Infrastructure Ombudsman: Mining Future Failure Concerns from Structural Disaster Response

¹Rochester Institute of Technology

{mac9908, sd3528, naveen.sharma, axkvse}@rit.edu

Abstract

Current research concentrates on studying discus-1 sions on social media related to structural fail-2 ures to improve disaster response strategies. How-3 ever, detecting social web posts discussing con-4 cerns about anticipatory failures is under-explored. 5 If such concerns are channeled to the appropriate 6 authorities, it can aid in the prevention and miti-7 gation of potential infrastructural failures. In this 8 paper, we develop an infrastructure ombudsman 9 - that automatically detects specific infrastructure 10 concerns. Our work considers several recent struc-11 tural failures in the US. We present a first-of-its-12 kind dataset of 2,662 social web instances for this 13 novel task mined from Reddit and YouTube. 14

15 **1** Introduction

On January 28, 2022, at 6.39 a.m. EST, the Fern Hollow 16 Bridge in Pittsburgh, Pennsylvania collapsed. Due to the 17 timing of the failure, thankfully, fewer vehicles were on the 18 bridge and only ten people were injured with no fatalities. 19 Pittsburgh, also known as the City of Bridges, was getting 20 ready for a visit from President Biden that day. Biden vis-21 ited the collapse site and assured federal assistance to rebuild 22 the bridge on the spot. This infrastructural failure, coinciding 23 with a high-profile political visit and a push towards passing 24 the Build Back Better infrastructure bill, attracted consider-25 able media attention to the flailing infrastructural health in 26 the US. 27

As we were sifting through the social web discussions sur-28 rounding this issue, broad themes such as words of compas-29 sion for the victims and typical responses in social web polit-30 ical discourse such as political name-calling, conspiracy the-31 ories, and partisan mud-slinging emerged. However, apart 32 from these expected social web reactions, we noticed a small 33 minority of interactions that talked about anticipatory failures 34 of other bridges in the US. Table 1 lists a few illustrative ex-35 amples. 36

The comments in Table 1, if mined efficiently and surfaced to the appropriate authorities, can present an effective path to intercept structural failure concerns. Failures that are yet to happen - but a responsible citizen is worried that they might.

Social Media Post
There is a bridge in Lowell Massachusetts, it goes over
the Merrimack river and it is rusted strait through. It
won't be long before we suffer major injuries because
that bridge is always bumper to bumper traffic!
I'm surprised the New Kensington bridge wasn't the
first to go. Haha. Terrible condition.
The bridge on 1-81 that spans the Potomac between WV
and MD that could be next. It has a lot of 18 wheelers
beating it up. I kept a hammer in my car to get out if it
collapsed when I was on it.

Table 1: Illustrative examples indicating concerns over other vulnerable bridges following the collapse of the Fern Hollow Bridge in Pittsburgh.

Vulnerabilities that they perhaps noticed before, but the sud-41 den, exogenous shock in the form of a structural failure gives 42 them an outlet to voice their concerns. Extant research on dis-43 aster response discourse focuses on a diverse set of tasks that 44 include: efficient distribution of relief [Varga et al., 2013], 45 crisis management, handling emergencies etc [Horita et al., 46 2017]. Much of this research has focused on natural disas-47 ters such as typhoons [Zou et al., 2023], earthquakes [Sakaki 48 et al., 2012], floods [Feng et al., 2020], and accidents with 49 severe fatalities [Liu et al., 2020]. Analyzing social media re-50 sponses to structural failures such as bridges or building col-51 lapses is rather new and looking beyond the immediate, and 52 focusing on future potential crises has no prior literature to 53 the best of our knowledge. 54

This paper¹ presents *infrastructure ombudsman*, an auto-55 mated social media listener tool that surfaces infrastructure 56 concerns. Via a curated corpus of 2,662 instances (271 posi-57 tives and 2391 negatives)², we demonstrate that state-of-the-58 art NLP and AI methods can be harnessed to build such tools 59 that can aid humans in identifying infrastructure concerns. 60 While the Pittsburgh Bridge collapse was an exogenous shock 61 that got people discussing other potential vulnerabilities, it is 62 possible that our findings point to a broader human pattern 63 where structural collapses trigger similar thoughts about an-64

¹Paper has been published in The ACM Web Conference (WWW) 2024

²https://github.com/towhidabsar/InfrastructureOmbudsman



Figure 1: Dataset Creation Pipeline

ticipatory failures. Our dataset thus considers several recent infrastructure failures in the US and combines effective rare-

class mining methods to present a more holistic resource tothis novel task.

69 Our contributions are the following:

Novel task: detecting infrastructure concerns: We define
a new task of detecting infrastructure concerns from social
web posts. Our task presents a marked departure from existing disaster response literature where our goal is to identify
citizen concerns about possible, future structural failures.

New resource: We release a dataset of 2,662 instances 75 (271 positives and 2,391 negatives) of infrastructure concerns 76 mined from millions of Reddit posts and YouTube comments. 77 We present a suite of strong baselines trained on this dataset 78 and demonstrate the feasibility of automatically detecting in-79 frastructure concerns with reasonably high precision and re-80 81 call. On unseen data, we demonstrate that our content classifier can effectively aid humans in detecting potential candi-82 dates that express infrastructure concerns. 83

84 2 Dataset

We curate the first dataset for this novel task by collecting so-85 cial media instances from Reddit and YouTube by employing 86 a multi-step pipeline to filter the data and identify the most 87 likely candidates expressing infrastructure concerns. The 88 process begins with keyword filtering using a set of search 89 terms related to various structural failures. This high-recall 90 approach [Halterman et al., 2021; Dutta et al., 2022] yields 91 an initial corpus of over 140,326 Reddit posts and 416,009 92 YouTube comments. Next, we prune the data using natu-93 ral language inference (NLI) where we only retain instances 94 where an off-the-shelf NLI system predicts entailment be-95 tween the post and the hypothesis "There is a growing in-96 frastructure concern somewhere." The filtered data is then 97 annotated using a large language model (PaLM 2) [Chowdh-98 ery et al., 2022] in a zero-shot setting. We designed a prompt 99

to guide the annotation process [Ziems et al., 2023], aligning 100 the model's responses with the desired labels. This machine 101 annotation step further narrows the dataset to 243 Reddit ex-102 amples and 2,419 YouTube examples. Each step of our data 103 creation process, visualized in Figure 1, focused on filtering 104 down the sample to just the most likely positive instances 105 of anticipatory infrastructure concerns. After collecting an 106 initial dataset of over 2 million social media discussions 107 from both Reddit and YouTube, we filtered it down to over 108 10,251 positives with textual entailment. Another round of 109 LLM-based zero-shot annotation lowered our sample size to 110 \sim 3.000 possible positives. The first round of partisan MTurk 111 annotation yielded 1,000 possible comments labeled as future 112 infrastructure concerns by at least two crowd-sourced annota-113 tors. The two expert annotators then reduced this to just 271 114 highly confident positive examples, compared to 2,391 con-115 firmed challenging negative examples (as it passed through 116 our three filters including crowdsourced workers), for a total 117 annotated corpus of 2,662 data points. 118

3 Experiments

A key contribution of our work is developing natural language 120 processing methods to automatically identify discussions ex-121 pressing anticipatory infrastructure concerns. We frame this 122 as a binary text classification problem, where the goal is to 123 determine whether a given text contains evidence of concerns 124 about potential infrastructure failures. In what follows, we 125 design and evaluate a suite of classifiers trained on this bi-126 nary classification task. 127

119

128

3.1 Zero-Shot Classification

We select three well-known LLMs (PaLM 2 [Chowdhery 129 et al., 2022], GPT-3.5-Turbo [noa,], and Mistral 130 AI [Jiang et al., 2023]) to perform zero-shot classification 131 as our baselines. As already mentioned, this prompt has been 132 designed following best practices suggested by [Ziems et al., 133 2023]. We present the social web post following the prompt. 134

Model	Precision	Recall	F1	Accuracy
BERT _{nomask}	0.79±8e-5	0.82±1e-4	0.81±5e-4	0.92±1e-5
BERT _{mask}	0.80±2e-4	0.83±2e-4	0.81±1e-4	0.93±2e-5
RoBERTa _{nomask}	0.78±8e-5	0.83±1e-4	0.80±3e-5	0.92±2e-5
RoBERTa _{mask}	0.82±1e-5	0.83±3e-4	0.82±9e-5	0.93±3e-5
LLAMA2 _{nomask}	0.83±2e-5	0.78±5e-5	0.80±2e-5	0.93±4e-5
LLAMA2 _{mask}	0.81±1e-4	0.79±1e-4	0.80±1e-4	0.93±4e-5
Mistral _{nomask}	0.77±9e-6	0.75±5e-5	0.76±8e-5	0.91±6e-5
Mistral _{mask}	0.83±2e-5	0.79±1e-4	0.80±5e-5	0.93±3e-5
Mistral _{zero}	0.54	0.57	0.31	0.32
GPT3.5-Turbo _{zero}	0.53	0.57	0.35	0.38
PaLM2 _{zero}	0.58	0.63	0.59	0.80

Table 2: The results of the classifier performances on our dataset. We split our corpus into training and validation sets (70:30). All values are the macro average results evaluated across a 5-fold run over the validation set except for the zero-shot classification. We report the mean along with the variance.

135 3.2 Supervised Classification

We consider multiple classifier models including smaller lan-136 guage models with around 700M parameters (BERT [Devlin 137 138 et al., 2018] and ROBERTA [Liu et al., 2019]) as well as large 139 language models with 7B parameters (Mistral and Llama 2 [Touvron *et al.*, 2023]). We fine-tuned each model on our 140 dataset using a 70/30 train/test split for 5 epochs using Adam 141 optimizer [Kingma and Ba, 2017]. We consider standard ma-142 chine learning performance metrics precision, recall, F1, and 143 accuracy on the held-out validation set. 144

Most of the instances present in our dataset mention physi-145 cal locations. All these locations can be clubbed into a place-146 holder <LOCATION> which might benefit a text classifier 147 not to attend to irrelevant information. To isolate the impact 148 of locale-specific references, we ran experiments under two 149 settings: (1) masked locations, and (2) no masked 150 locations. Under the masked location setting, we 151 utilized named entity recognition (NER) to identify locations 152 and geopolitical entities like cities and states within each 153 comment using spaCy [Montani et al., 2023], replacing them 154 with the generic token <LOCATION>. With no masked 155 locations, we left the comments unmodified. This en-156 157 abled us to evaluate whether classifiers rely heavily on local references to identify infrastructure concerns, or can infer 158 these concerns solely from high-level semantic content. Eval-159 uating performance with masked and unmasked locations can 160 reveal opportunities to improve model robustness. Our exper-161 iments aim to determine the degree to which anticipatory in-162 frastructure concerns can be detected from language patterns 163 alone, without relying on localization signals. 164

165 4 Results

Table 2 summarizes our supervised solutions' performance. We observe all supervised solutions attaining reasonable precision, recall, and F1 score. We do not observe any acrossthe-board discernible benefit in masking or not masking the location information indicating that our models are most likely learning from the semantic content of the infrastructure concern rather than being fixated on location mentions.

173 In contrast with the fine-tuned models, we observe that the

zero-shot models perform considerably poorly, with PaLM 2 174 emerging as the winner among the zero-shot LLMs. This is 175 possibly due to the nuanced nature of our task. 176

177

194

4.1 Performance In The Wild

To be practically useful, the infrastructure ombudsman, i.e., our content classifier, will have to effectively mine infrastructure concerns in the wild. To evaluate in-the-wild performance, we employ the top-performing classifier from Table 2 (ROBERTa_{mask}) to previously unseen $D_{in-the-wild}$.

Running on this data, our classifier flagged 2,116 comments as positive anticipatory infrastructure concerns, and 7,884 as negative. To estimate precision and recall on this unlabeled set, we manually annotated a random sample of 100 predicted positives and 100 predicted negatives.

On these manually annotated samples, our classifier achieved a macro precision of 0.82, recall of 0.85, F1 score of 0.85, and an accuracy of 0.82, in identifying true infrastructure concerns outside of the training data. This demonstrates promising generalization to unseen data, with a high precision indicating most flagged instances were true positives. 188 199 190 192 193

5 Conclusion and Discussions

This paper presents a novel direction in mining anticipatory 195 concerns from social media following structural failures (e.g., 196 bridge or building collapses). To this end, we present the 197 first dataset resource to perform this novel task. We present a 198 suite of strong baselines relying on the recent advancements 199 in large language models. Our automated anticipatory infras-200 tructure concern mining tool, dubbed infrastructure ombuds-201 man, performs effectively in the wild. 202

Beyond infrastructure concerns: Our study reveals an in-203 teresting pattern in human behavior where we observe that 204 people often discuss potential vulnerabilities following a 205 structural failure. There could be broader generalizability to 206 this phenomenon. Following a brawl in a bar that results in a 207 shooting incident can trigger discussions on other bars where 208 violence might happen. Similar methods can be employed to 209 intercept such concerns. 210

References 211

- [Chowdhery et al., 2022] Aakanksha Chowdhery, Sharan 212 Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, 213 Adam Roberts, Paul Barham, Hyung Won Chung, Charles 214
- Sutton, Sebastian Gehrmann, and others. Palm: Scal-215
- ing language modeling with pathways. arXiv preprint 216
- arXiv:2204.02311, 2022. 217
- [Devlin et al., 2018] Jacob Devlin, Ming-Wei Chang, Ken-218 ton Lee, and Kristina Toutanova. Bert: Pre-training of 219
- deep bidirectional transformers for language understand-220
- ing. arXiv preprint arXiv:1810.04805, 2018. 221
- [Dutta et al., 2022] Sujan Dutta, Beibei Li, Daniel S. Nagin, 222 and Ashiqur R. KhudaBukhsh. A Murder and Protests, the 223 Capitol Riot, and the Chauvin Trial: Estimating Disparate 224 News Media Stance. In Proceedings of the Thirty-First 225 IJCAI, pages 5059-5065, Vienna, Austria, July 2022. In-226
- ternational Joint Conferences on Artificial Intelligence Or-227 ganization. 228
- [Feng et al., 2020] Yu Feng, Claus Brenner, and Monika 229 Sester. Flood severity mapping from Volunteered Geo-230 graphic Information by interpreting water level from im-231 ages containing people: A case study of Hurricane Harvey. 232 ISPRS Journal of Photogrammetry and Remote Sensing, 233 169:301-319, 2020. Publisher: Elsevier. 234
- [Halterman et al., 2021] Andrew Halterman, Katherine 235 Keith, Sheikh Sarwar, and Brendan O'Connor. Corpus-236 Level Evaluation for Event QA: The IndiaPoliceEvents 237 Corpus Covering the 2002 Gujarat Violence. In Findings 238 of the Association for Computational Linguistics: ACL-239 IJCNLP 2021, pages 4240-4253, Online, August 2021. 240 Association for Computational Linguistics. 241
- [Horita et al., 2017] Flávio EA Horita, João Porto de Albu-242 querque, Victor Marchezini, and Eduardo M Mendiondo. 243 Bridging the gap between decision-making and emerging 244 big data sources: An application of a model-based frame-245 work to disaster management in Brazil. Decision Support 246 Systems, 97:12-22, 2017. Publisher: Elsevier. 247
- [Jiang et al., 2023] Albert Q. Jiang, Alexandre Sablay-248 rolles, Arthur Mensch, Chris Bamford, Devendra Singh 249 Chaplot, Diego de las Casas, Florian Bressand, Gi-250 anna Lengyel, Guillaume Lample, Lucile Saulnier, 251 Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, 252 253 Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7B, October 2023. 254 arXiv:2310.06825 [cs]. 255
- [Kingma and Ba, 2017] Diederik P. Kingma and Jimmy Ba. 256 Adam: A Method for Stochastic Optimization, January 257 2017. arXiv:1412.6980 [cs]. 258
- [Liu et al., 2019] Yinhan Liu, Myle Ott, Naman Goyal, 259 Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, 260 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 261 RoBERTa: A Robustly Optimized BERT Pretraining Ap-262
- proach. 2019. Publisher: arXiv Version Number: 1. 263
- [Liu et al., 2020] Tiezhong Liu, Huyuan Zhang, and Hubo 264 Zhang. The impact of social media on risk communica-265 tion of disasters-a comparative study based on sina weibo 266

blogs related to tianjin explosion and typhoon pigeon. In-267 ternational journal of environmental research and public 268 health, 17(3):883, 2020. Publisher: MDPI. 269

- [Montani et al., 2023] Ines Montani, Matthew Honnibal, 270 Matthew Honnibal, Adriane Boyd, Sofie Van Landeghem, 271 and Henning Peters. explosion/spaCy: v3.7.2: Fixes for 272 APIs and requirements, October 2023. 273
- [noa,] GPT-3.5 Turbo fine-tuning and API updates.
- [Sakaki et al., 2012] Takeshi Sakaki, Makoto Okazaki, and 275 Yutaka Matsuo. Tweet analysis for real-time event de-276 tection and earthquake reporting system development. 277 IEEE transactions on knowledge and Data Engineering, 278 25(4):919-931, 2012. Publisher: IEEE. 279
- [Touvron et al., 2023] Hugo Touvron, Louis Martin, Kevin 280 Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, 281 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, 282 Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Can-283 ton Ferrer, Moya Chen, Guillem Cucurull, David Es-284 iobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian 285 Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, 286 Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan 287 Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, 288 Isabel Kloumann, Artem Korenev, Punit Singh Koura, 289 Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana 290 Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, 291 Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin 292 Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, 293 Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael 294 Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh 295 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, 296 Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, An-297 gela Fan, Melanie Kambadur, Sharan Narang, Aurelien 298 Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas 299 Scialom. Llama 2: Open Foundation and Fine-Tuned Chat 300 Models, July 2023. arXiv:2307.09288 [cs]. 301
- [Varga et al., 2013] István Varga, Motoki Sano, Kentaro 302 Torisawa, Chikara Hashimoto, Kiyonori Ohtake, Takao 303 Kawai, Jong-Hoon Oh, and Stijn De Saeger. Aid is out 304 there: Looking for help from tweets during a large scale 305 disaster. In Proceedings of the 51st annual meeting of the 306 association for computational linguistics (volume 1: Long 307 papers), pages 1619-1629, 2013. 308
- [Ziems et al., 2023] Caleb Ziems, William Held, Omar 309 Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. Can 310 Large Language Models Transform Computational Social 311 Science?, April 2023. arXiv:2305.03514 [cs]. 312
- [Zou et al., 2023] Lei Zou, Danqing Liao, Nina SN Lam, 313 Michelle A Meyer, Nasir G Gharaibeh, Heng Cai, Bing 314 Zhou, and Dongying Li. Social media for emergency res-315 cue: An analysis of rescue requests on Twitter during Hur-316 ricane Harvey. International Journal of Disaster Risk Re-317 duction, 85:103513, 2023. Publisher: Elsevier. 318

274