

# The Loud Minority: How a Few Frequent Commenters Shape Digital Discourse with Hostility

Anonymous ACL submission

## Abstract

Digital platforms were expected to foster broad participation in public discourse, yet online engagement remains highly unequal and underexplored. This study examines the digital participation divide and its link to hostile engagement in news comment sections. Analyzing 260 million comments from 6.2 million users over 13 years on *Naver News*, South Korea's largest news aggregation platform, we quantify participation inequality using the Gini and Palma indexes and estimate hostility levels with a BERT-based deep learning model. The findings reveal a highly skewed participation structure, with a small group of frequent users dominating discussions, particularly in Politics and Society and widely read stories. Participation inequality spikes during presidential elections, and frequent commenters are significantly more likely to post hostile content, suggesting that a vocal, and often hostile, minority disproportionately shapes digital discourse. By leveraging individual-level digital trace data, this study provides empirical insights into the behavioral dynamics of online participation inequality and its broader implications for digital public discourse.

## 1 Introduction

Digital platforms were once expected to foster broad and equitable participation in public discourse (Papacharissi, 2004). However, growing evidence suggests that online engagement remains highly unequal, with a small fraction of users dominating digital conversations, potentially skewing public discourse (e.g., Van Mierlo, 2014; Gasparini et al., 2020; Carron-Arthur et al., 2014; Baqir et al., 2023; Antelmi et al., 2019). The '90-9-1' principle, although not rigorously tested, suggests a significant disparity in online participation, where 90% of users ('lurkers') primarily observe without participating, 9% ('contributors') engage occasionally, and a mere 1% ('superusers') generate the majority

of online content (Nielsen, 2006).

This study examines the digital participation divide and its relationship with hostile engagement in online news discussions. Using a 13-year dataset from *Naver News*, South Korea's largest news aggregation platform, we analyze 260 million comments from 6.2 million users to assess the participation inequality between frequent and infrequent commenters in news comment sections and its connection with content hostility. We employ the Gini and Palma indexes to quantify participation disparities and apply a BERT-based deep learning model to classify comment hostility levels.

The findings reveal a highly unequal participation structure, with a small number of frequent users contributing disproportionately to news comment sections. This participation divide is particularly pronounced in political news domains and in a more widely read news stories. Notably, participation inequality spikes during presidential elections, suggesting that major political events exacerbate engagement disparities. Moreover, these frequent commenters are significantly more likely to post hostile content, including both uncivil and hateful content, indicating that digital discourse is shaped by a vocal, and often hostile, minority.

By leveraging individual-level digital trace data, which offers a rare opportunity to observe engagement disparities at a granular level, this study provides empirical insights into the behavioral mechanisms underlying digital discourse inequalities and their broader implications for online public discourse and public opinion formation.

## 2 Digital divide and Online Hostility

Research on digital participation has long documented significant disparities across online platforms. Contrary to early expectations that digital spaces would foster widespread civic participation (Papacharissi, 2004), the "90-9-1" principle sug-

gests that 90 percent of users passively consume content, 9 percent contribute occasionally, and only 1 percent generate the majority of online content (Nielsen, 2006). Although comprehensive research on this inequality remains scarce, several studies confirm that only a small fraction of users actively participate in digital spaces (e.g., Van Mierlo, 2014; Gasparini et al., 2020; Carron-Arthur et al., 2014; Baqir et al., 2023; Antelmi et al., 2019).

The inequality of digital participation nevertheless remains largely unexplored. Most studies on the digital divide have focused on disparities in physical access to digital systems (Chaqfeh et al., 2023) or differences in digital skills and literacy (Hargittai, 2018; Hargittai and Shaw, 2015), with far less attention given to other dimensions of digital inequality (Korovkin et al., 2023; Scheerder et al., 2017; Van Dijk, 2006). Thus, there is limited understanding of the extent of participation inequality among individuals who have access to digital platforms but engage with them to varying degrees.

Prior research also suggests that digital participation inequality may be linked to a higher likelihood of hostile engagement. Hostility or incivility in online spaces have been widely documented, particularly in political discussions and news comment sections (e.g., Coe et al., 2014; Humprecht et al., 2020; Rowe, 2015; Santana, 2014; Rossini, 2022). In online comment sections, frequent users are more likely to post hostile content. For example, research on Facebook found that highly engaged users exhibit greater levels of toxicity in their comments (Kim et al., 2021a). Similarly, studies on news comment sections indicate that hostility tends to cluster among the most active participants (Humprecht et al., 2020; Rowe, 2015), potentially shaping broader public perceptions of digital discourse. The potential association between frequent commenting and hostile content may be driven by anger, a high-arousal emotion that is strongly linked to greater engagement and participation (Berger, 2011; Brady et al., 2017; Crockett, 2017; Hasell and Weeks, 2016; Masullo et al., 2021; Valentino et al., 2011). This pattern is particularly pronounced in partisan digital environments, where hostility toward out-groups generates higher engagement than in-group favoritism (Rathje et al., 2021; Yu et al., 2024). Masullo et al. (2021) further suggests that anger increases the likelihood of users actively expressing their opinions online, regardless of the opinion climate they encounter.

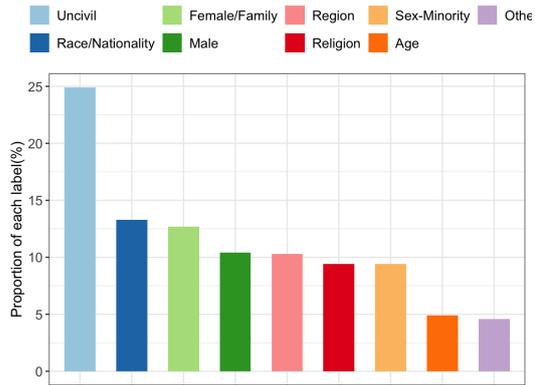


Figure 1: Distribution of Hateful and Uncivil Sentences. ‘Civil’ sentences are excluded in this figure. We allow for overlapping counts here. If a sentence has two labels, it will be counted once for each label

Building on these insights, this study advances research on the digital divide by bridging two critical aspects of online engagement—digital participation inequality and online hostility—that have not been systematically examined together. By leveraging individual-level news comment behavior data over a 13-year period, this study provides a rare opportunity to examine both the severity of the participation divide between frequent and infrequent users and whether this divide is indeed linked to hostile engagement.

### 3 Data

#### *Naver News*

South Korea is one of the most digitally connected countries in the world, boasting the highest percentage of high-speed broadband connections among OECD nations (Pak et al., 2021). In addition, in this country, online news consumption is overwhelmingly concentrated on news aggregator platforms rather than individual news websites. According to a global comparison of 46 countries, South Korea had the highest rate of news consumption via news aggregators and the lowest rate via direct access to news websites in 2021 (Oh et al., 2021). Among these platforms, *Naver News* stands as the most dominant, reflecting its unparalleled role in shaping the country’s digital news ecosystem. Over 90 percent of Koreans use *Naver* as their primary search engine, and 87 percent rely on *Naver News* for their online news consumption (Kim et al., 2021b). This shows the inequality of digital access is at least very little at play.

This minimal digital access inequality ensures

167 that disparities in online engagement are not driven  
168 by differences in basic access to digital infrastruc-  
169 ture but rather by individual preferences and be-  
170 havioral choices. Unlike in countries where digital  
171 divides are primarily shaped by disparities in inter-  
172 net access, South Korea presents a unique context  
173 where virtually all users have the opportunity to en-  
174 gage with news content online, allowing for a more  
175 precise examination of participation inequality in  
176 digital discourse.

177 The platform, *Naver News* offers users free ac-  
178 cess to news content from major news outlets in the  
179 country. A key feature of the platform is its in-link  
180 system, which enables users to read full articles  
181 and comment on them directly within *Naver*, rather  
182 than being redirected to the original news websites.  
183 This design eliminates the need for users to create  
184 accounts on multiple media sites, effectively cen-  
185 tralizing news consumption and discussion within  
186 a single platform.

187 The comprehensive scope of *Naver News* and  
188 its centralized commenting system make its data  
189 particularly valuable for studying digital participa-  
190 tion and hostile engagement at the individual level.  
191 Because South Korea has minimal barriers to inter-  
192 net access, participation disparities on the platform  
193 likely reflect user preferences rather than structural  
194 access limitations. Moreover, *Naver News* data al-  
195 lows for tracking individual commenting behavior  
196 over time, providing a rare opportunity to examine  
197 participation patterns based on frequency of use.

## 198 News Comment Data

199 From Naver News, we collected approximately 260  
200 million comments along with unique user identi-  
201 fiers from January 2008 to September 2020. Dur-  
202 ing this period, *Naver News* published a daily  
203 list of the top 30 most-read articles ("Ranking  
204 News") across six news domains: Politics, Soci-  
205 ety, Economy, World, IT/Science, and Life/Culture,  
206 totaling 180 articles per day. The dataset com-  
207 prises 802,946 articles from 141 news outlets,  
208 with 260,203,552 comments posted by approxi-  
209 mately 6,170,121 unique users. On average, each  
210 article received 324 comments.

## 211 Hate Speech Data

212 To classify hostility in news comments, we trained  
213 a BERT based deep-learning model using the *Ko-*  
214 *rean Unsmile Dataset*, a hate speech dataset pro-  
215 vided by Smilegate-AI (Kim, 2022). The dataset  
216 defines hateful expressions as those involving hos-

217 tile speech, ridicule, caricature, or prejudice against  
218 specific social groups, including explicit references,  
219 stereotype reinforcement, or conventional assump-  
220 tions about targeted groups.

221 Each comment is assigned multiple labels from  
222 ten categories, making the dataset multi-class and  
223 multi-labeled. Categories include *Civil* (devoid  
224 of hate speech), *Uncivil* (disparaging language or  
225 personal attacks), and various hate speech types  
226 targeting race/nationality, region, gender, religion,  
227 age, and sexual minorities.

228 One limitation of this dataset is the potential mis-  
229 classification of neutral comments as hateful. For  
230 example, a benign statement referencing a group  
231 may be incorrectly flagged as hate speech. To miti-  
232 gate this issue, we supplemented the dataset with  
233 additional neutral sentences following Kang et al.  
234 (2022).

235 In the training dataset, uncivil content is the most  
236 frequent category (24.5%), followed by hateful con-  
237 tent targeting race/nationality (13%), female/family  
238 (12%), male (11%), region (10%), religion (9%),  
239 sex minority (9%), and age (4.8%). Figure 1 illus-  
240 trates the label distribution.

## 241 4 Methods

### 242 Measuring Participation Inequality

243 To assess user engagement levels, we first ranked  
244 all users in the dataset based on the number of  
245 comments they posted, with the most active com-  
246 menters placed at the top. This ranking allowed us  
247 to classify users into different engagement groups,  
248 which were then used to compare hostility levels  
249 in their comments. Our analysis primarily focuses  
250 on the top 10% of the most active commenters,  
251 comparing them to the bottom 40% of commenters,  
252 who exhibit significantly lower engagement.

253 To quantify participation inequality among these  
254 user groups, we employed two widely used eco-  
255 nomic disparity metrics: the Gini index and the  
256 Palma index (Atkinson et al., 1970; Kakwani,  
257 1977), both of which have been applied in prior re-  
258 search to assess engagement inequalities in digital  
259 spaces (Glenski et al., 2020).

260 The Gini index measures the overall dispersion  
261 of participation levels, reflecting how unequally  
262 comments are distributed among users. A higher  
263 Gini index indicates greater inequality in engage-  
264 ment. However, the Gini index has notable lim-  
265 itations in interpretation. Two distributions with  
266 identical Gini values can have different underlying

267	structures, making it difficult to capture whether	<b>Measuring Comment Hostility</b>	313
268	disparities are driven by the most or least active		
269	users. Additionally, the Gini index is more sensi-	To assess levels of comment hostility, we conducted	314
270	tive to changes in the middle of the distribution but	a content analysis of comments from both heavy	315
271	less responsive to variations at the top and bottom.	(top 10%) and light (bottom 40%) commenters, as	316
272		defined by the Palma index. Within the top 10%	317
273	To address these limitations, we incorporate	group, we further distinguished the extreme top	318
274	the Palma index, which specifically measures the	1% from the remaining users, as a small subset of	319
275	ratio of participation between the top 10% of com-	commenters appeared significantly more frequently	320
276	menters and the bottom 40%. An increasing Palma	than others.	321
277	index indicates that the most active users are gain-		
278	ing even greater dominance over the least active	As an initial step, we trained <i>KC-BERT</i> , a BERT-	322
279	users, highlighting the skewed nature of digital par-	-based deep-learning model (Lee, 2020), using	323
280	ticipation. Unlike the Gini index, the Palma index	the hate comment data described earlier. Following	324
281	provides a clearer interpretation of who dominates	model training, we selected the best-performing	325
282	the discourse in digital spaces and to what extent.	version and applied it to a 1% stratified sample	326
283		of comments from each engagement group. The	327
284	We applied these two metrics across different	re-trained model assigned a hostility score to each	328
285	time periods, news domains, and news popularity	comment, and for simplicity, we categorized each	329
	rankings, depending on the specific analytical focus	comment based on its highest-scoring label while	330
	of each part of the study.	discarding those with all label scores below 0.5.	331
286			
	<b>Measuring Contribution to Inequality</b>	To facilitate analysis, we collapsed the ten pre-	332
287	After calculating the inequality metrics, we as-	defined hate speech categories into three broader	333
288	sess whether the observed disparities are primarily	classifications: <i>civil</i> , <i>uncivil</i> , and <i>hateful</i> . Specifi-	334
289	driven by frequent or infrequent commenters using	cally, comments containing general profanity and	335
290	the relative mean deviation (RMD). This metric is	personal attacks were classified as <i>uncivil</i> , while	336
291	mathematically defined as follows:	those with derogatory language targeting specific	337
		groups (e.g., race, gender, religion, region, etc.)	338
292		were categorized as <i>hateful</i> . Comments devoid	339
	$RMD_{ig} = \frac{N_i - \mu_g}{\mu_g} \quad (1)$	of such language were considered as <i>civil</i> . We	340
293		then compared the distribution of comment types	341
294	where $i$ represents an individual user, $g$ denotes	across user engagement groups, employing a chi-	342
295	the news domain. $N_i$ is the number of comments	squared proportion test to determine whether differ-	343
296	posted by user $i$ , and $\mu_g$ represents the average	ences in hostility levels between user groups were	344
297	number of comments per user in news domain $g$ .	statistically significant.	345
298			
299	The RMD serves as a counterfactual measure to	<b>5 Participation Inequality</b>	346
300	evaluate participation inequality. In a scenario		
301	where all users contributed an equal number of	Descriptive statistics on participation levels indi-	347
302	comments, the comment space would exhibit per-	cate a stark digital participation gap (Figure 3). On	348
303	fectly equal participation. This hypothetical equal	average, the top 10% of frequent commenters ac-	349
304	participation level is represented by $\mu_g$ . By com-	count for nearly half of all comments in news com-	350
305	paring each user’s actual comment count to $\mu_g$ , the	ment sections (50.11%), while the least active half	351
306	RMD quantifies how much more or less each user	(bottom 50%) contributes only 14.99% of total	352
307	contributes relative to this counterfactual equality.	comments over the years. The figure clearly illus-	353
308		trates a consistent and substantial divide in digital	354
309	This metric allows us to determine whether in-	participation, where a small subset of users dis-	355
310	equality is driven by frequent commenters posting	proportionately dominates the conversation. This	356
311	significantly more than expected or by infrequent	imbalance underscores the motivation for our study,	357
312	commenters contributing far less than the counter-	highlighting the need to investigate the structural	358
	factual amount. In doing so, it provides a clearer	disparities in online engagement.	359
	picture of how participation disparities emerge in		
	online discussions.		



Figure 2: Share of Comments by Top 10% and Bottom 40% Group

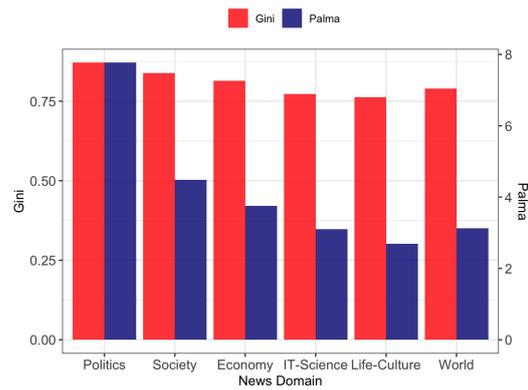


Figure 3: The Gini and Palma index Over Time by News Domain

### Participation Inequality by News Domain and Popularity

To further examine this divide, we quantified participation inequality within the news ecosystem using the Gini index and the Palma index. We then compared participation inequality (a) across six news domains (*Politics, Society, Economy, World, IT/Science, and Life/Culture*) and (b) at varying levels of news popularity. Note that *Naver News* publishes a daily list of the 30 most-read articles, referred to as ‘*Ranking News*.’ To measure news popularity, we used these rankings, with 1st representing the least popular and 30th the most popular article of the day. We then calculated Gini and Palma indexes for different news stories based on their popularity ranks to assess how inequality changes across news interest levels.

Figure 3 illustrates participation inequality across different news domains, showing that political news exhibits the highest levels of inequality compared to other categories. Both Gini and Palma indexes reveal that Politics consistently stands out as the most unequal domain, indicating that discussions in political news sections are dominated by a small subset of highly active commenters. Society and Economy also exhibit relatively high participation inequality, though to a lesser extent than Politics. In contrast, domains such as Life/Culture and IT/Science display lower levels of inequality, suggesting that discussions in these categories are more evenly distributed among users.

Figure 4 presents participation inequality as measured by the Palma index (Panel A) and the Gini index (Panel B) across different levels of news popularity. Across all domains, both indexes show a clear upward trend, indicating that as a news story becomes more popular, participation inequality increases. This pattern suggests that highly popular

articles tend to be dominated by a small group of frequent commenters, while less popular articles see a more balanced distribution of participation. Among the different news domains, Politics and Society, again, consistently exhibit the highest levels of inequality across all levels of popularity, reinforcing the idea that digital participation gaps are most pronounced in politically charged discussions.

Taken together, these findings suggest that participation inequality is not only domain-specific but also influenced by news popularity. The more widely read an article is, the more concentrated the conversation becomes among a small subset of highly active users, particularly in Politics and Society.

### User Contribution to Participation Inequality

To assess which user groups contribute most to participation inequality, we analyzed Relative Mean Deviation (RMD) scores. While the Palma and Gini indices measure overall inequality, they do not reveal how different user groups contribute to these disparities. RMD addresses this gap by indicating how much each group’s participation deviates from a hypothetical benchmark of perfect equality, where all users contribute an equal number of comments within a given news domain and news popularity level. A value of 0 represents perfect equality, while negative values indicate lower-than-expected participation, and positive values indicate excessive participation relative to the equality benchmark.

Figure 5 presents RMD scores across different user groups, segmented into ten participation levels to capture finer distinctions beyond the broad bottom 40% and top 10% classifications. The figure shows that the least active commenter groups

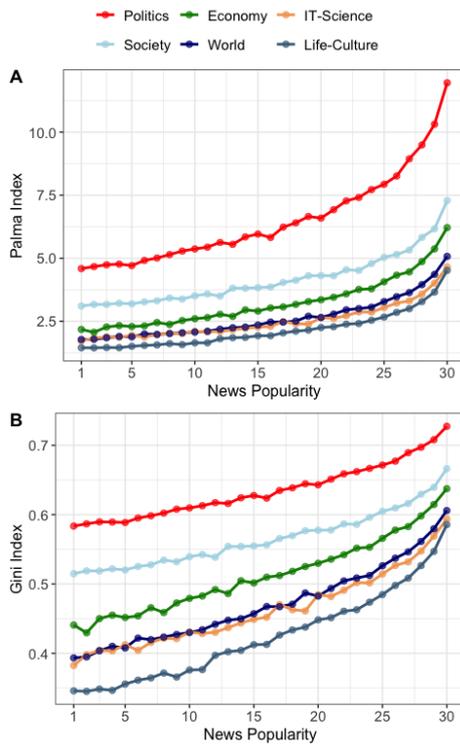


Figure 4: The Palma(Panel A) and Gini(Panel B) index by News Popularity

(Bottom 10% to Top 30-20%) cluster around zero, indicating that their participation closely aligns with the expected equal participation benchmark. In contrast, there is a progressive and disproportionate increase in deviation among more active users, with the top 1% of commenters exhibiting the highest deviation. The top 1% of users have an RMD between 23 and 30, compared to an average deviation of 3 among other active groups, demonstrating their outsized influence on digital discourse.

These findings underscore two key aspects of participation inequality. First, they indicate that the observed participation gap is primarily driven by highly active users posting disproportionately more comments, rather than infrequent users posting significantly fewer comments. This suggests that participation inequality is a function of over-contribution by a small subset of users rather than disengagement by the majority. Second, there is a sharp divide even among active commenters, particularly between the top 1% and the rest, highlighting that the most extreme contributors play a dominant role in shaping discussions. This suggests that online discourse is not only concentrated among a small subset of users but is further skewed by an even smaller group of hyper-active commenters,

reinforcing the severe imbalances in digital participation.

### Participation Inequality and Political Events

Beyond these structural patterns, we now examine how participation inequality fluctuates in response to major political events, particularly during South Korea's electoral cycles and one of the most significant political events of the study period—the 2017 impeachment of President Park Geun-hye.

Figure 6 illustrates the Gini and Palma indices in the weeks leading up to three key political events: the 2012 and 2017 presidential elections and the 2016 impeachment of the president. The trends suggest that participation inequality intensifies as major political events approach, with both indices showing a marked increase in the final weeks leading up to each event. This pattern indicates that a small subset of highly active users becomes even more dominant in news comment sections during politically charged periods, further exacerbating the imbalance in online discourse. These findings suggest that political events act as catalysts for deepening participation inequality, amplifying the influence of highly engaged users while sidelining less active participants.

## 6 Comment Hostility

Previous studies suggest that more active users in comment sections are more likely to exhibit hostility. To examine this, we conducted a computational content analysis to assess the levels of hostility in comments posted by different user groups.

For this analysis, we focused on three distinct commenter groups, ranked by their commenting activity: (1) the top 1% most active commenters, (2) the next most active group (top 10% - < top 1%), and (3) the bottom 40% least active commenters. It is important to note that the top 1% and top 10% - < top 1% are distinct groups, unlike the broader categories used in prior analyses. Given the unique behavior of the most active users, as shown in the participation inequality results, we isolated the top 1% separately to better capture the extreme engagement patterns of this highly active subset. For each group, we randomly selected 1% of comments from the raw dataset for analysis. These comments were then classified as either (1) civil, (2) uncivil, or (3) one of eight types of hateful comments using a deep learning classifier trained on a large dataset of labeled comments.

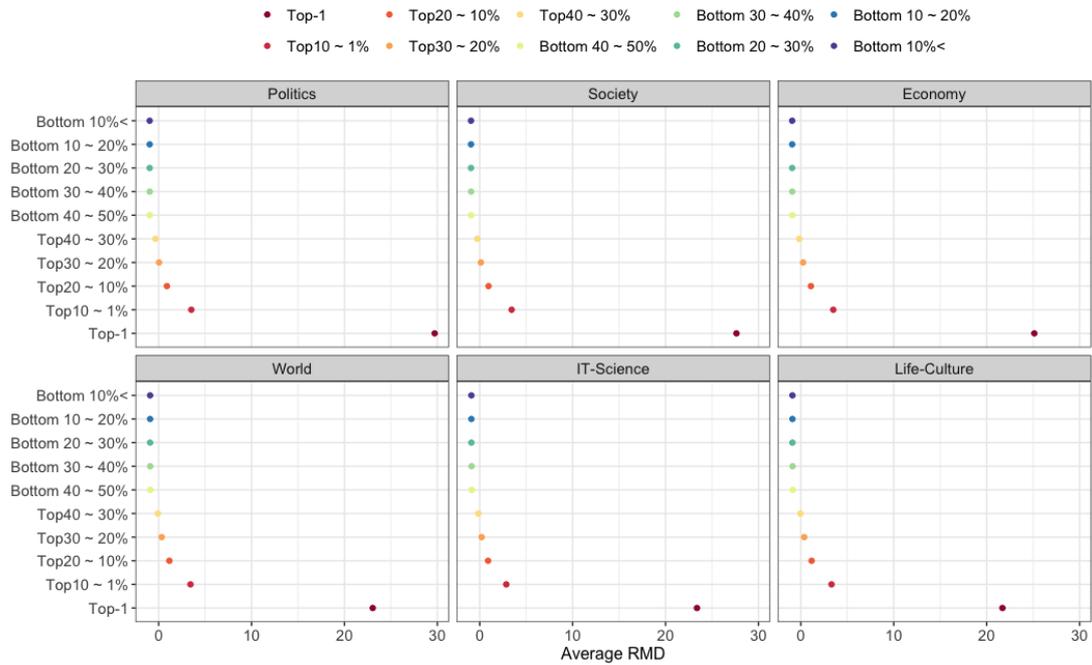


Figure 5: Average Relative Mean Deviation by News Domain

Figure 7 presents the distribution of comment categories across these three user groups. As expected, the most frequent commenters—the top 1% and top 10% - < top 1%—are significantly more likely to post *uncivil* comments compared to the less active bottom 40% (chi-square: 85.761,  $p < 0.001$  for the comparison between bottom 40% and top 10% - < top 1%, and chi-square: 71.764,  $p < 0.001$  for the comparison between bottom 40% and top 1%).

Regarding *hateful* content, the divide in online hostility extends even among active users: the top 1% is significantly more likely to post hateful comments than the top 10% - < top 1% (chi-square: 139.19,  $p < 0.001$ ). This finding further reinforces the digital participation divide, showing that not only do a small number of users dominate discussions, but they also tend to engage in higher levels of incivility and hate speech.

The disparity in hostility between active and inactive groups is still evident when examining differences across news domains. As shown in Figure 8, the gaps in both uncivil and hateful comment proportions are particularly pronounced in the Politics domain, suggesting that highly engaged users are especially likely to contribute hostile discourse in political discussions.

## 7 Conclusion

This study underscores the stark participation inequality in online news comment sections, where a small but highly active subset of users disproportionately shapes digital discourse. Analyzing 260 million comments over 13 years on Naver News, we find that this participation gap is particularly pronounced in political news discussions and highly popular news stories, intensifying during major political events such as presidential elections. The analysis also reveals that the most active commenters contribute disproportionately to the overall volume of engagement, further amplifying their influence. Moreover, these frequent commenters are significantly more likely to engage in hostile discourse, posting both uncivil and hateful content at higher rates than less active users. This suggests that online discussions are not only dominated by a small fraction of users but are also skewed toward a more hostile or hateful discourse.

These findings carry important implications for digital public discourse and online platform governance. The dominance of a small, often hostile group in comment sections raises concerns about the representativeness of online discussions and their potential to skew public perceptions. Platforms aiming to foster healthier discourse may need to consider interventions that encourage broader participation while mitigating the outsized influ-

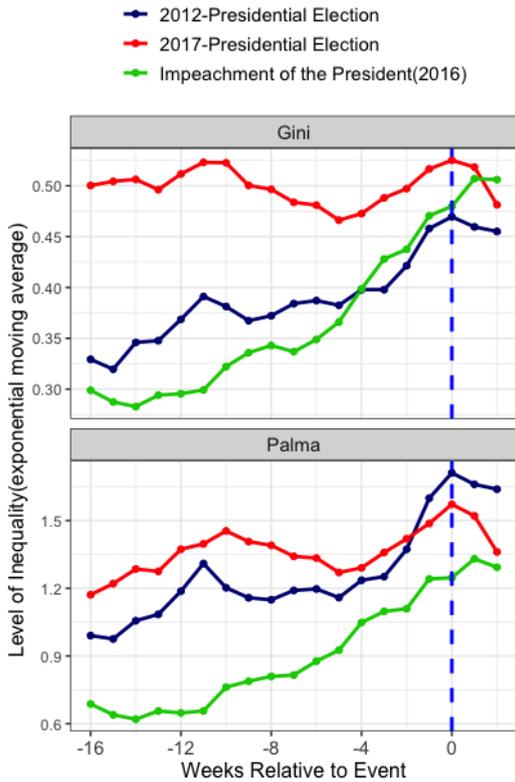


Figure 6: Participation Inequality Leading Up to Presidential Elections and the 2017 Impeachment

ence of highly engaged yet hostile users. Future research should further explore the causal mechanisms behind these dynamics and investigate potential strategies to counteract digital participation disparities and online hostility.

## 8 Limitations

While this study provides valuable insights into digital participation inequality and hostile discourse, it has several limitations that should be addressed in future research.

First, although our findings reveal a significant disparity in hostility between active and inactive user groups, further analysis is needed to understand the underlying linguistic mechanisms driving this disparity. Specifically, a more granular examination of how hostile language is constructed and varies between these groups would provide deeper insights. However, this presents a methodological challenge due to the complex structure of the Korean language. Korean allows for the creation of new words through character combinations, often leading to non-standard lexical variations in online discussions. This makes tokenization particularly difficult, as conventional NLP methods may fail

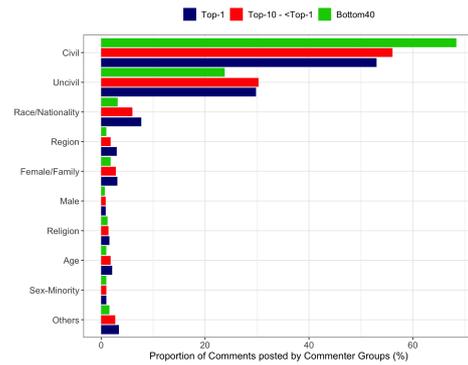


Figure 7: Hate Comment Classification Result by Percentile User Group

to capture these variations accurately. Additionally, detecting hostility—especially hateful content targeting specific sociopolitical groups—is further complicated by implicit and coded expressions that may not contain overt hate speech terms but still convey derogatory or exclusionary meanings. This linguistic flexibility enables users to mask hostility, making deep-learning-based classification models prone to under-detection of such content. Addressing this issue requires more sophisticated linguistic processing techniques, such as context-aware tokenization models, morphological analysis tailored to Korean online discourse, and adversarial training methods that can better capture implicit hostility. Future research should refine these approaches to improve the precision of hostility detection, particularly for nuanced forms of incivility and hate speech.

Second, our study does not establish a direct causal relationship between participation inequality and online hostility. While our findings suggest that hostility is more prevalent among highly active users, we have not explicitly tested whether increasing inequality drives greater hostility or if other factors mediate this relationship. As participation inequality intensifies—especially during politically charged periods—aggressive discourse may become more concentrated among dominant commenters. However, our dataset is limited to observational digital trace data, which primarily captures user behaviors, comment timing, and content but does not account for underlying psychological or social motivations. Future research should explore experimental methods to better understand the causal links between participation inequality and online hostility.

Despite these limitations, this study offers a found-

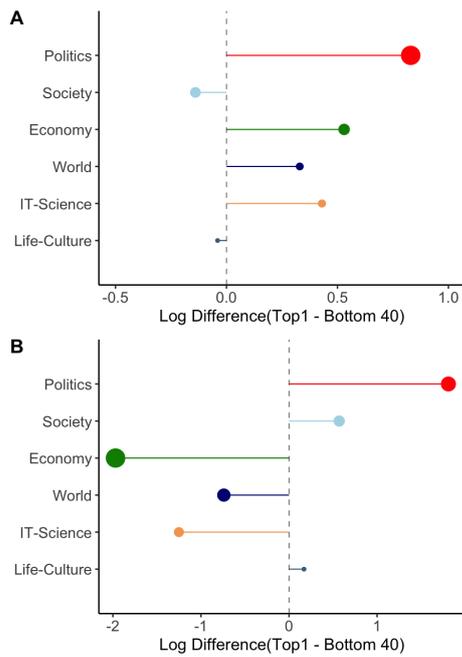


Figure 8: Log Difference in the Proportion of Uncivil (Panel A) and Hateful Comments (Panel B) between Extreme (Top 1%) and Inactive (Bottom 40%) User Group Across News Domains. Point sizes indicate the absolute difference in proportion.

dational analysis of how a vocal minority shapes digital discourse through both disproportionate engagement and increased hostility. Addressing these challenges in future research will be crucial for developing more effective moderation strategies and fostering healthier online discussions.

## References

Alessia Antelmi, Delfina Malandrino, and Vittorio Scarano. 2019. *Characterizing the Behavioral Evolution of Twitter Users and The Truth Behind the 90-9-1 Rule*. In *Companion Proceedings of The 2019 World Wide Web Conference*, pages 1035–1038, San Francisco USA. ACM.

Anthony B Atkinson et al. 1970. On the measurement of inequality. *Journal of economic theory*, 2(3):244–263.

Anees Baqir, Yijing Chen, Fernando Diaz-Diaz, Serkan Kiyak, Thomas Louf, Virginia Morini, Valentina Pansanella, Maddalena Torricelli, and Alessandro Galeazzi. 2023. *Beyond Active Engagement: The Significance of Lurkers in a Polarized Twitter Debate*. Preprint, arXiv:2306.17538. [physics].

Jonah Berger. 2011. Arousal increases social transmission of information. *Psychological science*, 22(7):891–893.

William J Brady, Julian A Wills, John T Jost, Joshua A Tucker, and Jay J Van Bavel. 2017. Emotion shapes the diffusion of moralized content in social networks. *Proceedings of the National Academy of Sciences*, 114(28):7313–7318.

Bradley Carron-Arthur, John A. Cunningham, and Kathleen M. Griffiths. 2014. *Describing the distribution of engagement in an Internet support group by post frequency: A comparison of the 90-9-1 Principle and Zipf’s Law*. *Internet Interventions*, 1(4):165–168. Publisher: Elsevier.

Moumena Chaqfeh, Rohail Asim, Bedoor AlShebli, Muhammad Fareed Zaffar, Talal Rahwan, and Yasir Zaki. 2023. *Towards a World Wide Web without digital inequality*. *Proceedings of the National Academy of Sciences*, 120(3):e2212649120.

Kevin Coe, Kate Kenski, and Stephen A Rains. 2014. Online and uncivil? patterns and determinants of incivility in newspaper website comments. *Journal of communication*, 64(4):658–679.

Molly J Crockett. 2017. Moral outrage in the digital age. *Nature human behaviour*, 1(11):769–771.

Mattia Gasparini, Robert Clarisó, Marco Brambilla, and Jordi Cabot. 2020. *Participation Inequality and the 90-9-1 Principle in Open Source*. In *Proceedings of the 16th International Symposium on Open Collaboration*, pages 1–7, Virtual conference Spain. ACM.

Maria Glenski, Svitlana Volkova, and Srijan Kumar. 2020. *User Engagement with Digital Deception*. In Kai Shu, Suhang Wang, Dongwon Lee, and Huan Liu, editors, *Disinformation, Misinformation, and Fake News in Social Media*, pages 39–61. Springer International Publishing, Cham. Series Title: Lecture Notes in Social Networks.

Eszter Hargittai. 2018. The digital reproduction of inequality. In *The inequality reader*, pages 660–670. Routledge.

Eszter Hargittai and Aaron Shaw. 2015. Mind the skills gap: the role of internet know-how and gender in differentiated contributions to wikipedia. *Information, communication & society*, 18(4):424–442.

Ariel Hasell and Brian E Weeks. 2016. Partisan provocation: The role of partisan news use and emotional responses in political information sharing in social media. *Human Communication Research*, 42(4):641–661.

Edda Humprecht, Lea Hellmueller, and Juliane A. Lischka. 2020. *Hostile Emotions in News Comments: A Cross-National Analysis of Facebook Discussions*. *Social Media + Society*, 6(1):205630512091248.

Nanak C Kakwani. 1977. Applications of lorenz curves in economic analysis. *Econometrica: Journal of the Econometric Society*, pages 719–727.

707	TaeYoung Kang, Eunrang Kwon, Junbum Lee,	Ian Rowe. 2015. Civility 2.0: A comparative analysis of	762
708	Youngeun Nam, Junmo Song, and JeongKyu Suh.	incivility in online political discussion. <i>Information,</i>	763
709	2022. Korean online hate speech dataset for mul-	<i>communication &amp; society</i> , 18(2):121–138.	764
710	tilabel classification: How can social science aid		
711	developing better hate speech dataset? <i>Preprint,</i>	Arthur D Santana. 2014. Virtuous or vitriolic: The	765
712	arXiv:2204.03262.	effect of anonymity on civility in online newspa-	766
		per reader comment boards. <i>Journalism practice,</i>	767
713	Jin Woo Kim, Andrew Guess, Brendan Nyhan, and Ja-	8(1):18–33.	768
714	son Reifler. 2021a. The distorting prism of social		
715	media: How self-selection and exposure to incivility	Anique Scheerder, Alexander Van Deursen, and Jan	769
716	fuel online comment toxicity. <i>Journal of Communi-</i>	Van Dijk. 2017. Determinants of internet skills, uses	770
717	<i>cation</i> , 71(6):922–946.	and outcomes. a systematic review of the second-and	771
		third-level digital divide. <i>Telematics and informatics,</i>	772
718	Seonghyun Kim. 2022. Korean unsmile dataset:	34(8):1607–1624.	773
719	Human-annotated multi-label korean hate speech		
720	dataset. <a href="https://github.com/smilegate-ai/korean_unsmile_dataset">https://github.com/smilegate-ai/</a>	Nicholas A Valentino, Ted Brader, Eric W Groenendyk,	774
721	<a href="https://github.com/smilegate-ai/korean_unsmile_dataset">korean_unsmile_dataset</a> .	Krysha Gregorowicz, and Vincent L Hutchings. 2011.	775
		Election night’s alright for fighting: The role of emo-	776
722	Youngjoo Kim, Yoonjin Shin, Hayoung Sim, Yoonjae	tions in political participation. <i>The journal of politics,</i>	777
723	Jang, and Park Mingyoo. 2021b. <i>Media Users in Ko-</i>	73(1):156–170.	778
724	<i>rea 2021</i> . Technical report, Korea Press Foundation.		
		Jan AGM Van Dijk. 2006. <i>Digital divide research,</i>	779
725	Vladimir Korovkin, Albert Park, and Evgeny Kaganer.	<i>achievements and shortcomings.</i> <i>Poetics</i> , 34(4-	780
726	2023. Towards conceptualization and quantification	5):221–235. Publisher: Elsevier.	781
727	of the digital divide. <i>Information, Communication &amp;</i>		
728	<i>Society</i> , 26(11):2268–2303.	Trevor Van Mierlo. 2014. <i>The 1% rule in four digi-</i>	782
		<i>tal health social networks: an observational study.</i>	783
729	Junbum Lee. 2020. Kcbert: Korean comments bert. In	<i>Journal of medical Internet research</i> , 16(2):e2966.	784
730	<i>Annual Conference on Human and Language Techno-</i>	Publisher: JMIR Publications Inc., Toronto, Canada.	785
731	<i>logy</i> , pages 437–440. Human and Language Techno-		
732	<i>logy</i> .	Xudong Yu, Magdalena Wojcieszak, and Andreu Casas.	786
		2024. Partisanship on social media: In-party love	787
733	Gina M Masullo, Shuning Lu, and Deepa Fadnis. 2021.	among american politicians, greater engagement with	788
734	Does online incivility cancel out the spiral of silence?	out-party hate among ordinary users. <i>Political Be-</i>	789
735	a moderated mediation model of willingness to speak	<i>havior</i> , 46(2):799–824.	790
736	out. <i>New Media &amp; Society</i> , 23(11):3391–3414.		
		<b>A Appendix</b>	791
737	Jakob Nielsen. 2006. The 90-9-1 rule for partici-	<b>A.1 Descriptive Statistics for the Comment</b>	792
738	pation inequality in social media and online com-	<b>Dataset</b>	793
739	munities. <a href="https://www.nngroup.com/articles/participation-inequality/">https://www.nngroup.com/articles/</a>	<b>Change in the Size of Comment Space</b>	794
740	<a href="https://www.nngroup.com/articles/participation-inequality/">participation-inequality/</a> .	The size of the comment space has grown rapidly	795
741	Accessed: 2024-01-06.	over the years (Figure 9), and since our analysis	796
		focuses only on articles that received comments,	797
742	Se-Uk Oh, Ahran Park, and Jinho Choi.	we exclude users who did not engage in posting	798
743	2021. Digital news report in korea 2021.	comments. This means we overlook the 90% of	799
744	<a href="https://www.kpf.or.kr/front/research/selfDetail.do?seq=592216">https://www.kpf.or.kr/front/research/</a>	users, often referred to as "Lurkers" in the 90-9-1	800
745	<a href="https://www.kpf.or.kr/front/research/selfDetail.do?seq=592216">selfDetail.do?seq=592216</a> .	principle.	801
		<b>Distribution of the Number of Comments</b>	802
746	Mathilde Pak, Christophe André, and Jinwoan Beom.	Online comment space is highly skewed. The his-	803
747	2021. DIGITALIZATION IN KOREA: A PATH	togram in Figure 10 indicates that the majority of	804
748	TO BETTER SHARED PROSPERITY? Technical	users post one or two comments. When dealing	805
749	report, Korea Economic Institute of America.	with a highly skewed distribution, it is generally	806
		more appropriate to consider specific percentiles,	807
750	Zizi Papacharissi. 2004. Democracy online: Civility,	as there is a significant difference in values between	808
751	politeness, and the democratic potential of online	the top and the bottom of the distribution. Hence,	809
752	political discussion groups. <i>New media &amp; society,</i>	this paper compares only top 10% and bottom 40%	810
753	6(2):259–283.	groups.	811
754	Steve Rathje, Jay J Van Bavel, and Sander Van Der Lin-		
755	den. 2021. Out-group animosity drives engage-		
756	ment on social media. <i>Proceedings of the National</i>		
757	<i>Academy of Sciences</i> , 118(26):e2024292118.		
758	Patrícia Rossini. 2022. <i>Beyond Incivility: Understand-</i>		
759	<i>ing Patterns of Uncivil and Intolerant Discourse in</i>		
760	<i>Online Political Talk.</i> <i>Communication Research,</i>		
761	49(3):399–425.		

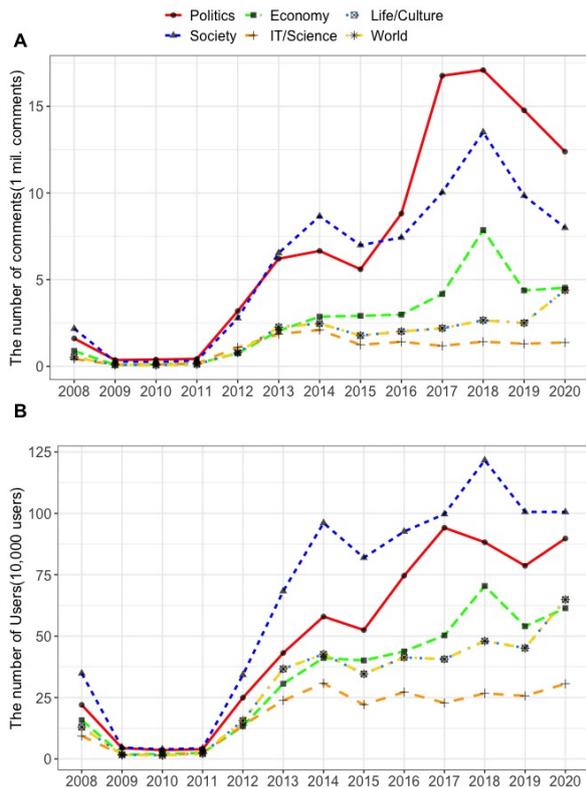


Figure 9: Change in the size of comment space: A. Change in the number of comments over time. B. Change in the number of users over time

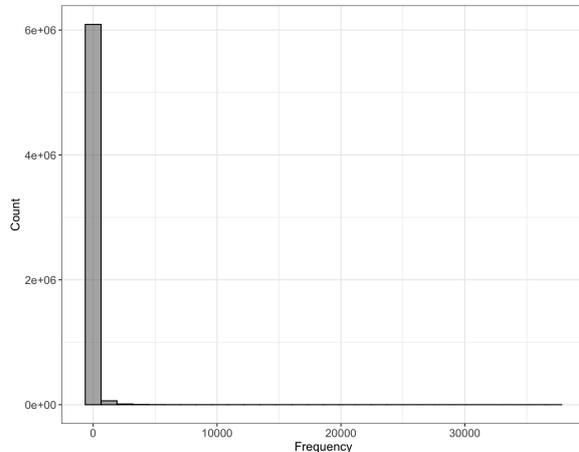


Figure 10: Histogram of the Comment Frequency

## A.2 Training Performance(KC-BERT)

For training KC-BERT, we primarily trained two models: KC-BERT Base and KC-BERT Large. The KC-BERT Large model is larger than KC-BERT Base, with significantly more parameters.

To identify the model with the best performance, we conducted experiments using various hyperparameters, such as learning rate and batch size. Additionally, recognizing that the construction of the train/validation set could influence performance, we repeated the experiments with different configurations of the train/validation split.

The table 1 presents the best performance achieved by each model, with KC-BERT Base yielding slightly better results. The table 2 displays the classification performance for KC-BERT Base

Model	LRAP
KC-BERT base	<b>0.923</b>
KC-BERT large	0.92

Table 1: Model performance based on label ranking average precision(LRAP)

Category	Precision	Recall	F1-Score	Support
0	0.82	0.72	0.77	423
1	0.87	0.81	0.84	341
2	0.87	0.81	0.84	326
3	0.85	0.76	0.80	436
4	0.87	0.83	0.85	160
5	0.89	0.87	0.88	387
6	0.88	0.89	0.89	319
7	0.93	0.17	0.29	148
8	0.72	0.57	0.64	832
9	0.93	0.92	0.93	3990

Table 2: Classification Performance based on Precision, Recall, and F1-Score