

Moderating Democratic Discourse with LLMs

Aaditya Bhatia^{1,2} and Gita Sukthankar^{2(\boxtimes)}

¹ United States Military Academy, West Point, NY 10996, USA aaaditya.bhatia@westpoint.edu
² University of Central Florida, Orlando, FL 32816, USA {aadityabhati,gita.sukthankar}@ucf.edu

Abstract. Many social media platforms promote political polarization by creating online echo chambers where people are only exposed to information confirming their beliefs. Newer systems such as Polis and Kialo aim to foster constructive conversations and teach critical reasoning skills. However, these platforms rely heavily on human moderators to manage discussions effectively. This paper examines the effectiveness of large language models (LLMs) as moderators on Polis, an opensource, real-time system designed for democratic discourse. We evaluate the F1 score of various prompting techniques at classifying five Polis datasets labeled by human moderators. Our findings indicate that LLMs are robust to different prompting strategies and produce minimal false positives. While LLMs come with certain risks, we argue that they can be valuable tools to support human moderators, enabling broader participation in democratic discourse.

Keywords: Content moderation \cdot LLMs \cdot Computational democracy

1 Introduction

Strong democracies thrive on active citizen engagement in constructive debate. *Collective intelligence*, the combined power of diverse perspectives working together, offers solutions to complex issues [9]. Public participation allows policymakers to tap into the "wisdom of the crowds" for better decision-making [10]. Polis is one of the platforms, developed by the Computational Democracy Project¹ to harness this collective wisdom. This open-source tool facilitates realtime online discussions, gathering and analyzing citizen viewpoints directly.

Public deliberation tools face hurdles in handling the sheer volume of public opinion data. This data can be messy, filled with personal views, and may lack strong evidence. Additionally, these platforms are susceptible to manipulation by those spreading misinformation. Thus, most platforms rely heavily on human moderators to ensure that political discourse remains productive, making it more difficult to scale policy discussions to include a significant percentage of the

¹ https://compdemocracy.org/.

[©] The Author(s), under exclusive license to Springer Nature Switzerland AG 2024 R. Thomson et al. (Eds.): SBP-BRiMS 2024, LNCS 14972, pp. 123–132, 2024. https://doi.org/10.1007/978-3-031-72241-7_12

citizency. Moreover, there is a hidden cost—the potential for mental health issues like PTSD and depression in moderators who are regularly exposed to nefarious content [2].

This paper evaluates the usage of LLMs for content moderation on Polis and examines the role of different prompting strategies on F1 score, false positive rate, and confidence level. Our experiments on five different Polis datasets show that content moderation is fairly robust to different prompting strategies, exhibits a low false positive rate, but is not reliable enough to entirely supplant human moderators. This is unsurprising since moderation on Polis also involves removing redundant ideas from discussions. To enhance Polis's content moderation capabilities, we advocate for a human-machine collaborative approach that integrates large language models (LLMs) into the moderation pipeline.

2 Related Work

2.1 Deliberative Democracy

constitutes discourse that focuses on evidence and reasoning, which encourages participants to reflect on various perspectives and form an informed opinion. It assumes that through rational discourse, participants can arrive at decisions that are more legitimate and informed [18]. According to Fishkin [4], the practical realization of deliberative democracy faces several challenges. While effective deliberation can be facilitated in small groups, scaling this to a larger population becomes difficult. Maintaining inclusivity and equality in discussions is especially critical to ensure that all voices are heard and considered. Most importantly, the moderation efforts needed to structure and manage complex discussions are essential to keep them focused and productive, and this time commitment grows significantly with the number of participants. Hadfi et al. [8], proposed the usage of conversational agents to promote constructive discussion in democratic forums; however their agent was proactively driving discussion towards consensus, rather than performing content moderation.

Technological Platforms. Rather than using existing social media platforms, Klein [10–12] championed the creation of technological platforms that facilitate the aggregation, organization, and analysis of collective inputs, ensuring that the deliberative process is efficient, scalable, and inclusive. The Polis platform was specifically designed to host public deliberation and has been used by several organizations and governments in Austria, Taiwan, New Zealand, and the US [16]. Discourse on Polis is structured around conversations, created by an owner who sets the discussion topic. Participants submit comments on the topic, which are then routed to other participants to vote on. The current version of Polis relies more heavily on voting rather than natural language processing and uses matrix factorization to understand group opinions based on participants' votes on comments. The authors of the Polis platform themselves highlighted some of the ways that LLMs could promote scalable deliberation and also potential risks related to LLMs [17]. Their report was cautiously optimistic about the usage of LLMs on Polis and identified summarization, topic modeling, reporting, vote prediction, and content moderation as potential applications. However, they remained concerned about the problem of hidden biases inherited from training data.

Algorithmic Content Moderation. Algorithmic content moderation includes "systems that classify user-generated content based on either matching or prediction, leading to a decision and governance outcome." [5] Grimmelman presented a taxonomy of moderation practices for building online communities, including exclusion, pricing, organization, and norm-setting [6]. Our research is an example of moderation through organization, shaping the information flow between content producers and consumers. Much of the research in this area has focused on the detection of hate speech [14], removal of copyrighted content [20], and age-appropriate content moderation [1]. However, Polis does not contain the same type of inappropriate content found on platforms like Twitter and YouTube, since there is no pathway for content monetization. Instead, it is more similar to content moderation on Reddit, where the aim is to remove material that does not adhere to specific guidelines determined by the moderator. However, some of the Polis content moderation guidelines are mainly meant to facilitate voting on policy issues, which is not a consideration for most platforms.

3 Method

3.1 Comment Moderation

For our set of experiments on comment moderation, we use the Polis moderation guidelines to design prompts for various language models.² This method aims to identify and label irrelevant and overly complex statements. Specifically, each statement undergoes individual analysis by the language model for classification purposes. The efficacy of this strategy is evaluated by comparing the outcomes of the language model's spam detection against a gold standard-moderation data previously labeled by the moderators of each Polis study dataset. This comparison seeks to ascertain the spam detection accuracy across various language models.

To manage text generation from our language models, we use the *guidance* framework originally developed by Microsoft. This library represents a unique programming paradigm that enhances control and efficiency for a language model by constraining generation through regular expressions and context-free grammar. Developers can freely add text to the context window at any point between text generations, effectively interleaving control and generation seam-lessly using traditional programming paradigms such as conditionals and loops.

The Polis project has proposed the following moderation guidelines. Each organization conducting a user study ultimately decides its policy for moderation but generally follows these guidelines.

² Our code is available at: https://github.com/aadityabhatia/polis-argmap.

- Spam: Comments devoid of relevance to the discussion.
- Duplicative: Comments restating a previously made point.
- Complex: Comments articulating multiple ideas or problems.

Our experiment considers the effect of several variables on moderation outcomes, including:

- 1. Class labels used for classification by the language model, varying across experiments between a simple set (ACCEPT, UNSURE, REJECT) and a more detailed set (ACCEPT, UNSURE, SPAM, IRRELEVANT, UNPRO-FESSIONAL, SCOPE, COMPLEX).
- 2. A set of examples provided to the language model to help it better understand the task.
- 3. Chain-of-thought reasoning techniques that ask the model to state a predecision "thought" statement and potentially enhance the model's reasoning abilities, as demonstrated by the ReAct strategy [19].
- 4. Application of additional reasoning techniques, like statement deconstruction, to dissect and understand statements in depth.

We took steps to ensure that the model is not too strict in its moderation, since there is more harm done from falsely rejecting a statement than from falsely accepting it. A statement that is allowed into discussion will still be voted upon and moderated by the community. For example, a statement that adds no value to the discussion is likely to be treated as noise and see a strong disagreement. However, a statement that is incorrectly rejected is less likely to be seen again. Moderating legitimate statements discourages community participation and damages the trust people place in the system. This is also a place where inherent biases in language models can surface. Therefore, we prefer false positives to false negatives.

We structured our methodological approach to leverage the capabilities of transformer-based language models, specifically focusing on their application in comment moderation tasks. Using the *transformers* library and specifically chat-tuned versions of Mixtral model, we defined the task using instructional prompts and user statements for the model to respond to in a chat paradigm. We carefully crafted our inputs to align with our experimental objectives by providing context, task and expected output format. The use of the *guidance* framework allowed us to retain control over most aspects of the model's output, ensuring adherence to our structured generation format.

To uncover the latent structures within the discussion statements, we employed the BERTopic library [7], leveraging an advanced pipeline based on transformer embeddings and class-based Term Frequency-Inverse Document Frequency (TF-IDF) statistics. This process starts with generating statements embeddings, reducing their dimensionality using Uniform Manifold Approxi mation and Projection (UMAP) [15], and identification of clusters using Hierar chical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN). The resulting clusters were then analyzed for their most representative keywords and assigned descriptive labels for easy identification and interpretation using an LLM. Figure 2 shows an example topic distribution projected on a 2D plane using PCA (Fig. 1).

user	List the 8 most important areas of improvements from these statements, each on a new line.
assistant	- Cooperation between WKU and Bowling Green for community involvement and internships
	- Accessible mentorship programs for young adults and teenagers
	- Community enrichment classes and activities, including woodworking, welding, gardening, and home beautification
	- More after-school activity centers for youth with extended hours of operation
	- Increased number of concerts held in Bowling Green
	- Better incorporation of university resources in community, including Kentucky Museum
	- Public adult swim facilities for year-round physical use
	- More activities and places to go for teenagers
user	AREA OF IMPROVEMENT: Cooperation between WKU and Bowling Green for community involvement and internships
assistant	PROBLEMS IDENTIFIED: Lack of cooperation and partnership between WKU and Bowling Green, limited opportunities for community involvement and internships for students
	ACTIONABLE SOLUTIONS: Strengthening mutually beneficial partnerships between WKU and Bowling Green, increasing cooperation for community involvement and internships
	ARGUMENT: Strengthening WKU-Bowling Green partnerships will create more community involvement and internship opportunities, benefiting both students and the local community.
	ARGUMENT LABEL: Boosting student opportunities and community involvement
user	AREA OF IMPROVEMENT: Community enrichment classes and activities, including woodworking, welding, gardening, and home beautification
assistant	PROBLEMS IDENTIFIED: Limited community enrichment classes and activities in Bowling Green, particularly in woodworking, welding, gardening, and home beautification
	ACTIONABLE SOLUTIONS: Establishing and promoting community enrichment classes and activities that cater to various interests and hobbies
	ARGUMENT: Introducing diverse community enrichment classes and activities will cater to various interests, enhancing the overall quality of life in Bowling Green.
	ARGUMENT LABEL: Diversifying community enrichment
user	AREA OF IMPROVEMENT: More after-school activity centers for youth with extended hours of operation
	PROBLEMS IDENTIFIED: Insufficient after-school activity centers for youth with limited hours of operation in Bowling Green
	ACTIONABLE SOLUTIONS: Developing and expanding after-school activity centers with extended hours of operation
assistant	ARGUMENT: Providing more after-school activity centers with extended hours will offer youth safe, engaging, and educational spaces beyond traditional school hours.
	ARGUMENT LABEL: Expanding youth opportunities

Fig. 1. The set of instructions and prompts used to generate insights using *Guidance* framework. Text highlighted in green was produced by LLM, while everything else was programmatically inserted into the LLM context. Guidance enables developers to intertwine text and generation, allowing precise control of language model contexts.

Datasets. Our research focuses on two sets of events from the Polis dataset: *american-assembly.bowling-green* and *scoop-hivemind*. The first one originates from a conversation facilitated by the American Assembly in Bowling Green, Kentucky. This dataset explores local and regional perspectives regarding the city's well-being and probes the narrative of a divided America. It comprises 896 comments, of which 289 were moderated out, leaving a substantial 607 comments for analysis. This dataset is particularly valuable for understanding community priorities and perceptions at a local level.

The second set, *scoop-hivemind*, stems from multiple conversations conducted by New Zealand's Public Engagement Projects (PEP) in partnership with the news outlet Scoop regarding issues of national significance. It consists of a total of 752 comments submitted by 96 people, with 294 comments moderated out, resulting in 458 accepted comments. Within this group, the *biodiversity* and *freshwater* datasets provide insights into protecting and restoring New Zealand's biodiversity and preserving freshwater resources, an area of global environmental concern that drives significant policy decisions, while *taxes* and *affordable-housing* address the socio-economic challenges faced by the public. These datasets include detailed voting data accounting for each *agree* or *disagree* vote cast by the participants.

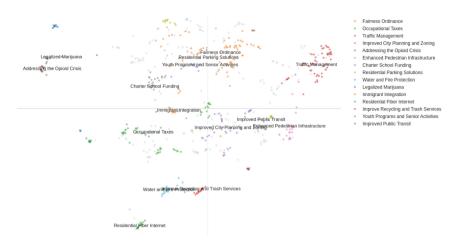


Fig. 2. Distribution of statements in the *american-assembly.bowling-green* dataset colored by topic projected in a 2D plane using PCA.

4 Experiments

Table 1 shows eight experimental configurations tested for comment moderation. For the baseline, we used a simple instructional prompt to categorize comments into three labels: ACCEPT, UNSURE, or REJECT. When the model rejected a statement, we further classified it as either SPAM or COMPLEX. Subsequently, the configurations were expanded to a more complex seven-class system, providing detailed instructions for each category to enhance the model's decisionmaking precision and explainability. The *semantic extraction* technique focused on the content of each comment to identify the problem being addressed, the proposed solution and the number of ideas introduced. We require each comment to mention at least one problem or solution and no more than one unique idea. Combining it with chain-of-thought reasoning, this approach deviated from standard moderation guidelines and instead judged comments based on their relevance to the ongoing conversation. Chain-of-thought reasoning allowed the model to articulate its thought before each decision, aiming for a higher moderation accuracy and transparency. We made extracted ideas and thought statements available as a part of our results for increased explainability.

Evaluation Metrics. To evaluate these experiments, we selected datasets with high-quality statements that align closely with Polis moderation guidelines. Raw statements were used with no pre-processing, mirroring real-world conditions where moderators make quick decisions without access to the full dataset. Effectiveness was measured using the F1 score, false positive rate, and the rate at which the model selected UNSURE over ACCEPT or REJECT. The F1 score is a balanced measure that considers both the precision and recall of the classification process, particularly useful when the costs of false positives and false

Configuration	Target Classes	Chain-of-Thought	Semantic Extraction
3	3	No	No
7	7	No	No
3T	3	Yes	No
7T	7	Yes	No
3E	3	No	Yes
7E	7	No	Yes
3ET	3	Yes	Yes
7ET	7	Yes	Yes

Table 1. Summary of experimental configurations for comment moderation. $\mathbf{E} =$ semantic extraction and $\mathbf{T} =$ chain of thought reasoning.

negatives differ significantly. It is especially useful when dealing with imbalanced datasets where positive cases, which in our case are the comments to be rejected, are significantly less common than the negative ones. It is calculated using the formula

$$F_1 = 2 \cdot \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \times TP}{2 \times TP + FP + FN}$$

Throughout these experiments, we carefully considered the implications of false positives on the moderation outcome, with a particular emphasis on minimizing false positives to foster inclusive community discussions.

5 Results

Table 2 shows the F1 score, false positive rate and number of comments classified by the LLM as uncertain across the five Polis datasets. The weighted average shows the average across all datasets, normalized by the size of the dataset. Our results show that the simple three-class baseline in which the model is simply asked to classify a statement as ACCEPT, REJECT, or UNSURE outperforms the other prompting strategies. The more sophisticated prompting strategies, chain of thought (T) and semantic extraction (E) did not consistently improve the F1 score over the baseline. The overall F1 score was not sufficiently high to make the LLM alone a convincingly good replacement for a human Polis moderator.

A second question is whether an LLM can work well in tandem with a human moderator. All the model variants exhibit low false positive rates, making them well suited for this application. Both chain of thought and semantic extraction tend to make the model more unsure, which is less problematic when the LLM is not expected to be the final arbiter. The LLM struggles with duplicate detection, which is more important for Polis since the moderators don't want to keep issuing redundant votes. Detecting duplicates requires a large context window and increases the memory footprint linearly with the number of comments, often causing the model to run out of GPU memory or overflow past its maximum context window length. A more effective approach would involve clustering comments using their text embeddings and detecting semantically identical statements, which is a promising technique for implementing a moderation system.

Kolla et al. recently published a paper on the usage of GPT-3.5 for moderating Reddit content [13]. Although they employed different prompting strategies, they reported similar performance trends in terms of false negative and true positive rates. One issue that they noted is that the LLM cannot be easily queried about its confidence level and sometimes reversed its decision upon further queries. Hence, we believe that our simpler strategy of including UNSURE as a possible class produces superior results.

Metric	Dataset	Configuration							
		3	7	3T	$7\mathrm{T}$	3E	$7\mathrm{E}$	$3 \mathrm{ET}$	7ET
$\mathrm{F1}\uparrow$	american-assembly.bowling-green	0.26	0.08	0.21	0.19	0.15	0.11	0.18	0.12
	scoop-hivemind.biodiversity	0.33	0.09	0.25	0.21	0.20	0.22	0.25	0.13
	scoop-hivemind.freshwater	0.26	0.07	0.16	0.29	0.16	0.06	0.18	0.12
	scoop-hivemind.taxes	0.10	0.04	0.06	0.07	0.10	0.19	0.26	0.17
	scoop-hivemind. affordable-housing	0.11	0.12	0.15	0.12	0.10	0.12	0.24	0.20
	Weighted Average	0.24	0.08	0.19	0.18	0.15	0.14	0.20	0.13
$\mathrm{FPR}\!\!\downarrow$	american-assembly.bowling-green	0.06	0.02	0.04	0.05	0.02	0.01	0.02	0.02
	scoop-hivemind.biodiversity	0.05	0.03	0.07	0.04	0.06	0.06	0.07	0.03
	scoop-hivemind.freshwater	0.10	0.04	0.20	0.02	0.16	0.13	0.27	0.09
	scoop-hivemind.taxes	0.02	0.01	0.16	0.04	0.12	0.04	0.10	0.05
	scoop-hivemind. affordable-housing	0.06	0.02	0.15	0.04	0.10	0.01	0.12	0.06
	Weighted Average	0.06	0.02	0.08	0.04	0.05	0.03	0.06	0.03
Unsure	american-assembly.bowling-green	0.01	0.05	0.05	0.04	0.02	0.02	0.02	0.03
	scoop-hivemind.biodiversity	0.02	0.08	0.07	0.07	0.04	0.04	0.06	0.11
	scoop-hivemind.freshwater	0.04	0.07	0.12	0.12	0.10	0.10	0.04	0.12
	scoop-hivemind.taxes	0.07	0.11	0.03	0.11	0.09	0.09	0.09	0.13
	scoop-hivemind. affordable-housing	0.04	0.05	0.05	0.09	0.02	0.08	0.01	0.10
	Weighted Average	0.02	0.06	0.06	0.06	0.03	0.04	0.03	0.07

 Table 2. Comment moderation results averaged over all datasets. Weighted average normalizes the results by the size of the dataset.

6 Conclusion

In this study, we assessed how well large language models (LLMs) perform as content moderators for Polis, a platform designed to facilitate open and democratic discussions. We tested different prompting techniques, such as those that focus on extracting meaning (semantic extraction) and revealing the reasoning process (chain-of-thought reasoning). Our findings indicate that for the Mixtral model, prompts focused on simple classification yielded the best results.

The use of artificial intelligence (AI) tools in public debates and policymaking has the potential to significantly change how we understand and address social issues. LLMs can lighten the load on human moderators by automating some content review tasks, paving the way for democratic discourse to truly flourish at scale. Our research suggests that LLMs are best suited to work alongside human moderators as part of a larger moderation system. This is because they have a low rate of incorrectly flagging comments, making it less likely that valid content will be removed unnecessarily.

Limitations There is a growing emphasis on addressing algorithmic bias and ethical considerations in these methodologies. As LLMs are trained on extensive datasets, there is a risk of inheriting biases present within the data. Contemporary methods are frequently focused on mitigating these biases to ensure that the insights generated for policymaking are equitable and representative of diverse viewpoints [8]. Another critical aspect is the scalability and computational efficiency of these approaches, particularly vital when addressing global-scale issues with large number of participants. The application of LLMs and AI-driven tools needs to be weighed against computational costs and the practicality of implementing these solutions on a large scale [3].

Acknowledgements. This material is based upon work partially supported by the Defense Advanced Research Projects Agency (DARPA) under Contract No. W911NF-20-1-0008. The conclusions and opinions expressed in this research paper are those of the authors and do not necessarily reflect the official policy or position of the U.S. Government or Department of Defense.

References

- Ahmed, S.H., Khan, M.J., Sukthankar, G.: Enhanced multimodal content moderation of children's videos using audiovisual fusion. In: The International FLAIRS Conference Proceedings, vol. 37, no. 1 (2024)
- 2. Arsht, A., Etcovitch, E.: The human cost of online content moderation. Harv. J. Law Technol. (2018)
- Ding, S., Ito, T.: Self-Agreement: A Framework for Fine-tuning Language Models to Find Agreement among Diverse Opinions (May 2023). https://doi.org/10. 48550/arXiv.2305.11460, arXiv:2305.11460 [cs]
- Fishkin, J.S.: When the People Speak: Deliberative Democracy and Public Consultation. Oxford University Press, Oxford, August 2011. https://doi.org/10.1093/ acprof:osobl/9780199604432.001.0001, https://academic.oup.com/book/12596

- Gorwa, R., Binns, R., Katzenbach, C.: Algorithmic content moderation: technical and political challenges in the automation of platform governance. Big Data Soc. 7(1), 2053951719897945 (2020)
- 6. Grimmelmann, J.: The virtues of moderation. Yale JL Tech. 17, 42 (2015)
- 7. Grootendorst, M.: BERTopic: neural topic modeling with a class-based TF-IDF procedure, March 2022. arXiv:2203.05794 [cs]
- Hadfi, R., Ito, T.: Augmented democratic deliberation: can conversational agents boost deliberation in social media? In: Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems, pp. 1794–1798. AAMAS '22, International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, May 2022
- Kirschner, P.A., Buckingham Shum, S.J., Carr, C.S., Diaper, D., Sanger, C. (eds.): Visualizing Argumentation: Software Tools for Collaborative and Educational Sense-Making. Computer Supported Cooperative Work, Springer, London (2003). https://doi.org/10.1007/978-1-4471-0037-9
- Klein, M.: How to Harvest Collective Wisdom on Complex Problems: An Introduction to the MIT Deliberatorium (2012). https://www.semanticscholar. org/paper/How-to-Harvest-Collective-Wisdom-on-Complex-An-to-Klein/ 79541bad78d11f448f94c2d3fd83d304dcf2f733
- 11. Klein, M.: Crowd-Scale Deliberation For Complex Problems: A Progress Report, March 2022. https://doi.org/10.2139/ssrn.4049680, https://papers.ssrn.com/abstract=4049680
- Klein, M., Iandoli, L.: Supporting Collaborative Deliberation Using a Large-Scale Argumentation System: The MIT Collaboratorium, February 2008. https://doi. org/10.2139/ssrn.1099082, https://papers.ssrn.com/abstract=1099082
- Kolla, M., Salunkhe, S., Chandrasekharan, E., Saha, K.: LLM-Mod: can large language models assist content moderation? In: Extended Abstracts of the CHI Conference on Human Factors in Computing Systems, pp. 1–8 (2024)
- MacAvaney, S., Yao, H.R., Yang, E., Russell, K., Goharian, N., Frieder, O.: Hate speech detection: challenges and solutions. PLoS ONE 14(8), e0221152 (2019)
- McInnes, L., Healy, J., Melville, J.: UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction, September 2020. https://doi.org/10.48550/ arXiv.1802.03426, arXiv:1802.03426 [cs, stat]
- 16. Small, C.: Polis: Scaling Deliberation by Mapping High Dimensional Opinion Spaces. RECERCA. Revista de Pensament i Anà lisi, July 2021. https://doi.org/10.6035/recerca.5516, https://www.e-revistes.uji.es/index. php/recerca/article/view/5516
- Small, C.T., et al.: Opportunities and Risks of LLMs for Scalable Deliberation with Polis, June 2023. arXiv:2306.11932 [cs]
- Williams, A.E.: Are wicked problems a lack of general collective intelligence? AI Soc. 38(1), 343–348 (2023). https://doi.org/10.1007/s00146-021-01297-8
- Yao, S., et al.: ReAct: synergizing reasoning and acting in language models, March 2023. arXiv:2210.03629 [cs]
- Zhang, D.Y., Badilla, J., Tong, H., Wang, D.: An end-to-end scalable copyright detection system for online video sharing platforms. In: 2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), pp. 626–629 (2018)