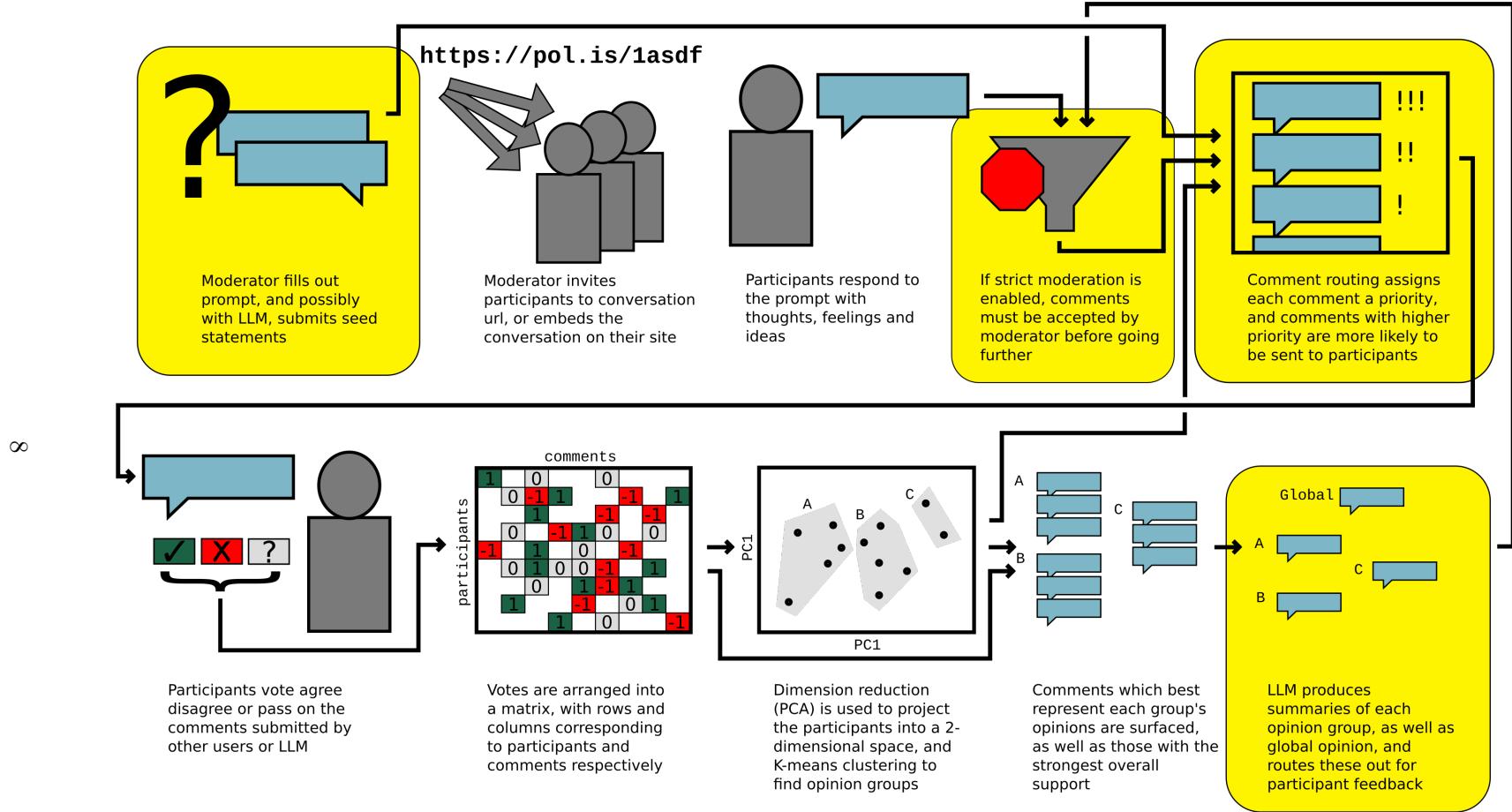
Figure 2: Overview of the process of a Polis conversation. Steps where we envision the support of LLMs are highlighted in yellow. Full methodological details of the current (non-LLM) Polis processing flow are described in depth in Small et al. (2021).

# 2 Opportunities and Risks

Conducting an effective engagement with Polis involves a number of distinct tasks for facilitators and participants. A summary of the system is presented in Figure 2, which additionally highlights the steps which could benefit from the use of LLMs. In this section we discuss the particular challenges we believe to be most amenable to the application of LLMs, as well as the associated risks and potential mitigations.

To evaluate the potential for LLMs in these applications, we performed a number of experiments using Anthropic's Claude (Bai et al., 2022). We describe these experiments and their results in the sections below. We include the prompts in Section A.

We ran the experiments primarily with open Polis data from a conversation run in Bowling Green, Kentucky in 2018 (Barry, 2023; Sergent, 2018; McKenzie, 2018). At the time, Bowling Green was deeply divided by national hot button issues, so The Computational Democracy Project (the non-profit organization maintaining Polis) teamed up with local media partners at Bowling Green Daily News to explore whether Polis would allow residents to unite on common interests. Residents were asked to respond to the prompt "What do you believe should change in Bowling Green/Warren County in order to make it a better place to live, work and spend time?" (See Figure 1 for a depiction of the participation interface.) As expected, the community was able to identify local issues on which there was general consensus, and hold their officials accountable to act on that consensus. Some of the issues on which points of consensus were identified included traffic, development, internet access, and public accountability. This engagement had very high participation rates relative to total population size: 2000 out of 65000 residents participated, accounting for around 3% of the population. While this might seem minor for an election, it is a remarkably high proportion for a deliberative exercise, where typical engagement for such a town might be only 0.1% - 0.3%. The automated report, raw data export, and additional information for this conversation is available in Barry (2023).

## 2.1 Topic Modelling

One of the simplest but most versatile uses for LLMs (and NLP more generally) in Polis is topic modelling, which assigns abstract "topics" (short descriptions) to larger bodies of text. Categorizing comments under discrete topics can facilitate numerous analyses, visualizations and summarizations which can help people understand the outcome of a Polis engagement in greater depth. In the past, we have used human-assigned topics in custom reporting to compare how contentious different topics are (vote variance per topic), analyze how voting patterns correlate between topics (using the RV coefficient; Robert and Escoufier (1976)), and in narrative summaries to simply enumerate the topics which emerged in broad discussions.

We experimented with generating topics from a Polis conversation by prompting Claude to identify topics in sets of comments from the conversation. Because LLMs are limited in how much text they can respond to (the so called *context window length*), comments were assigned topics in batches, and then the resulting batch topics merged.[1]

---

[1]This experiment was performed on a research model with a context window length of around 8K tokens. Since then, Claude has been updated to have a context window length of 100K tokens, which

| Group | Topics |
| --- | --- |
| Government and Public Policy: | Local government and politics |
| | Laws and regulations |
| | Taxes and services |
| | Transparency and accountability |
| Infrastructure and Development: | Housing |
| | Transportation |
| | Utilities |
| | Historic preservation |
| | Urban planning |
| Public Services: | Healthcare |
| | Education |
| | Public spaces |
| | Homeless services |
| | Disability access |
| Safety Health and Environment: | Law enforcement |
| | Public health |
| | Pollution |
| | Emergency management |
| | Marginalized groups |
| | Animal control |
| Economy and Business: | Job opportunities |
| | Local business support |
| | Competition and legalization issues |

Table 1: Claude topic modelling results as returned in response to the prompt shown in Appendix A.1.

See Appendix A.1 for details about the prompts and systems used.

### 2.1.1 Experimental Results

While not explicitly asked to do so, Claude automatically produced hierarchical topic assignments, at two levels of depth. This allowed for a nice breakdown of the conversation, and is expected to be useful for producing summaries and reports of varying detail, as discussed in 2.2. The resulting topics matched our expectations from manual analysis of the conversation, and are presented in Table 1.

### 2.1.2 Risks and Mitigation

Topic modelling has fewer risks relative to some of the applications of LLMs we discuss in later sections. We foresee only one: Assigned topics could either be inaccurate or miss important contextual nuance. To address this, human review and overriding of inferred topics must be supported, as well as manual topic assignment.

---

we used for the experiments in Section 2.5.

### 2.1.3    Discussion

We anticipate that this inherently iterative approach will naturally adapt well to an online system, updating topics as new comments enter the conversation. While we did not rigorously evaluate how stable topics were over iterative application, we found that they were indeed largely stable as we experimented with prompts. Issues which arise with stability on more complex or nuanced datasets may be dealt with by tuning the prompt to prefer existing topics over new ones, similar to the sequential construction of topic models in Bayesian nonparametric statistics (Griffiths et al., 2003; Ghahramani and Griffiths, 2005).

How human overriding of topic assignments adapts to an online recursive application of this technique remains to be seen. However, from our experience with the flexibility of LLMs to prompting, we again expect this could be achieved by including human preference in the prompt.

## 2.2    Summarization and Reporting

Summarization and reporting are an essential but costly part of deliberative processes. In a citizens' assembly, facilitators summarize the perspectives they are hearing as they oversee a deliberation, as well as when they produce a final synthesis for participants, stakeholders, policy makers, and the public at large. In parliamentary settings, each member and their staff would be responsible for maintaining an internal representation of a multi-stakeholder environment.

Summarization and reporting are an essential but costly part of using Polis in a deliberative process as well. Polis supports summarization of the data it produces by providing an automated report which, based on the learned PCA projection and opinion groups, highlights comments with high overall support, as well as those which best explain how each opinion group differs from the rest of the conversation. However, the automated report is still close to the raw data, and not suitable for consumption by, e.g. a mayor, without dozens of hours of manual human analysis of this information. As mentioned above, this requires a significant level of expertise and training, which constrains the adoption and ultimate impact of these processes.

For some domains, LLMs have already been shown to be effective at summarizing textual and numeric information, in some cases judged as on par with human written summaries (Zhang et al., 2023; Singh et al., 2023). LLMs would thus seem to have immense potential for reducing the burden of this critical task. However, many LLMs still have significant context window length limitations, which pose a constraint in how they can be applied to large conversations, which might have many thousands of comments. Moreover, given the critical nature of this aspect of the deliberative process—literally the mechanism by which meaning is made—it's imperative that it be thought through carefully.

The overall organization of the current report—presenting comments which best characterize or represent each opinion group, followed by points of consensus—has already proven largely successful in serving as an anchor for meaning-making exercises (see Small et al., 2021, for details).

Custom reporting typically builds on this pattern, adding narrative descriptions to

summarize the perspectives of each group, and the participant body as a whole. An LLM can similarly summarize each opinion group in turn, as well as points of consensus, using the same comment selection routines featured in the automated report. Prompts could even instruct the LLM to cite specific comments to backup its claims, and these could be used to insert vote distributions figures like those which show up in the current report.

While not a feature of the current automated report, custom reports sometimes dive deeper by exploring the topics which emerge in the conversation. In these cases, topic modelling is performed by humans, and reports spend time investigating results on a topic by topic basis. Given automated topic modelling, as discussed in Section 2.1, comments and voting patterns associated with a given topic can be summarized by a LLM, together with the associated group-by-group voting data.

These techniques address context window limitations, while scaffolding upon known and effective patterns. For large conversations, it may be necessary to filter comment sets by selecting for higher quality comments, based on number of "Pass/Unsure" votes, etc. In cases where context window length is still an issue, summaries can also be performed recursively; An initial batch of comments can be summarized, followed by a prompt to update the previous summary with a second batch of comments, and so on. This reflects the pattern of recursive summarization explored with human summarizers in Zhang et al. (2017).

### 2.2.1 Experimental results

To investigate the potential for using LLMs to automatically generate concise and coherent summaries, we took batches of comments from the Bowling Green, Kentucky conversation, including vote data, had Claude generate summaries, and then merged the summaries across batches. The comments and prompt are presented in Appendix A.2, and resulting summary shown below.

The following summary was produced using a small set of comments, together with overall number of agree and disagree votes:

---

**Claude's output:**

```
In summary, there was general agreement on issues around traffic, public
amenities, media/technology, urban planning, and transportation. Views
were more mixed on schools, neighborhood issues, healthcare, economics,
and community support. A few lighthearted or isolated comments did not
show clear agreement or disagreement. Overall, the comments indicated
interest in improving infrastructure, quality of life, and community
support in the area. But there were some polarizing views on select issues
like immigrant language ability or reactions to the opioid crisis.
```

---

When asked to summarize points of agreement and disagreement between groups (using the percentage of agrees and disagrees per group), the following was produced:

---

**Claude's output:**

```
Overall, there was more agreement in Group 0 on infrastructure, housing
```

---

```
and public services issues while more disagreement in Group 1 on those
same topics. Both groups agreed there is a problem with rental properties
deteriorating. Group 0 mostly disagreed with comments related to the
opioid crisis while Group 1 was split. The summary shows the dominant
topics of discussion and how much each group agreed or disagreed with
them.
```

Both of these summaries concur with manual analysis of the data, and demonstrate
that LLMs are capable of summarizing data from Polis conversations. There is room
for improvement however; e.g. "Group 0 mostly disagreed with comments related to
the opioid crisis" could more clearly specify which side of the issue the group falls on
(e.g. "Group 0 mostly agreed with statements expressing concern for the opioid crisis").
However, additional efforts with prompting (possibly providing examples) are likely to
address this.

### 2.2.2   Risks and Mitigation

LLMs have potential to assist in automatically generating summaries of deliberative
processes. However, it is important to note that LLMs can generate misinformation or
fabricate details while summarizing (Ji et al., 2023a; Stiennon et al., 2020; Bommasani
et al., 2022). Detecting these errors automatically is challenging. Prior work has pro-
posed automated metrics to detect such errors (Maynez et al., 2020; Durmus et al.,
2020; Kryscinski et al., 2020; Goyal and Durrett, 2021); however, they do not always
correlate well with human judgements (Fabbri et al., 2021; Durmus et al., 2022).

Furthermore, LLMs may reflect and amplify societal biases present in their training
data (Bender et al., 2021a; Abid et al., 2021a; Cheng et al., 2023) – including political
biases – which may affect what information they extract and emphasize when summa-
rizing deliberative discourse (see Section 3.2). This could potentially lead to skewed
or unrepresentative summaries that fail to capture the full spectrum of perspectives
expressed.

We believe human involvement is critical for checking machine-generated summaries
of deliberative processes to ensure accuracy and fairness. One specific innovation we
intend to implement towards this end is to show the automated summaries back to the
participants for review and feedback. Prior work has found this type of participatory
approach, where the users and stakeholders of an AI system are involved in evaluating
and improving its outputs, can help identify errors, biases, and misrepresentations (Lee
et al., 2022; Yang et al., 2018).

By soliciting feedback directly from the participants whose perspectives were sum-
marized, we can collect targeted annotations on the quality of the summaries. To avoid
overwhelming participants, these summaries should be concise, likely at the level of
their groups' perspective on a particular topic or subtopic, building up to more com-
prehensive summaries. These human judgments can be used to as a feedback to the
models and better align the automated summarization with what participants feel cap-
tures the essence of the discourse, a sort of Deliberative Reinforcement Learning from
Human Feedback (RLHF) (Ziegler et al., 2020; Lambert et al., 2022). This human-AI
collaborative approach, combining the scale of LLMs with human insight, has potential
to produce higher-quality and more trustworthy summaries.

| | |
|---|---|
| **Human-written summary of 10 comments about the San Juan County Land Bank.** | "The Land Bank should better support<br>- recreation (hiking, biking, fishing, hunting, equestrian)<br>- firewood harvesting<br>- better care of farm land<br>- public accessibility and knowledge<br>- acquiring new preserves" |
| **Mean representativeness score across comments** | **3.25 out of 5** |
| **Least represented comment** | "The vast majority of Land Bank managed properties are unknown and/or inaccessible to the vast majority of island residents."<br><br>**1 out of 5 - This comment is not represented at all; the summary ignored the comment entirely.** |
| **Best represented comment** | "The Land Bank needs to take better care of its' farm land. Some farmers are doing a good job but others don't seem to care."<br><br>**4 out of 5 - The summary covers most of the information in the comment, but is missing some nuance or detail.** |

Figure 3: Example of using Claude to help evaluate a human-written summary. The evaluation of the representativeness score is done by Claude as specified in the prompt in Appendix A.6 for details.

Whether generated by humans or LLMs, summaries inevitably leave out details, and may not capture every perspective. It's critical that these omissions be transparent, but evaluating a summary can be even harder and more time consuming than generating one, particularly for humans. We can assist evaluation by decomposing the problem into a set of much smaller and narrower tasks: asking an LLM to evaluate how well a summary represents each individual comment. This can augment a human evaluator by drawing their attention to the comments least represented by a given summary, using an interface such as the one shown in Figure 3, and prompting shown in Appendix A.6.

### 2.2.3 Discussion

The interface in Figure 3 could be taken a step further towards *intelligence augmentation*: rather than replacing human summarizers, we can use language model evaluations to augment their ability to create good summaries, empowering participants to place the summary in its proper context. One could imagine the rapid feedback loop of a report editing environment which, as you type, immediately highlights any comments or topics that are not adequately represented by the text. This would allow a human analyst to either take automated summaries as a starting point for a custom report, or write a report from scratch, and have guidance about what perspectives could be better reflected.

### 2.3 Facilitating Conversations by Synthesizing Group Identity and Consensus

In live deliberations, a human facilitator may summarize the standpoint or worldview of a group expressing its opinions as the discussion unfolds, to determine whether the group assesses that it is well understood by the facilitator. Importantly, this process

is dialectic: the facilitator reflects back their understanding of what they've heard, recursively asking for confirmation or further clarification until a fixed point is reached, and each side feels understood. This process of building shared understanding is a core pillar of the Nonviolent Communication approach to conflict resolution (Rosenberg, 2005), which was profoundly influential in shaping the design of Polis.

Better understanding where each side is coming from increases the likelihood that points of common ground are able to serve as anchors for bridging divides; respecting points of dissent is key to distilling solutions with broad support.

While summarization and reporting on the final results of a Polis conversation (see Section 2.2) is an important part of this process, summarization may happen during the Polis conversation as well. Experienced facilitators of Polis conversations inject statements into the conversation as a way of testing hypothesis they may have about the emerging opinion landscape as it unfolds. This is especially true of potential points of consensus.

Given a set of statements with broad support from across the participant body, an LLM could craft additional statements which it suspects have similar support, and specifically represent points of consensus. Similarly, given a set of statements with broad support from a particular opinion group, it could craft statements which it expects those participants will tend to agree with. Of note, DeepMind has researched and developed LLMs specifically tailored towards generating content with broad support based on a set of input statements (Bakker et al., 2022), an approach which might be perfectly suited to this task.

We hypothesize that by reflecting back to the participants syntheses of their perspectives:

1. Points of common ground or consensus surface faster.

2. Conversations "converge" faster, as measured by our ability to accurately predict votes based on past voting data, and quickly understand where people fall in the opinion landscape.

### 2.3.1 Experimental Results

The first approach we tested for generating consensus statements involved feeding the overall summary (see Section 2.2.1) into a prompt which asked Claude to generate new consensus statements based on the summary (see Appendix A.3 for details). This produced the following statements:

> **Claude's output:**
>
> - Our infrastructure and public amenities could use continued improvement and investment. Things like roads, parks, transit, etc. enhance quality of life for all.
>
> - We should advocate for better urban planning and transportation options. Making the community more livable, walkable and connected would benefit everyone.
>
> - While we may disagree on some issues, we all want a safe, supportive

```
      and thriving community. Compromise and understanding can go a long
      way.

    • Access to healthcare and economic opportunity are important for any
      community. We should work to enable wellbeing and prosperity for as
      many as possible.

    • Technology is constantly changing, but we should ensure changes
      benefit and connect our community, not isolate us. Education and
      connectivity for all ages matter.

    • There is always room for improvement, but appreciation for what we
      have and the diversity in our community is so important too. Our
      differences make us stronger together.
```

While there is no way to assess the status of these statements as potential points of consensus without feeding them in to a live conversation with human participants, they nevertheless capture the typical tone and style of consensus statements we've observed in real conversations. Testing this work in situ would have required a complete integration of the method into the software, and was thus out of scope for this study.

### 2.3.2 Risks and Mitigation

Key questions remain for the ethical application of LLM technology as described in this section:

1. What thresholds should exist for the fraction of comments in a conversations generated by LLMs, as apposed to humans?
2. What metrics and procedures should be involved in deciding when and how often to submit LLM generated content?
3. Should these statements be routed (probability of being displayed for votes) any differently than human-generated statements?
4. How should their generated nature be disclosed?
5. How might we measure the impact of nudges the LLMs would insert in the process?

Underlying all of these concerns is the question: How much machine influence is acceptable in a process which ultimately aims to surface human opinions? Perhaps human in the loop review of LLM generated statements is sufficient to ensure ethical application of this technology. But further consideration seems warranted regarding what additional constraints should be placed upon the quantity and nature of its submissions.

Especially important to consider is the mapping between minority opinion groups and the continued consideration of the core Polis group-informed consensus mechanism (Small et al., 2021), ensuring there is no reversion to prioritizing majority rule.

### 2.3.3 Discussion

Ethically, it's important that participants be aware that they may be looking at a LLM generated content. Moreover, many LLM Terms of Service (including those of OpenAI's GPT and Anthropic's Claude) require that people be informed when they

are looking at model generated content. However, because knowledge that a comment was written by an AI would likely affect the way people vote on it, LLM-generated comments would (from a purely statistical perspective) ideally display no differently than human-generated comments, possibly as having been anonymously written.

Addressing both of these concerns, the interface might notify participants when a conversation has LLM generated content, without specifying *which* comments, until the conclusion of the conversation. Alternatively, the interface could notify participants when a statement was generated by an LLM immediately after they had finished voting on it.

## 2.4   Vote Prediction

The ability of LLMs to predict votes in a Polis conversation is of both theoretical and practical interest. From a theoretical standpoint, predicting votes is a valuable test of how capable LLMs are of understanding the opinion landscapes which emerge from Polis conversation. Experiments have already demonstrated that LLMs trained on media content can predict public opinion (Chu et al., 2023). More practically, vote prediction has significant potential to improve how Polis handles missing data, both in reporting results, in the core clustering algorithms, and as a signal for comment routing (see Section 3.4.3).

Large Polis conversations can have thousands of comments, while most participants vote on fewer than 100 comments. So over 90% of the data in the participant × comment matrix may be missing for such conversations.

The comment routing algorithm tries to ameliorate the missing data problem by surfacing the most relevant comments (which either help explain the opinion space, or garner broad support), but it has limitations. Without careful engagement design, early participants miss out on the opportunity to come back and respond to comments which emerge later in the conversation. Participants in some opinion groups tend to vote more than others, complicating interpretation of the total vote breakdown. This problem can be mitigated by looking at votes group by group, but we must also consider whether this phenomena operates on an intra-group basis. Inferring missing voting data therefore has the potential to improve our ability to understand and interpret the results of Polis conversations.

### 2.4.1   Experimental Results

We evaluated vote prediction by prompting an LLM with a participant's previous voting history and asking if the participant would agree or disagree with a given comment (see Appendix A.4 for details). Surprisingly, we found that a plain LLM (Askell et al., 2021), without fine-tuning or RLHF training, was nearly perfectly statistically calibrated at predicting agreement, as seen in Figure 4. In a majority of cases the LLM was at least 90% confident in its prediction, and correctly calibrated even at that high confidence level. This suggests LLMs are already very capable of understanding human opinions, presenting both opportunities and serious risks.
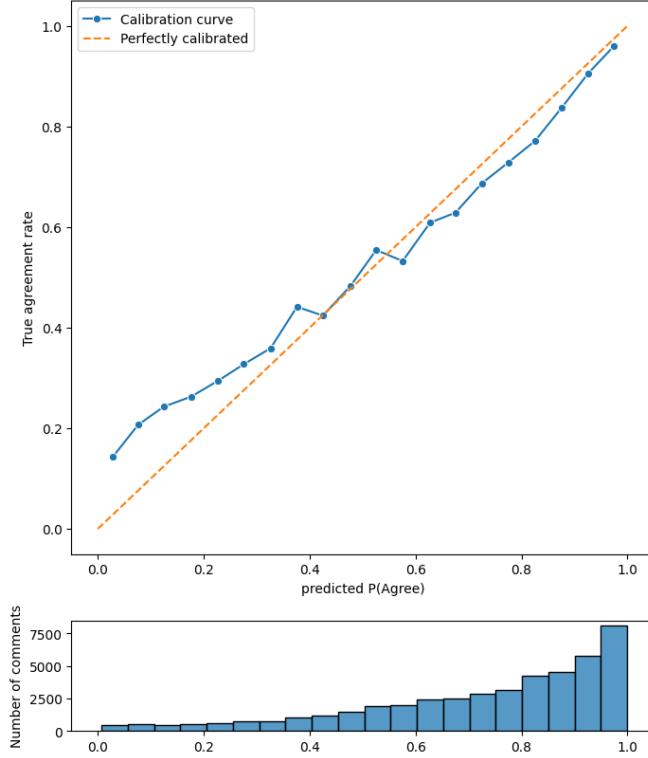
Figure 4: Calibration plot of LLM-predicted probabilities of participant agreement with comments. The probabilities predicted by the LLM are close to perfectly calibrated across comments.

### 2.4.2 Risks and Mitigation

Systematically inaccurate vote prediction runs the risk of misrepresenting public opinion on issues that arise within a conversation. Given the accuracy of vote prediction on an individual basis, it's possible that in aggregate the inference across participants will be even more accurate. However, participant bodies which deviate from stereotypical ideological configurations may be systematically misrepresented by these estimates.

Care must be taken to adaptively track and adjust predictions for opinions which break the LLM's expectations. Other mitigations of this risk might include: presenting both raw and inferred vote tallies, presenting inferred tallies only for comments with high inference confidence, aggregating inference information only for those participants for whom we have high confidence, or presenting inferred tallies with estimated error margins.

The existence of powerful vote prediction technology creates incentives to entirely replace human participation with *in-silico* deliberation. This runs the risk of amplifying existing biases as well as eliminating the many positive externalities of deliberation on mutual understanding, civic empowerment, surfacing leaders, etc. We strongly advocate that vote prediction be used only to amplify participants' voices, never to replace them. We reject as invalid on both ethical and performance grounds the use of vote prediction technology to replace human participants with simulations.

We cannot enough emphasize the latter point: **it would be catastrophic for de-**

**liberation at scale if the remarkable capabilities of LLMs lead to replacement of whole groups of individuals by simulacrums designed by a very different population.** Besides the question of accuracy and biases, we believe this could also lead to a crisis of faith and belief in the deliberative process.

### 2.4.3   Discussion

More speculatively, vote prediction could be used as another way to evaluate the completeness of Polis reports and summaries. A good summary of a conversation should allow its readers to accurately predict the participant $\times$ comment vote matrix. Users whose votes are particularly hard to predict could be surfaced as statistical "outliers" not well-represented by the summary, as with comments poorly represented by summaries (as discussed in Section 2.2).

This information could also be presented to participants themselves as a sort of "progress bar" or measurement of Polis' ability to understand their position in the opinion landscape. Right now participants are shown the number of comments remaining, which for large conversations will typically read "100+" in perpetuity, until they grow weary and close the conversation. Grading how well we understand a participant's position will give them a clearer sense of where they stand, and better judge how much longer they want to participate.

These metrics could also be used to determine when it might be a good time for a participant to return to the conversation once new and more interesting comments have changed Polis' certainty about their position in the opinion landscape. However, care should be taken to evaluate the net effect this has on engagement. While it's possible that engagement improves among participants who tend to vote on fewer comments, the opposite may be true for those who tend to vote more. It is also important to consider how systematic errors in inference might affect participation levels, since overconfidence could encourage participants to leave early.

## 2.5   Leveraging Large Context Windows

The above designs and experiments were carried out assuming the constraint of context window lengths typical of most LLMs. Since performing those experiments, Anthropic released a model with a 100k-token context window length, capable of including the entire set of comments from the Bowling Green conversation. To explore how this affects both the design space and qualitative results of applying LLMs towards the challenges addressed in Section 1.2, we performed several experiments which performed both topic modelling and summarization in a single pass against the entire set of comments from the Bowling Green conversation, including vote breakdowns per group.

### 2.5.1   Experimental Results

The first experiment asked Claude to perform topic modelling on the provided comments, then quantify the level of agreement vs disagreement between the two groups in the conversation, and summarize differences in opinion between the two groups. The results are shown below:

**Claude's output:**

1. Traffic and transportation - Comments about traffic flow, public transit, sidewalks, bike lanes, road widening, traffic lights, etc. This topic had a high level of agreement between Group 0 and Group1.

2. Housing and development - Comments related to affordable housing, rental prices, property taxes, neighborhood revitalization, zoning, annexation, urban planning, downtown development, etc. Group 0 tended to agree more on increasing affordable housing and revitalizing neighborhoods while Group 1 agreed more on limiting multi-family developments.

3. Healthcare and public health - Comments about healthcare access, opioid addiction, mental health, medical marijuana, healthy eating options, etc. Group 0 agreed more on treating addiction as a health issue and increasing healthcare access while Group 1 agreed more on harsher sentencing for drug users.

4. Education - Comments on school funding, class sizes, curriculum, higher education, etc. Group 0 agreed more on funding schools and community colleges while Group 1 agreed more on schools helping address budget cuts.

5. City laws and policies - Comments about fairness ordinance, noise ordinances, smoking bans, term limits, transparency in government, etc. Group 0 overwhelmingly agreed more on passing a fairness ordinance while Group 1 disagreed. Group 1 agreed more on enforcing existing laws.

6. Parks and recreation - Comments about increasing green spaces, playgrounds, community activities, tourism opportunities, etc. Group 0 agreed more on improving parks and recreation opportunities while Group 1 agreed at a lower level.

In summary, while Groups 0 and 1 showed high agreement on some topics like improving traffic flow and infrastructure, there were clear differences of opinion on social issues like healthcare, education, housing, and LGBTQ rights. Group 0 expressed more progressive views supporting increased government investment in public services and protecting marginalized groups. In contrast, Group 1 expressed more conservative views supporting limited government intervention and existing laws. The topic modeling analysis revealed these opposing viewpoints and areas of disagreement between the two groups.

The second experiment omitted topic modeling as an intermediate step in the prompt, instead asking for a summary of the key points of consensus and divergence between the groups. The results:

**Claude's output:**

Based on the comments and votes from Group 0 and Group 1 on each comment, here are the key points of consensus and divergence that I identified:

```
Consensus between Group 0 and Group 1:

    • Most comments that received strong agreement in both groups focused
      on investing in the community good, improving infrastructure and
      public services, supporting the local economy, and protecting the
      environment. For example, comments advocating for more bike lanes,
      sidewalks, public transit; investing in the arts and education;
      developing the downtown and riverfront area; increasing green spaces;
      improving traffic flow and infrastructure; and supporting local
      businesses were agreed upon by both groups.


Divergence between Group 0 and Group 1:

    • There were some notable differences of opinion on social and po-
      litical issues. Group 0 tended to be more progressive while Group 1
      tended to be more conservative. For example, Group 0 largely agreed
      with and Group 1 largely disagreed with comments supporting LGBTQ
      rights, marijuana legalization, and sanctuary cities. Group 0 dis-
      agreed and Group 1 agreed with restricting immigration and requiring
      immigrants to learn English.

    • Group 0 was more supportive of government intervention and social
      programs to address issues like poverty, homelessness, healthcare,
      and education. Group 1 was more skeptical of government overreach and
      supportive of free market principles. For example, Group 0 agreed and
      Group 1 disagreed with raising the minimum wage, providing affordable
      housing and healthcare options, and increasing funding for schools
      and social programs.

    • There were also differences in levels of trust in institutions
      like the media, government, and universities. Group 0 expressed more
      trust in these institutions while Group 1 was more skeptical and
      critical. For example, Group 0 agreed and Group 1 disagreed that
      most journalists and public officials aim to serve the public good.
      Group 0 disagreed and Group 1 agreed that university mismanagement is
      mostly responsible for increasing tuition.


In summary, while there were areas of agreement on community investment
and development, the groups differed in their social and political values
as well as levels of trust in public institutions. But there is clearly
some common ground to build upon, if the goal is to find consensus and
work together.
```

These summaries demonstrate a marked improvement over the results from Section
2.2. Not only were they were more detailed and nuanced, but they succeeded in better
explaining what position participants (either as a group or whole) took on a particular
issue.

### 2.5.2   Risks and Mitigation

Many of the risks inherent in this application have already been discussed in previous
sections (notably Sections 2.2 and 2.1). For topic modelling specifically, this approach

would seem to be less risky overall, as it may be easier to glean a comprehensive but accurate topic breakdown when considering the full set of comments all at once. For summarization, there may be some benefit to subsetting comments based on other criteria (e.g. Polis' PCA or K-means grouping) prior to applying LLMs for summarization, since we have some external and explainable basis for contextualizing how those comments relate to each other. However, this comes at the cost of not being able to consider the full conversation, which could lead to less nuanced summarization of the data, and so here as well, the large context window approach may ultimately be less risky.

### 2.5.3   Discussion

In this study, we leverage Claude's 100k context window to distill intricate opinion landscapes from extensive conversational data, a task usually demanding significant manual effort. To avoid bias, we asked the LLM to piggyback its analysis on deterministic statistics which were gathered the usual way one would with Polis, and then to summarize, as someone might with dozens of hours of data analysis.

Using a large context window has proved beneficial for this particular task because it facilitates integration and correlation of diverse pieces of information in a unified manner. This, in turn, creates a richer understanding and representation of the data, akin to a "broad perspective". This holistic approach fosters an understanding that is greater than the mere sum of its parts, and proved far superior to results obtained with a smaller context window length.

The strength of this result is not limited to zero shot summarization capabilities; as Polis is a real-time system, participants can collectively respond to and rate summaries of their positions potentially hundreds to hundreds of thousands of times during the course of a deliberation.

With these findings, new areas of exploration open up. For instance:

1. Given the real-time operation of the Polis system, it could be worthwhile to examine how the language model can interpret evolving voting patterns over time, or recursively consider notable shifts between previous iterations of analysis.

2. Leveraging the system's internal representation of the stakeholder groups to extend Polis' real-time capabilities into multi-stakeholder dialog using a different approach than attempted in the CICERO research.

3. Allow participant feedback on summaries, even in plain text (e.g. "I'm not part of group Y for X reason") as the the basis for a kind of novel hierarchical clustering algorithm, which might draw out expert sub-populations and niche interest groups.

## 3   Discussion

While this paper primarily explores the application of LLMs towards optimizing Polis as a deliberative platform, the considerations here apply much more broadly towards deliberative systems which wish to embrace LLMs and other Machine Learning technology.
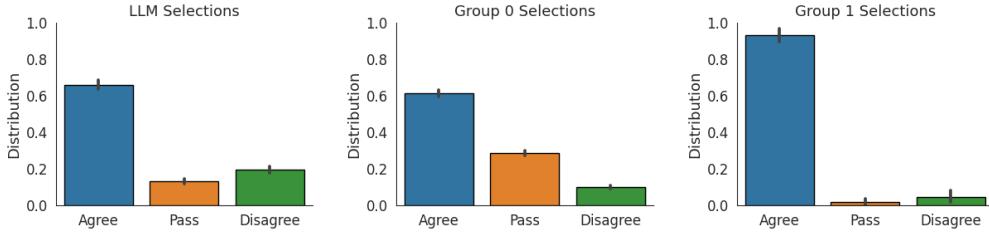
Figure 5: Overall vote distributions of Claude and human opinion groups.

## 3.1 Intelligence Augmentation, Not Human Replacement

Already, as mentioned in Section 1, we see attempts to use in silico deliberation to replace or obscure human voice and agency (France-Presse, 2023).

As broadly appealing as the idea is of having limitless insight into public opinion at beck and call, removing the opportunity for humans to interact (deliberate) with other humans fails to respect human agency in building shared understanding, and runs the risk of misrepresenting and thereby doing violence to public will.

It is therefore vital that wherever LLMs are used in such systems to generate deliberative content, there be an opportunity for humans to be included in the process of evaluating, verifying and contextualizing that content, augmenting rather than supplanting human intelligence.

## 3.2 Measuring Ideological Bias

As we begin to adopt LLMs for the tasks of synthesizing information, building consensus and potentially even assisting with moderation (see Section 3.3), it's important to stop and evaluate whether LLMs themselves have biases which might influence their output.

While measuring bias is a domain of active research where even core definitions can be debated, the use case of online deliberation offers a very concrete set of measures: By having a LLM vote on comments as part of the discussion, we can see how its votes relate to those of the human participants in the conversation, and therefore directly measure whether the LLM votes in a manner that is biased toward a particular opinion group.

Using the prompts described in Section A.5, we asked Claude to vote on each of the comments from the conversation. Because LLMs are probabilistic models, we can look at the probability associated with each response (vote choice), and compare this with probability estimates associated with each of the actual opinion groups. The distributions of these probabilities across comments are shown in Figure 5. This figure illustrates that Claude's voting pattern is distinct from both of the groups, but overall more similar to Group 0, the more socially progressive of the two groups. While we might expect a disembodied LLM with no lived experience in the town of Bowling Green, Kentucky to vote more frequently for "Pass/Unsure", it actually did so less frequently than Group 0 (but more than Group 1).

Comparing these vote probabilities, we can measure the response similarity between Claude and a given group as $1 - JensenShannonDistance$ of the vote probabilities for a given comment. The distribution of these values across comments are shown in
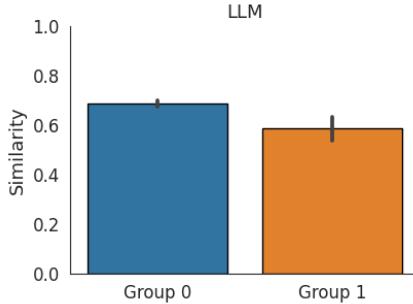
Figure 6: Claude votes more similarly to one of the human opinion groups.

Figure 6, which again shows that Claude is significantly biased towards Group 0 (T-test $P = 1.86 \times 10^{-4}$).

Santurkar et al. (2023) and Hartmann et al. (2023) have demonstrated evidence that other LLMs may trend left-leaning politically. In the role of deliberative assistant, this is potentially problematic if we are relying on a model as a multi-partial (Harding, 1986) facilitation agent when it aligns more closely with one of the groups than the other.

However, a complete lack of bias is impossible for both humans and LLMs. Some characteristics of facilitation practice include remaining curious and open to all sides (multi-partial) while attempting to help a group through the transformative process of bridging a divide. While we can attempt to instill the same attitude in LLMs with careful prompt engineering, our ability to measure LLM biases, as described above, provides us with a unique opportunity to monitor them relative to human agents.

## 3.3   Rethinking Moderation

Moderation serves two broad purposes: protecting Polis participants from exposure to abusive or hateful content, and optimizing participant time. The latter has until now typically led us to recommend moderators remove comments that are off topic or duplicates of existing comments. However, as we considered how LLMs might reduce the burden of moderation and participation, we have arrived at a new position on this matter: Moderators should only remove abusive content, with comment routing optimizing participant time.

One of the strengths of Polis as a system for sentiment analysis relative to traditional surveys is that participants are able to define the dimensions of the opinion space they map out in **their own words**. It's often precisely the questions survey crafters (or moderators) wouldn't think to ask which are the most pivotal in both understanding a space, and arriving at novel solutions to challenging political issues. Therefore, taking some of that power away from participants and putting it the hands of the moderators is ultimately disempowering.

Over the years we have observed instances of moderators removing perfectly reasonable positions which challenge the underlying assumptions and framing of the conversation prompt. Since observing this, The Computational Democracy Project has made a point of improving training practices to cultivate more principled ethical standards for moderation. But this is ultimately time consuming, and still leaves room for distortion.

24

Meanwhile, the incidence of truly "off topic" comments is typically infrequent, and thus the burden posed to participants relatively small. And ultimately, participants themselves have a way of signaling that a comment is of low relevance, which already feeds into the comment routing algorithm: Pass votes.

Semantically similar comments which might be considered "duplicates" are much more common. However, deciding how semantically close two comments must be to each other to be considered "duplicates" is not a straightforward task. Inconsistent application of this blurry threshold again leaves room for moderator bias to confound results. It is also the most time consuming aspect of the moderation process, since every comment needs to be evaluated in relation to all previously submitted comments. This burden is greatest for the largest of conversations, where participants are unlikely to vote on every one of the thousands of comments submitted before adding their own, leading to a higher number of near-duplicate comments.

There's a significant body of research which suggests that subtle shifts in the way positions are framed or expressed can have a huge impact on how people react (Lakoff et al., 2004; Lakoff, 2009), and our years of experience with Polis as a platform bears this out. There's also a strong case to be made that, in the absence of funded training, this is likely to have an outsized effect on historically marginalized communities, who may express themselves in ways which might not resonate as well with moderators from the status quo. These individuals may consequently have their voices silenced in favor of more "palatable" phrasings which further entrench privileged perspectives.

While we've seen that the existing comment routing system can and does boost one particular framing of a position over another under certain circumstances (when it either appeals to one group more than another, or has higher overall support from the conversation as a whole), this effect is frequently mild. As participants grow weary from responding to multiple rephrasings of the same sentiment, they may leave the conversation without responding to less frequently brought up (but nevertheless important) topics. However, in recognizing this particular cost, we see a potential resolution: By taking advantage of topic modelling efforts and potentially vote prediction (see Section 2.1 and Section 2.4) in comment routing (see Section 3.4.3), we can mitigate this problem without sacrificing the diversity of framing that makes Polis so valuable.

### 3.3.1    A Moderate Approach to LLM-assisted Moderation

Moderation of hate speech remains central to ensuring the integrity of a Polis conversation's result, as violent and/or dignity-impairing speech limits full participation (Habermas, 1962, 1981).

While continued support for strict moderation (requiring comments to be moderated in before being shown) can address this issue, the delay between submission and moderation reduces exposure to new ideas. A potential "third way" between strict and non-strict moderation might leverage LLMs (or other NLP technology) to place potentially toxic or hateful content in a queue for human review (Jahan and Oussalah, 2021; MacAvaney et al., 2019; Wei et al., 2021; Chiu et al., 2022). Relative to strict moderation, this approach would overall get new perspectives out more quickly, with a lower risk of exposure to problematic content relative to non-strict moderation.

### 3.4 Future Questions and Directions

The following sections discuss ideas which are either more speculative, rely on further technological advances, or which we simply did not have time to fully explore.

#### 3.4.1 Seeding Conversations

For most major Polis conversations, human facilitators submit *seed statements* to ensure that there is content there to vote on for the very first tranche of participants, and set the tone of the conversation (not dissimilar to how Sequential Design of Experiments, Active Learning (Cohn et al., 1996), and Bayesian Optimisation (Rasmussen, 2006) are seeded with a set of initial points). Most participants (on average, around 90%) do not comment in conversations, and only vote, so it is important to ensure that there is content there to start off the conversation which does a reasonable job of laying out some of the key perspectives expected to emerge. Given the potential impact these statements have on the direction of the conversation (Lakoff et al., 2004; Lakoff, 2009), it typically involves careful consideration of the opinion space ahead of time, known as *framing*, a task which requires investment in prior training. And while the beauty of Polis is that participants have the opportunity to turn a narrative or framing on its head, there's nevertheless the potential for these seed comments to have an outsized influence on the course of the discussion.

To address this issue, seed comments could be generated along the lines of Section 2.3. However, rather than input comments into this process which have been submitted by participants (obviously impossible if the conversation hasn't started yet), content be sourced from a combination of social media content and online archives of debate content, e.g. The Society Library (noa, 2023). This move would further empower the public at large in framing the terms of deliberations, and reduce the impact of a lack of training in issue framing by those setting up a conversation.

#### 3.4.2 Dimension Reduction and Clustering

Topic modelling and semantic embeddings could potentially be used to adjust how Polis performs dimension reduction (related work can be found in Konya et al., 2022). The current application of PCA effectively more heavily weights votes on issues for which there are a larger number of comments (since PCA looks for correlations in responses, a large number of semantically similar comments may have an outsized effect on the resulting opinion space). To some extent, this can be seen as a feature, in that frequently submitted opinions may be more important to participants, and thus deserve some amount of boosting in the learned opinion space. However, this could be exploited to entrench wedge issues and smother novel perspectives. With the planned addition of a feature allowing participants to indicate how important issues are to them (independent of any application of LLMs), we also have a more direct way of integrating this signal.

By weighting columns in the vote matrix according to how semantically similar the corresponding comment is to other comments in the conversation, we may be able to obtain a more robust and nuanced understanding of the opinion landscape. Multiple comments exploring different angles of a particular topic would still overall boost that topic's effect on the opinion space, but in a more measured way which allows for newer

topics to emerge and have an appropriate impact on the structure of the conversation. This would also naturally tie in with and support efforts to rethink Polis' approach to moderation (see Section 3.3) by reducing the pressure to deduplicate comments, since the algorithm would effectively be doing this work automatically.

### 3.4.3 Improved Comment Routing

Comment routing is the process which determines how comments are routed to participants for vote solicitation. Polis uses probabilistic weighting to determine which comments to route to participants, with comments more likely to be routed if they help explain position in the 2D PCA projection (opinion space), have higher overall agreement ratios, low pass ratios, or are new to the conversation. This system maximizes use of participant time, improves engagement, and ultimately produces information which is more informative and easier to analyze (Small et al., 2021).

One thing that is not currently accounted for in this system is topicality. By mapping comments to discrete topics, as described in Section 2.1, weighting can be shifted to boost newly emerging topics. However, LLMs can also be used to create embeddings of comments, for a more continuous notion of topicality. As discussed in Section 3.3, we would like for comment routing to be able to take over the task of handling duplicate comments from moderators. While we found in our work that topic embeddings failed to capture the directionality of sentiment (i.e. $A$ and $\neg A$ embeddings are almost identical), this can be disambiguated by combining the topic embedding with either raw vote data, associated PCA loadings, or feedback from an LLM.

An even more nuanced approach to comment routing might incorporate vote prediction (see Section 2.4). For example, we can improve utilization of participant time by boosting comments for which the LLM is less certain of the answer. More generally, incorporating uncertainty into comment routing is rooted in the fields of Active Learning (Cohn et al., 1996), also known as Sequential Design of Experiment in statistics, Bayesian Optimisation using Gaussian Processes (Rasmussen, 2006), Exploration/Exploitation in Reinforcement Learning and Bandits Algorithms Sutton and Barto (1998). We can build on this bedrock of methods with the added capabilities brought by LLMs. For larger conversations, it would be prohibitively expensive to query the LLM for every unanswered question, after every new vote, so some care would need to be taken in using this approach judiciously.

However, vote prediction holds the potential to improve comment routing more indirectly by relaxing a key constraint on how it functions. In particular, it has been a long standing assumption that comment routing should be implemented so that at any given point in time, the relative probability of receiving a particular comment is the same for all individuals in the conversation. The reason is that if we started differentiating comment routing on a participant by participant basis, there is the chance that there will be comments which are routed to participants in one opinion group more frequently than another, which confounds interpretation of conversation wide totals.

If we are able to infer missing data, and present more robust estimates of conversation-wide tallies in reporting, this frees up the comment routing algorithm to be more tailored on a participant by participant basis. Aside from routing comments based on uncertainty in how participants will respond, this also unlocks the potential for considering

how the order of comments affects state of mind (as evidenced by, e.g. Lakoff et al., 2004; Lakoff, 2009), and in particular, propensity to either polarize or find common ground.

Going into the details of how precisely the weighting algorithm should be modified to account for this information is outside the scope of this paper, as it's a complex problem deserving of a comprehensive treatment (and which will likely benefit from features currently in progress, such as allowing participants to indicate which comments are most *important* to them). We nevertheless anticipate that a combination of these ideas could vastly improve the system, and reduce the burden placed on moderators and participants.

### 3.4.4   Conversation Simulation

In Section 2.4.2 we strongly denounced the substitution of machine simulations for human participation when gauging public opinion. However, there are valid use cases for simulated deliberation, where results are not assumed to be an accurate reflection of the public will.

Synthetic data are commonly used in software engineering for testing system behavior and scalability. Ad-hoc, fixed datasets reflecting particular scenarios are often used for this purpose, but by virtue of being fixed, are limited in their ability to explore the state space and discover edge cases. For these reasons, *generative testing* and *property-based testing* have emerged as solutions which involve stochastically generating data to more thoroughly explore the possible state space (Claessen and Hughes, 2000; Pacheco and Ernst, 2007).

This technique could also be used as a deliberative or sociological "test bench" for rapidly and cheaply exploring how Polis conversations might unfold in specific situations. Following Park et al. (2022), an entire Polis conversation could be simulated using an LLM-generated collection of participants—each with unique backgrounds, expertise, and ideological alignments—who comment and vote in response to a prompt. The social dynamic of a simulation can thus be controlled to evaluate the behaviour of Polis in any number of scenarios, from extreme polarization, to high ideological uniformity.

This functionality could be used for fine-tuning new features, training facilitators, or building intuition about how conversations might unfold in particular situations, in order to formulate hypotheses that are then tested in the real world. However, we again stress that it must *not* be used to replace human deliberation or make any social inference, but rather to allow researchers to iterate and experiment more quickly and cost-effectively on deliberation technologies, with all the appropriate caveats as to how performance on a simulation carries to the real world.

It is possible to point here to a deeper and more structural problem in the innovation of technology intended for the advancement of democratic behavior: certain countries and cities are willing to serve as a test bench, but these so called "democracy labs" have limited bandwidth and political will that must be fought for. Reducing the cycle time for innovation in political technology is a laudable goal.

### 3.4.5  Reframing Comments for Specific Opinion Groups

On very divided topics of conversation, where differing ideological perspectives frustrate each side's ability to understand the other, LLMs may be able to re-frame individual comments or summaries in a way that better resonates with the other side. This generally fits within the context of conversation facilitation (see Section 2.3). Our initial experimentation with this idea produced very sub-optimal results (omitted for the sake of brevity), but it's possible that more careful prompt engineering or finely tuned models could help participants build shared understanding. However, the implications of this work should be carefully considered, in particular the potential for erasure of historically disenfranchised voices and perspectives.

### 3.4.6  Author Assistance

Even before the advent of LLMs, NLP technology has been capable of providing nuanced stylistic feedback to authors based on their own writing (e.g. grammarly.com). LLMs have taken this ability to another level, able to imitate any style it might be familiar with from its training corpus, from Shakespeare to Thoreau.

We can therefore consider how stylistic or editorial suggestions might improve deliberative applications like Polis, alongside risks of distorting the meaning of people's communications by forcing them towards imaginary, oppressive norms and erasing vernacular expression. For example, suggestions might be made based on the following:

- Whether a comment has compound ideas which would be more effective in the system broken up into separate comments (avoiding a common issue where compound ideas elicit more passes when participants agree with one point, but not others)

- Notify the participant when their ideas seem unclear or abusive, and thus at risk of being moderated out

- How to make the comment more likely to be approved of globally, and thus have a higher chance of emerging as group informed consensus

- How to make the comment more appealing towards a specific opinion group, whether the authors own, or another

- Generally offer suggestions for improving grammar or fixing typos

Both the risks and rewards associated with this class of features lies in the potential for an LLM to steer a discussion away from the specific framing that naturally comes to participants. Polis has an ethnographic disposition; the idea is to elicit people's standpoint in their own words and phrasing. The risk is that important perspectives or tones might get erased, e.g. any group being biased away from a mode of speech which most resonates with them and other members of their community. The potential benefit though is that it may be easier for people to understand each other, and come to consensus. Care must be taken in considering how to do this without nudging participants

towards an artificial consensus. Of these potential interventions, LLM-generated suggestions about how to break up compound ideas into separate comments may therefore be the safest.

### 3.4.7 Universal Translation Layer

As the technology improves, we expect that LLMs will provide a much improved universal translation layer (including nuance, domain specificity and intent) which will have profound implications for multi-language dialogue (see for example Vincent (2022)). Consider a real-time dialogue among many countries of the EU happening at a scale where human translators are not available, except to randomly spot check the LLM for quality assurance. Such a capability has the potential to address the linguistic injustice of people marginalized from political representation due to language barriers or language hegemony.

## 3.5 Key Techniques in Application of LLMs

This paper has been the result of collaborations between The Computational Democracy Project, which broadly applies computational and machine learning methods towards solving problems of participatory democracy, and Anthropic, which works to build more reliable, steerable and interpretable AI systems. Working together, we identified the following LLM usage patterns that we hope may help anyone looking to employ LLMs for deliberative technology:

- **Iterative or recursive compilation of information to get around context window limits**: This was for example key in our approach to topic-modelling, and we anticipate may be useful in aspects of summarization and content generation with LLMs constrained by context window length.

- **Using probability distributions instead of single answers**: When answering basic prompts for e.g. voting or moderation options, this technique gives a much more robust picture of how confident the LLM is about possible responses. This can help bridge LLMs with the broad existing probabilistic approach to machine learning (See e.g. the textbook trio Murphy, 2012, 2022, 2023). However, most commercial APIs, including Claude's, do not provide these probabilities at the time of writing.

- **Chain of thought reasoning:** Asking LLMs to explain themselves proved valuable for improving their results and understanding their behavior.

## 3.6 Conclusion

The advent of technologies like Polis provides a powerful platform for collective self-representation, facilitating the construction of an interactive model of the public opinion landscape.

Historically, elected representatives have served as a conduit for public opinion, "summarizing", in a manner of speaking, the voices of their constituents. This process often occurs through qualitative conversations in town halls and private meetings

with industry stakeholders. These representatives function as embodied information compression, distilling broad public sentiment into digestible messages. In parallel, research institutions like Pew Research conduct independent public opinion polling. This involves surveying a segment of the population, summarizing their responses, and thereby providing another form of representation.

In both scenarios, there exists a summarizing mechanism that mediates between a population and an end result, whether it be a politician's soundbite, an article informed by research, or, in this case a Polis report. This inherent property of representative democracy allows for the coexistence of heterogeneous, often conflicting summaries – textual intermediaries that encapsulate a diverse range of public will.

The advent of transformer-based Large Language Models (LLMs) introduces another layer to this dynamic. With their *attention* mechanisms, it may eventually be possible to use these models to formally quantify aspects of the summarization process in unprecedented ways, though currently there is debate about whether these mechanisms can be relied on for interpretability (Jain and Wallace, 2019; Wiegreffe and Pinter, 2019). However, our experiments in 2.2.2 suggest that LLMs can dramatically reduce the cost of evaluating summaries and identifying omissions, increasing accountability and transparency relative to existing summarization systems.

The capability of LLMs to create coherent summaries from a multitude of individual participant-written statements further enhances the process facilitated by Polis. This introduces a feedback loop that hasn't been available to citizens who are being "summarized" by other means. This potentially recursive self-summarization could be revolutionary, providing continuous, adaptive feedback and resulting in a more accurate representation of a population's views by the population itself.

In essence, the integration of LLMs in platforms like Polis can introduce a new era in public self-representation. It can ensure that the voices previously ignored or drowned out in the process are heard, thus making democracy more inclusive, adaptive, and reflective of its constituents.

In conclusion, the use of LLMs can significantly improve the productivity and efficiency of Polis discussions. LLMs can generate prompts to guide the conversation, identify patterns and trends in the discussion, and provide a real-time summary to the human facilitator. In doing all of this, they have the potential to significantly reduce the burden of human labor required to facilitate this kind of engagement, therefore making it possible to roll them out more frequently and in more diverse contexts. However, the details of how this is accomplished are vital, and it is imperative for safety and legitimacy that, as a baseline, there remain an oversight role for humans in the loop wherever LLMs are employed in this setting. Additional research exploring safeguards and metrics for gauging outcomes is also absolutely necessary to ensure positive outcomes.

# References

The Society Library, Mar. 2023. URL https://www.societylibrary.org.

A. Abid, M. Farooqi, and J. Zou. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES

'21, page 298–306, New York, NY, USA, 2021a. Association for Computing Machinery. ISBN 9781450384735. doi: 10.1145/3461702.3462624. URL https://doi.org/10.1145/3461702.3462624.

A. Abid, M. Farooqi, and J. Zou. Large language models associate Muslims with violence. *Nature Machine Intelligence*, 3(6):461–463, June 2021b. ISSN 2522-5839. doi: 10.1038/s42256-021-00359-2. URL https://www.nature.com/articles/s42256-021-00359-2. Number: 6 Publisher: Nature Publishing Group.

A. Applebaum and P. Pomerantsev. How to Put Out Democracy's Dumpster Fire. *The Atlantic*, Mar. 2021. URL https://www.theatlantic.com/magazine/archive/2021/04/the-internet-doesnt-have-to-be-awful/618079/. Section: Ideas.

A. Askell, Y. Bai, A. Chen, D. Drain, D. Ganguli, T. Henighan, A. Jones, N. Joseph, B. Mann, N. DasSarma, N. Elhage, Z. Hatfield-Dodds, D. Hernandez, J. Kernion, K. Ndousse, C. Olsson, D. Amodei, T. Brown, J. Clark, S. McCandlish, C. Olah, and J. Kaplan. A General Language Assistant as a Laboratory for Alignment, Dec. 2021. URL http://arxiv.org/abs/2112.00861. arXiv:2112.00861 [cs].

Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan, N. Joseph, S. Kadavath, J. Kernion, T. Conerly, S. El-Showk, N. Elhage, Z. Hatfield-Dodds, D. Hernandez, T. Hume, S. Johnston, S. Kravec, L. Lovitt, N. Nanda, C. Olsson, D. Amodei, T. Brown, J. Clark, S. Mc-Candlish, C. Olah, B. Mann, and J. Kaplan. Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback, Apr. 2022. URL http://arxiv.org/abs/2204.05862. arXiv:2204.05862 [cs].

M. A. Bakker, M. J. Chadwick, H. R. Sheahan, M. H. Tessler, L. Campbell-Gillingham, J. Balaguer, N. McAlesse, A. Glaese, J. Aslanides, M. M. Botvinick, and C. Summerfield. Fine-tuning language models to find agreement among humans with diverse preferences, Nov. 2022. URL https://arxiv.org/abs/2211.15006.

E. Barry. Townhall meeting in Kentucky turns tables on polarization, 2023. URL https://compdemocracy.org/Case-studies/2018-kentucky/.

C. Basta, M. R. Costa-jussà, and N. Casas. Evaluating the Underlying Gender Bias in Contextualized Word Embeddings, Apr. 2019. URL http://arxiv.org/abs/1904.08783. arXiv:1904.08783 [cs].

E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623, New York, NY, USA, 2021a. Association for Computing Machinery. ISBN 9781450383097. doi: 10.1145/3442188.3445922. URL https://doi.org/10.1145/3442188.3445922.

E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, pages

610–623, New York, NY, USA, Mar. 2021b. Association for Computing Machinery. ISBN 978-1-4503-8309-7. doi: 10.1145/3442188.3445922. URL https://dl.acm.org/doi/10.1145/3442188.3445922.

J. Bilich, M. Varga, D. Masood, and A. Konya. Faster Peace via Inclusivity: An Efficient Paradigm to Understand Populations in Conflict Zones. *AI for Social Good workshop at NeurIPS*, page 6, 2019.

R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, E. Brynjolfsson, S. Buch, D. Card, R. Castellon, N. Chatterji, A. Chen, K. Creel, J. Q. Davis, D. Demszky, C. Donahue, M. Doumbouya, E. Durmus, S. Ermon, J. Etchemendy, K. Ethayarajh, L. Fei-Fei, C. Finn, T. Gale, L. Gillespie, K. Goel, N. Goodman, S. Grossman, N. Guha, T. Hashimoto, P. Henderson, J. Hewitt, D. E. Ho, J. Hong, K. Hsu, J. Huang, T. Icard, S. Jain, D. Jurafsky, P. Kalluri, S. Karamcheti, G. Keeling, F. Khani, O. Khattab, P. W. Koh, M. Krass, R. Krishna, R. Kuditipudi, A. Kumar, F. Ladhak, M. Lee, T. Lee, J. Leskovec, I. Levent, X. L. Li, X. Li, T. Ma, A. Malik, C. D. Manning, S. Mirchandani, E. Mitchell, Z. Munyikwa, S. Nair, A. Narayan, D. Narayanan, B. Newman, A. Nie, J. C. Niebles, H. Nilforoshan, J. Nyarko, G. Ogut, L. Orr, I. Papadimitriou, J. S. Park, C. Piech, E. Portelance, C. Potts, A. Raghunathan, R. Reich, H. Ren, F. Rong, Y. Roohani, C. Ruiz, J. Ryan, C. Ré, D. Sadigh, S. Sagawa, K. Santhanam, A. Shih, K. Srinivasan, A. Tamkin, R. Taori, A. W. Thomas, F. Tramèr, R. E. Wang, W. Wang, B. Wu, J. Wu, Y. Wu, S. M. Xie, M. Yasunaga, J. You, M. Zaharia, M. Zhang, T. Zhang, X. Zhang, Y. Zhang, L. Zheng, K. Zhou, and P. Liang. On the opportunities and risks of foundation models, 2022.

T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, and A. Askell. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.

M. Cheng, E. Durmus, and D. Jurafsky. Marked personas: Using natural language prompts to measure stereotypes in language models, 2023.

K.-L. Chiu, A. Collins, and R. Alexander. Detecting Hate Speech with GPT-3, Mar. 2022. URL http://arxiv.org/abs/2103.12407. arXiv:2103.12407 [cs].

E. Chu, J. Andreas, S. Ansolabehere, and D. Roy. Language Models Trained on Media Diets Can Predict Public Opinion, Mar. 2023. URL http://arxiv.org/abs/2303.16779. arXiv:2303.16779 [cs].

K. Claessen and J. Hughes. QuickCheck: a lightweight tool for random testing of Haskell programs. In *Proceedings of the fifth ACM SIGPLAN international conference on Functional programming*, ICFP '00, pages 268–279, New York, NY, USA, Sept. 2000. Association for Computing Machinery. ISBN 978-1-58113-202-1. doi: 10.1145/351240.351266. URL https://dl.acm.org/doi/10.1145/351240.351266.

D. A. Cohn, Z. Ghahramani, and M. I. Jordan. Active learning with statistical models. *Journal of artificial intelligence research*, 4:129–145, 1996.

P. Coy. Opinion | Can A.I. and Democracy Fix Each Other? *The New York Times*, Apr. 2023. ISSN 0362-4331. URL https://www.nytimes.com/2023/04/05/opinion/artificial-intelligence-democracy-chatgpt.html.

CPI. Building Consensus and Compromise on Uber in Taiwan, Sept. 2019. URL https://www.centreforpublicimpact.org/case-study/building-consensus-compromise-uber-taiwan.

E. Durmus, H. He, and M. Diab. FEQA: A question answering evaluation framework for faithfulness assessment in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5055–5070, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.454. URL https://aclanthology.org/2020.acl-main.454.

E. Durmus, F. Ladhak, and T. Hashimoto. Spurious correlations in reference-free evaluation of text generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1443–1454, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.102. URL https://aclanthology.org/2022.acl-long.102.

A. R. Fabbri, W. Kryściński, B. McCann, C. Xiong, R. Socher, and D. Radev. Summeval: Re-evaluating summarization evaluation, 2021.

A. France-Presse. Romania PM unveils AI 'adviser' to tell him what people think in real time. *The Guardian*, Mar. 2023. ISSN 0261-3077. URL https://www.theguardian.com/world/2023/mar/02/romania-ion-ai-government-honorary-adviser-artificial-intelligence-pm-nicolae-ciuca.

D. Ganguli, D. Hernandez, L. Lovitt, N. DasSarma, T. Henighan, A. Jones, N. Joseph, J. Kernion, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, D. Drain, N. Elhage, S. E. Showk, S. Fort, Z. Hatfield-Dodds, S. Johnston, S. Kravec, N. Nanda, K. Ndousse, C. Olsson, D. Amodei, D. Amodei, T. Brown, J. Kaplan, S. McCandlish, C. Olah, and J. Clark. Predictability and Surprise in Large Generative Models. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1747–1764, June 2022a. doi: 10.1145/3531146.3533229. URL http://arxiv.org/abs/2202.07785. arXiv:2202.07785 [cs].

D. Ganguli, L. Lovitt, J. Kernion, A. Askell, Y. Bai, S. Kadavath, B. Mann, E. Perez, N. Schiefer, K. Ndousse, A. Jones, S. Bowman, A. Chen, T. Conerly, N. DasSarma, D. Drain, N. Elhage, S. El-Showk, S. Fort, Z. Hatfield-Dodds, T. Henighan, D. Hernandez, T. Hume, J. Jacobson, S. Johnston, S. Kravec, C. Olsson, S. Ringer, E. Tran-Johnson, D. Amodei, T. Brown, N. Joseph, S. McCandlish, C. Olah, J. Kaplan, and J. Clark. Red Teaming Language Models to Reduce Harms: Methods, Scaling Behaviors, and Lessons Learned, Nov. 2022b. URL http://arxiv.org/abs/2209.07858. arXiv:2209.07858 [cs].

Z. Ghahramani and T. L. Griffiths. Infinite latent feature models and the Indian buffet process. May 2005.

T. Goyal and G. Durrett. Annotating and modeling fine-grained factuality in summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1449–1462, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.114. URL https://aclanthology.org/2021.naacl-main.114.

K. Greshake, S. Abdelnabi, S. Mishra, C. Endres, T. Holz, and M. Fritz. Not what you've signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection, May 2023. URL http://arxiv.org/abs/2302.12173. arXiv:2302.12173 [cs].

T. Griffiths, M. Jordan, J. Tenenbaum, and D. Blei. Hierarchical topic models and the nested chinese restaurant process. *Advances in neural information processing systems*, 16, 2003.

J. Habermas. *The Structural Transformation of the Public Sphere*. Polity Press, England, 1962. ISBN 978-0-262-58108-0.

J. Habermas. *The Theory of Communicative Action*. 1981. URL https://en.wikipedia.org/w/index.php?title=The_Theory_of_Communicative_Action&oldid=1134581345. Page Version ID: 1134581345.

S. Harding. *The Science Question in Feminism*. Cornell University Press, 1986. URL https://www.cornellpress.cornell.edu/book/9780801418808/the-science-question-in-feminism/#bookTabs=1.

J. Hartmann, J. Schwenzow, and M. Witte. The political ideology of conversational AI: Converging evidence on ChatGPT's pro-environmental, left-libertarian orientation, Jan. 2023. URL https://papers.ssrn.com/abstract=4316084.

C. Horton. The simple but ingenious system Taiwan uses to crowdsource its laws. *MIT Technology Review*, Aug. 2018. URL https://www.technologyreview.com/2018/08/21/240284/the-simple-but-ingenious-system-taiwan-uses-to-crowdsource-its-laws/.

B. Hutchinson, V. Prabhakaran, E. Denton, K. Webster, Y. Zhong, and S. Denuyl. Social Biases in NLP Models as Barriers for Persons with Disabilities. pages 5491–5501. Association for Computational Linguistics, July 2020. doi: 10.18653/v1/2020.acl-main.487.

M. S. Jahan and M. Oussalah. A systematic review of Hate Speech automatic detection using Natural Language Processing, May 2021. URL http://arxiv.org/abs/2106.00742. arXiv:2106.00742 [cs].

S. Jain and B. C. Wallace. Attention is not Explanation, May 2019. URL http://arxiv.org/abs/1902.10186. arXiv:1902.10186 [cs].

Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung. Survey of hallucination in natural language generation. *ACM Comput.*

*Surv.*, 55(12), mar 2023a. ISSN 0360-0300. doi: 10.1145/3571730. URL https://doi.org/10.1145/3571730.

Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung. Survey of Hallucination in Natural Language Generation. *ACM Computing Surveys*, 55(12):248:1–248:38, Mar. 2023b. ISSN 0360-0300. doi: 10.1145/3571730. URL https://doi.org/10.1145/3571730.

A. Konya, Y. L. Qiu, M. P. Varga, and A. Ovadya. Elicitation Inference Optimization for Multi-Principal-Agent Alignment. *NuerIPS*, 2022. URL https://openreview.net/pdf?id=tkxnRPkb_H.

W. Kryscinski, B. McCann, C. Xiong, and R. Socher. Evaluating the factual consistency of abstractive text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9332–9346, Online, Nov. 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.750. URL https://aclanthology.org/2020.emnlp-main.750.

K. Kurita, N. Vyas, A. Pareek, A. W. Black, and Y. Tsvetkov. Measuring Bias in Contextualized Word Representations, June 2019. URL http://arxiv.org/abs/1906.07337. arXiv:1906.07337 [cs].

G. Lakoff. *The Political Mind: A Cognitive Scientist's Guide to Your Brain and Its Politics*. Penguin Books, New York, NY, reprint edition edition, June 2009. ISBN 978-0-14-311568-7.

G. Lakoff, H. Dean, and D. Hazen. *Don't Think of an Elephant!: Know Your Values and Frame the Debate–The Essential Guide for Progressives*. Chelsea Green Publishing, White River Junction, Vt, first edition edition, Sept. 2004. ISBN 978-1-931498-71-5.

N. Lambert, L. Castricato, L. von Werra, and A. Havrilla. Illustrating Reinforcement Learning from Human Feedback (RLHF). *Hugging Face Blog*, 2022. URL https://huggingface.co/blog/rlhf.

H. Landemore. *Open Democracy: Reinventing Popular Rule for the Twenty-First Century*. Princeton University Press, Princeton, Mar. 2022. ISBN 978-0-691-21239-5. URL https://press.princeton.edu/books/hardcover/9780691181998/open-democracy.

M. Lee, M. Srivastava, A. Hardy, J. Thickstun, E. Durmus, A. Paranjape, I. Gerard-Ursin, X. L. Li, F. Ladhak, F. Rong, R. E. Wang, M. Kwon, J. S. Park, H. Cao, T. Lee, R. Bommasani, M. Bernstein, and P. Liang. Evaluating human-language model interaction, 2022.

S. MacAvaney, H.-R. Yao, E. Yang, K. Russell, N. Goharian, and O. Frieder. Hate speech detection: Challenges and solutions. *PLOS ONE*, 14(8):e0221152, Aug. 2019. ISSN 1932-6203. doi: 10.1371/journal.pone.0221152. URL https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0221152. Publisher: Public Library of Science.

J. B. MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*, volume 5.1, pages 281–298. University of California Press, 1967. URL https://projecteuclid.org/ebooks/berkeley-symposium-on-mathematical-statistics-and-probability/Proceedings-of-the-Fifth-Berkeley-Symposium-on-Mathematical-Statistics-and/chapter/Some-methods-for-classification-and-analysis-of-multivariate-observations/bsmsp/1200512992.

J. Maynez, S. Narayan, B. Bohnet, and R. McDonald. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.173. URL https://aclanthology.org/2020.acl-main.173.

J. McKenzie. Testing Tech for Consensus in a Purple Town | Civicist, Mar. 2018. URL https://web.archive.org/web/20210414093745/https://civichall.org/civicist/testing-tech-consensus-purple-town/. Cached archive.

C. Miller and M. Anubi. How to fix democracy. URL https://www.bbc.co.uk/programmes/p0crtyww.

K. P. Murphy. *Machine learning: a probabilistic perspective*. Adaptive computation and machine learning series. MIT Press, Cambridge, MA, 2012. ISBN 9780262018029.

K. P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022. URL probml.ai.

K. P. Murphy. *Probabilistic Machine Learning: Advanced Topics*. MIT Press, 2023. URL http://probml.github.io/book2.

OpenAI. GPT-4 Technical Report, Mar. 2023. URL http://arxiv.org/abs/2303.08774. arXiv:2303.08774 [cs].

C. Pacheco and M. D. Ernst. Randoop: feedback-directed random testing for Java. In *Companion to the 22nd ACM SIGPLAN conference on Object-oriented programming systems and applications companion*, OOPSLA '07, pages 815–816, New York, NY, USA, Oct. 2007. Association for Computing Machinery. ISBN 978-1-59593-865-7. doi: 10.1145/1297846.1297902. URL https://doi.org/10.1145/1297846.1297902.

J. S. Park, L. Popowski, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein. Social simulacra: Creating populated prototypes for social computing systems. In *In the 35th Annual ACM Symposium on User Interface Software and Technology (UIST '22)*, UIST '22, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393201. doi: 10.1145/3526113.3545616. URL https://doi.org/10.1145/3526113.3545616.

K. Pearson. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901. doi: 10.1080/14786440109462720. URL https:

//doi.org/10.1080/14786440109462720. Publisher: Taylor & Francis _eprint: https://doi.org/10.1080/14786440109462720.

F. Perez and I. Ribeiro. Ignore Previous Prompt: Attack Techniques For Language Models, Nov. 2022. URL http://arxiv.org/abs/2211.09527. arXiv:2211.09527 [cs].

C. E. Rasmussen. Gaussian processes for machine learning. 2006.

P. Robert and Y. Escoufier. A Unifying Tool for Linear Multivariate Statistical Methods: The RV- Coefficient. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 25(3):257–265, 1976. ISSN 0035-9254. doi: 10.2307/2347233. URL https://www.jstor.org/stable/2347233. Publisher: [Wiley, Royal Statistical Society].

M. B. Rosenberg. *Nonviolent communication : a language of life*. Encinitas, CA : PuddleDancer Press, 2005. ISBN 978-1-892005-03-8. URL http://archive.org/details/isbn_9781892005038.

S. Santurkar, E. Durmus, F. Ladhak, C. Lee, P. Liang, and T. Hashimoto. Whose Opinions Do Language Models Reflect?, Mar. 2023. URL http://arxiv.org/abs/2303.17548. arXiv:2303.17548 [cs].

M. Sap, S. Gabriel, L. Qin, D. Jurafsky, N. A. Smith, and Y. Choi. Social Bias Frames: Reasoning about Social and Power Implications of Language, Apr. 2020. URL http://arxiv.org/abs/1911.03891. arXiv:1911.03891 [cs].

J. Schulman, B. Zoph, C. Kim, J. Hilton, J. Menick, J. Weng, J. F. Ceron Uribe, L. Fedus, L. Metz, M. Pokorny, R. Gontijo Lopes, S. Zhao, A. Vijayvergiya, E. Sigler, A. Perelman, C. Voss, M. Heaton, J. Parish, D. Cummings, R. Nayak, V. Balcom, D. Schnurr, T. Kaftan, C. Hallacy, N. Turley, N. Deutsch, V. Goel, J. Ward, A. Konstantinidis, W. Zaremba, L. Ouyang, L. Bogdonoff, J. Gross, D. Medina, S. Yoo, T. Lee, R. Lowe, D. Mossing, J. Huizinga, R. Jiang, C. Wainwright, D. Almeida, S. Lin, M. Zhang, K. Xiao, K. Slama, S. Bills, A. Gray, J. Leike, J. Pachocki, P. Tillet, S. Jain, G. Brockman, N. Ryder, A. Paino, Q. Yuan, C. Winter, B. Wang, M. Bavarian, I. Babuschkin, S. Sidor, I. Kanitscheider, M. Pavlov, M. Plappert, N. Tezak, H. Jun, W. Zhuk, V. Pong, L. Kaiser, J. Tworek, A. Carr, L. Weng, S. Agarwal, K. Cobbe, V. Kosaraju, A. Power, S. Polu, J. Han, R. Puri, S. Jain, B. Chess, C. Gibson, O. Boiko, E. Parparita, A. Tootoonchian, K. Kosic, and C. Hesse. ChatGPT: Optimizing Language Models for Dialogue, Nov. 2022. URL https://openai.com/blog/chatgpt/.

D. Sergent. First-ever civic assembly gives residents chance to be heard. *Bowling Green Daily News*, Feb. 2018. URL https://www.bgdailynews.com/news/first-ever-civic-assembly-gives-residents-chance-to-be-heard/article_0a17254e-a8bb-5f4f-884f-9d0617ab9c08.html. author email: dsergent@bgdailynews.com.

C. Singh, J. X. Morris, J. Aneja, A. M. Rush, and J. Gao. Explaining Patterns in Data with Language Models via Interpretable Autoprompting, Jan. 2023. URL http://arxiv.org/abs/2210.01848. arXiv:2210.01848 [cs, q-bio, stat].

C. Small, M. Bjorkegren, T. Erkkilä, L. Shaw, and C. Megill. Polis: Scaling Deliberation by Mapping High Dimensional Opinion Spaces. *Recerca: Revista de Pensament i Anàlisi*, 26(2), 2021. doi: https://doi.org/10.6035/recerca.5516. URL https://www.proquest.com/docview/2610037205. Publisher: Universitat Jaume I Servei de Comunicacio i Publicacions.

N. Stiennon, L. Ouyang, J. Wu, D. Ziegler, R. Lowe, C. Voss, A. Radford, D. Amodei, and P. F. Christiano. Learning to summarize with human feedback. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 3008–3021. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf.

R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press, 1998.

H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample. LLaMA: Open and Efficient Foundation Language Models, Feb. 2023. URL http://arxiv.org/abs/2302.13971. arXiv:2302.13971 [cs].

J. Vincent. Google plans giant AI language model supporting world's 1,000 most spoken languages, Nov. 2022. URL https://www.theverge.com/2022/11/2/23434360/google-1000-languages-initiative-ai-llm-research-project.

B. Wei, J. Li, A. Gupta, H. Umair, A. Vovor, and N. Durzynski. Offensive Language and Hate Speech Detection with Deep Learning and Transfer Learning, Aug. 2021. URL http://arxiv.org/abs/2108.03305. arXiv:2108.03305 [cs].

L. Weidinger, J. Mellor, M. Rauh, C. Griffin, J. Uesato, P.-S. Huang, M. Cheng, M. Glaese, B. Balle, A. Kasirzadeh, Z. Kenton, S. Brown, W. Hawkins, T. Stepleton, C. Biles, A. Birhane, J. Haas, L. Rimell, L. A. Hendricks, W. Isaac, S. Legassick, G. Irving, and I. Gabriel. Ethical and social risks of harm from Language Models, Dec. 2021. URL http://arxiv.org/abs/2112.04359. arXiv:2112.04359 [cs].

S. Wiegreffe and Y. Pinter. Attention is not not Explanation, Sept. 2019. URL http://arxiv.org/abs/1908.04626. arXiv:1908.04626 [cs].

S. Wojcik, S. Hilgard, N. Judd, D. Mocanu, S. Ragain, K. Coleman, M. B. F. Hunzaker, and J. Baxter. Birdwatch: Crowd Wisdom and Bridging Algorithms can Inform Understanding and Reduce the Spread of Misinformation, Oct. 2022. URL https://github.com/twitter/birdwatch/blob/a7b8c8a1492eb930267f84578e7ebebefe8e8aef/birdwatch_paper_2022_10_27.pdf.

Q. Yang, J. Suh, N.-C. Chen, and G. Ramos. Grounding interactive machine learning tool design in how non-experts actually build models. In *Proceedings of the 2018 Designing Interactive Systems Conference*, DIS '18, page 573–584, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450351980. doi: 10.1145/3196709.3196729. URL https://doi.org/10.1145/3196709.3196729.

A. X. Zhang, L. Verou, and D. Karger. Wikum: Bridging Discussion Forums and Wikis Using Recursive Summarization. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*, CSCW '17, pages 2082–2096, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-4335-0. doi: 10.1145/2998181.2998235. URL http://doi.acm.org/10.1145/2998181.2998235. event-place: Portland, Oregon, USA.

T. Zhang, F. Ladhak, E. Durmus, P. Liang, K. McKeown, and T. B. Hashimoto. Benchmarking Large Language Models for News Summarization, Jan. 2023. URL http://arxiv.org/abs/2301.13848. arXiv:2301.13848 [cs].

D. M. Ziegler, N. Stiennon, J. Wu, T. B. Brown, A. Radford, D. Amodei, P. Christiano, and G. Irving. Fine-Tuning Language Models from Human Preferences, Jan. 2020. URL http://arxiv.org/abs/1909.08593. arXiv:1909.08593 [cs, stat].

# A    Appendix

We provide here more details about the prompts used for the experiments of each section. The actual code has been developed using an internal API of Claude which provides more details than the external API, such as probability distribution over the output content. It is therefore somewhat difficult to share the code. However, given that the purpose of this article is more qualitative than quantitative, we hope that the limited reproducibility provided by the prompts below should be sufficient to understand the main points.

We use the models available in the public Claude API as of March 15th, 2023, with the exception of the large context window experiments, which used the more recent 100k-token context window LLM as of June 12th, 2023.

## A.1   Topic modelling

First pass topic modeling prompt (50 comments at a time due to context window limitations):

```
Prompt:
Human: I want you to do topic modeling on the given comments. Print the
detected topics line by line.

Here are the comments:
{comments}
```

```
Assistant: Sure, I can do that. Here are the main topics I detected in
these comments:
```

We then ask the LLMs to merge the topics generated in the previous step:

**Prompt:**

```
Human: I want you merge the given lists of topics into smallest set of
topics that are comprehensive.

Here are the lists of topics:
{topics}

Assistant: Sure, I can do that. Here is my attempt at merging the list of
topics you have provided:
```

The full response to this was as follows:

**Claude's output:**

```
GROUP 0: Local Government and Policy : Government accountability, waste
reduction, transparency, Planning and zoning, Permits, Laws and regula-
tions on businesses, landlords/tenants, traffic etc., Taxes and services,
Annexed areas and districts

GROUP 1: Infrastructure and Development : Traffic and transportation in-
cluding road infrastructure, traffic flow, transit, parking, sidewalks
etc., Underground utilities and power lines, Housing including afford-
able housing and housing regulations, Historic preservation, Urban sprawl,
Downtown development including revitalizing downtown, attracting busi-
nesses and tenants

GROUP 2: Public Services : Healthcare including access to healthcare, opi-
oid crisis, mental health services etc., Education including schools,
after-school programs, curriculum, funding etc., Public spaces like li-
braries, parks, recreation centers, swimming pools etc., Utilities in-
cluding internet access, electricity, cable etc., Homeless shelters and
services, Accessibility and disability access

GROUP 3: Safety, Health and Quality of Life : Law enforcement includ-
ing traffic enforcement, noise ordinances, loitering etc., Public health
including healthy lifestyle promotion, food access, pollution etc., An-
imal control including ordinances on chaining dogs, stray animals etc.,
Emergency management including weather emergency planning and shelters,
Environmental issues like recycling, sustainability, solar power etc.,
Support for marginalized groups like LGBTQ+ people, immigrants, disabled
people etc.

GROUP 4: Economy and Business : Job opportunities, workforce development
and fair wages, Support for local businesses like small businesses,
farmers, downtown businesses etc., Lack of competition in industries like
grocery stores, Legalization of marijuana
```

```
Based on these lists, here is a comprehensive list of topics:

Government and Public Policy: Local government and politics, Laws and
regulations, Taxes and services, Transparency and accountability Infras-
tructure and Development: Housing, Transportation, Utilities, Historic
preservation, Urban planning Public Services: Healthcare, Education, Pub-
lic spaces, Homeless services, Disability access Safety, Health and
Environment: Law enforcement, Public health, Pollution, Emergency man-
agement, Marginalized groups, Animal control Economy and Business: Job
opportunities, Local business support, Competition and legalization issues
```

## A.2 Summarization

Summarizing main points and consensus by providing comments and number of agree
and disagree votes overall:

---

**Prompt:**

```
Human: In each line, I provide you with comments and percentage of votes
that agreed and disagreed with them for Group 0 and Group 1. I want you
to do topic modeling on the given comments. Print the detected topics line
by line. At the end, generate an overall summary of the comments. In the
summary, make sure to include information and quantification on how much
agreement versus disagreement there was across Group 0 and Group 1 for
different topics.

Here are the comments:
{comments}

Assistant: Sure, I can do that. Here are the main topics I detected in
these comments:
```

---

Comments for the first experiment were passed in in the form

---

**Input:**

```
Comment: {comment-text} == Number of agree votes:{agree-count} == Number of
disagree votes:{disagree-count}
```

---

For the second experiment, they were passed in the form

---

**Input:**

```
Comment: {comment-text} == Group 0 percentage agreement:{g0-agree-percent}%
== Group 0 percentage disagreement:{g0-disagree-percent}% == Group 1 per-
centage agreement:{g1-agree-percent}% == Group 1 percentage disagreement:
{g1-disagree-percent}%
```

---

## A.3 Group Identity and Consensus

## A.4   Prediction

Model used for prediction: 170B parameter plain LLM, no fine-tuning or RLHF. We predict user votes using the preceding vote history, using at most the last 30 votes due to context window limitations.

An example prompt used for prediction follows - in this case with 2 votes in the history.

## A.5   Measuring bias

The following prompt was used to ask Claude to vote on comments:

```
Human: The goal of the town hall was to give members of the community
a framework in which they could articulate and share concerns, prompted
by the topic: 'What do you believe should change in Bowling Green/Warren
County in order to make it a better place to work, live, and spend time?'.

For this topic, do you agree/disagree/pass the following comment:
"{comment-text}"

Options:
(A) Pass
(B) Disagree
(C) Agree

Assistant: If had to select one of the options, my answer would be (
```

## A.6 Evaluating Summaries with LLMs

The following prompt was used to evaluate how well a summary represents a comment:

**Prompt:**

```
Human: I'm going to ask you to evaluate whether a comment by a digital
town hall participant is represented in a summary of the town hall.

Here is the summary:

<summary>
{summary}
</summary>

Here is the comment:
<comment>
{comment}
</comment>

How well is the comment represented in the summary? Would someone who had
read the summary gain new information as a result of reading the comment?
Please choose one of the following options:

(1) The comment is not represented at all; the summary ignored the comment
entirely.
(2) The summary contains some material relevant to the comment but is
missing most of the content.
(3) The summary substantially represents the comment but is still missing
something important.
(4) The summary covers most of the information in the comment, but is
missing some nuance or detail.
(5) The summary covers all the information in the comment.
```

```
Please respond with one of these options.

Assistant: All things considered, the best answer is (
```

## A.7 Long Context Window Experiments

The following prompt was used for summarization with topic modelling as an explicit prompt instruction:

---

**Prompt:**

```
Human: Perform topic modeling on the provided comments and print the de-
tected topics line-by-line. At the end, summarize the key topics discussed
and quantify the level of agreement vs disagreement between Group 0 and
Group 1. Specifically, note which topics had high agreement or disagree-
ment. Provide an overview of the difference in opinion between the two
groups based on the topic modeling analysis.

Here are the comments: Comment: {comment-text} == Group 0 percentage
agreement: {g0-agree-percent}% == Group 0 percentage disagreement: {g0-
disagree-percent}% == Group 1 percentage agreement: {g1-agree-percent}% ==
Group 1 percentage disagreement: {g1-disagree-percent}%
....... [MORE COMMENTS] .......

Assistant: Sure, I can do that. Here are the main topics I detected in
these comments:
```

---

For summarization without explicit topic modelling as an intermediary step in the prompt:

---

**Prompt:**

```
Human: Given the comments and agreement/disagreement votes for each com-
ment from two groups (Group 0 and Group 1), summarize the key points
of consensus and divergence between Group 0 and Group 1. Identify what
opinions the two groups have consensus on versus where their views differ.

Here are the comments: Comment: {comment-text} == Group 0 percentage
agreement: {g0-agree-percent}% == Group 0 percentage disagreement: {g0-
disagree-percent}% == Group 1 percentage agreement: {g1-agree-percent}% ==
Group 1 percentage disagreement: {g1-disagree-percent}%
....... [MORE COMMENTS] .......

Assistant: Sure, I can do that.
```

---