Reducing Processing Time and Enhancing Classification Performance: Shortening Strategies for German Public Contributions

Anonymous ACL submission

Abstract

Public participation, the process of 2 voluntary engagement of citizens in urban 3 decision-making, requires a significant 4 amount of time and human resources. 5 Therefore, automatization of the evaluation 6 is essential. Classification of the citizens' 7 proposals is one of the prevalent analytical 8 tasks. Through the years many studies have 9 worked on automatization techniques of 10 this procedure and Natural Language 11 Processing (NLP) methods are among the 12 most effective ones. Nevertheless, most 13 developed techniques despite promising 14 results are optimized for the English 15 language. Moreover, the NLP pre-trained 16 models such as BERT have limitations in 17 the length of the texts they can process. 18 Hence, this paper focuses on the abstractive 19 summarization of the public proposal in 20 German and considers two different 21 shortening techniques (truncation and 22 summarization). The main aim is to explore 23 how pre-trained models such as the BERT 24 perform in the classification of summarized 25 German language contributions. For this 26 purpose, the German BERT model, which 27 is fine-tuned on the MLSUM DE dataset, 28 and the multilingual BART model are 29 considered for text summarization. The 30 results revealed that applying shortening 31 techniques on long contributions reduces 32 the model development time by an average 33 of 48% on CPU and 36% on GPU while 34 improving performance. Moreover, the 35 multilingual BART model works slightly 36 better than the BERT model fine-tuned on 37 the MLSUM DE dataset. 38

39 Introduction

41 academics, regulators, and governments, is one of 81 existing research is limited to English, with only a

42 the key instruments of democracy. Public 43 contribution helps citizens' concerns and thoughts 44 to be heard and avoids elitism (Rowe and Frewer, ⁴⁵ 2004). Even though consideration of the public ⁴⁶ voice provides more information to the authorities, 47 it requires a longer processing time (Lee and Kim, 48 2014). Indeed, processing and analyzing the 49 collected data is among the key challenges that 50 policy-makers in the public voice consideration 51 face (Romberg and Escher, 2023, Simonofski, 52 Fink, and Burnay, 2021, Arana-Catania et al., 53 2021). To protect the ultimate goals of public 54 participation and the core standards of a democratic ⁵⁵ society, all contributions should be treated equally. ⁵⁶ This emphasizes the requests for specific analytical 57 methods to ease inferring meaningful information 58 from an overwhelming collection of public 59 statements (Lee and Kim, 2014, Romberg, Mark, 60 and Escher, 2022, Romberg and Escher, 2022).

Previous experiences have demonstrated that 61 62 when the amount of data is manageable, manual 63 processing despite its inefficiency is preferred 64 (Simonofski, Fink, and Burnay, 2021). However, 65 the advances in technology especially in Artificial 66 Intelligence (AI) and particularly Natural 67 Language Processing (NLP) have made a huge 68 difference (Romberg and Escher, 2023). NLP 69 analyses linguistic data and enables computers to 70 understand, interpret, and generate human ⁷¹ language (Egger and Gokce, 2022). Despite 72 notable advances in NLP techniques, such as the 73 rise of Pre-trained Language Models (PLMs), 74 Romberg and Escher (2023) in a review of methods 75 for computational text analysis remark that public 76 contribution texts exhibit distinct characteristics 77 compared to other domains such as news or social 78 media contents (Romberg and Escher, 2023). 79 Another significant gap in the study of public 40 Public participation, an attractive topic for 80 contributions is language. The majority of the ⁸² few addressing non-English public proposal texts ¹³¹ 2

83 (Romberg and Escher, 2023, Romberg, 2023, Balta

⁸⁴ et al., 2019, Romberg and Conrad, 2021, Romberg, ¹³² 2.1

85 2022). Hence further studies are required to 133 Public participation, alternatively called citizen ⁸⁶ generate more robust and reliable models.

Automatic Text Summarization (ATS) is 135 influence 87 ⁸⁸ another field of NLP, which despite its high ₁₃₆ (Schroeter et al., 2016). However, processing all 89 potential has been left under-explored, especially in 137 the collected statements to infer meaningful ⁹⁰ the study of non-English texts (Aruneshwari et al., 138 information manually, while respecting the 91 2024, Alcantara et al., 2023). ATS is a method for 139 principles of democracy is a challenging task 92 generating a coherent shorter version of a text 140 (Karic et al., 2024, Arana-Catania et al., 2021). ⁹³ document while preserving the key information 141</sup> Hence, to achieve this goal (Kwon et al., 2006), ⁹⁴ from the original text (Aruneshwari et al., 2024). 142 (Arguello et al., 2008), (Habernal & Gurevych, 95 ATS is classified into two categories. (1) The 143 2016), and (Cardie et al., 2008) applied SVM 96 abstractive text summarization, in which the 144 models, (Konat et al., 2016) considered graph-97 system tries to understand the text and provides a 145 based analytics, and with the progress of AI-based 98 shorter novel narrative of the original text while 146 technologies in recent years such as Natural ⁹⁹ preserving the key points (Alcantara et al., 2023). ₁₄₇ Language Processing (NLP) techniques, (Fierro et 100 (2) The extractive summarization, which is more 148 al., 2017) considered fastText and Deep Averaging 101 flexible where the system first identifies the key 149 Networks (DAN) to automatically analyze or 103 104 Verma and Verma, 2020, Talib, 2021).

105 106 summarization is in the classification of the texts. 154 mining and compared the performance of SVM, 107 However, this area is in its early stages and the 155 Random Forest (RF), and k-Nearest-Neighbor 109 English language. (Scialom et al., 2020) mentions 157 Escher, 2022) studied in the domain of argument 110 that the dominance of the English language, lack of 158 mining in German introduced a new dataset of multilingual data, and the use of pre-trained models 159 German language arguments and (Romberg, Mark, 112 on the pivot language (English) resulted in a 160 and Escher, 2022) compared the performance of 113 significant gap in the performance of the 161 SVM, Maximum Entropy (MaxEnt), Naïve Bayes 114 classification studies. To the best of our knowledge, 162 (NB), and BERT classifiers in single and multi-115 this is the first study in summarization of public 163 label classification of contributions. 116 participation in the German language. Our 117 contribution is mainly on understanding how 164 2.2 ¹¹⁸ summarization affects the models' performance in ¹⁶⁵ Text Summarization has been always an interesting 119 detecting argumentation structures.

The rest of this paper is organized as follows. 120 121 Section 2 highlights the importance of analyzing 122 public participation and mentions the previously 123 explored techniques. This section also briefly 124 reviews the Automatic Summarization techniques 125 and models. Section 3 introduces the dataset. 126 Section 4 comprehensively explains the developed 127 methods in this study. In Section 5, the obtained 128 results are presented and are discussed in Section 129 6. Finally, Section 7 summarizes the findings and 130 mentions the topics for further studies.

Literature Review

Evaluation of Public Participation

134 participation, is the volunteer action of citizens to the decision-making processes sentences (texts), then extracts, orders, and returns 150 classify a large number of public comments and them in a condensed form (Alcantara et al., 2023, 151 arguments on proposed regulations. (Liebeck et al., 152 2016) is one of the few researchers in this domain One of the main applications of text 153 who developed a new German corpus for argument majority of the studies are still focused on the 156 (kNN) classification models. (Romberg and

Automated Text Summarization

166 topic for researchers. In 2010 with the introduction 167 of neural networks and later the development of 168 Sequence-to-Sequence (Seq2Seq), Bidirectional 169 Encoder Representations from Transformers 170 (BERT), and Generative Pre-trained Transformer 171 (GPT) models, more attention was attracted to 172 automatic text summarization (Khan et al., 2023). 173 Pre-trained models have significantly improved 174 NLP capabilities (Alcantara et al., 2023). BERT is 175 one of the bidirectional models pre-trained over 176 texts from Wikipedia and has shown relatively 177 good performance in abstractive text 178 summarization. Some other known encoder and 179 decoder pre-trained models are BART developed 180 by Facebook, GPT developed by OpenAI, and 181 RoBERTa and DistilBERT both as a refined 182 version of BERT (Alcantara et al., 2023, Syed, 183 Gaol, & Matsuo, 2021). Chen et al. (2021) by 184 highlighting the huge size of famous pre-trained 185 models such as GPT and their extensive 186 computation cost, tried to save the training cost of 187 large models by transferring the learned knowledge 188 of the t smaller rained model to the large model (Chen et al., 2021). They called the model 189 BERT2BERT. 190

Pre-trained models despite significant 191 ¹⁹² performance have some limitations. The majority 193 of them such as BERT are trained over English 194 language documents which leads to lower 195 performance on non-English documents. In order 196 to overcome this limitation, multilingual versions 197 of some of the models are introduced. mBERT, the 198 multilingual version of BERT, is trained for the top 100 languages with the longest documents on 199 Wikipedia (GitHub, 2022). This model supports 200 the German language. German BERT2BERT fine- 234 Thridataset considered for this study is a collection 202 tuned on MLSUM DE for summarization is 235 of public@ontributions from these different sities in another pre-trained BERT model for the German 236 Germany Q2nd has been used for automated topic 203 204 language trained on German OpenLegalData, and news articles. The 238 processing and processing processin 205 206 207 Hugging Face website API for pre-trained models. 241 October 2017. Detailed information on the 208 209 MLSUM is the first large-scale MultiLingual 242 collection of contributions and labeling can be 209 Will Solve is the first large scale introduced in 2020 243 found in (Romberg and Escher, 2022). 210 SUMmarization dataset introduced in 2020 243 found in (Romberg and Escher, 2022). 211 (Scialom et al., 2020). The dataset contains over 244 The total number of collected contributions is 212 1.5 million article/summary pairs in five languages, 245 3139 of which 2314 were received from Bonn, 366 Minimum field aread 450 from Moers 22 Eight 213 including German, and can be used to evaluate the 246 from Elangenfeld, 3and 459 from Moers33Eight ²¹⁴ summarization models (Scialom et al., 2020).

215 2.3 **Research Gap**

217 analyzing public contributions through text- 251 there are 231 contributions from Bonn, 46 218 processing techniques. However, most of these 252 contributions from Moers, and 36 contributions 219 studies have focused on the English language, and 253 from Ehrenfeld which belonged to two categories. 220 there is a lack of exploration into the fine-tuning of 254 Of the defined categories, "Cycling Traffic 221 transformer models (such as BERT) and the 255 Management" and "Cycle Path Quality" have 222 performance of other transformer architectures 256 received the most contributions. As claimed by 223 (Romberg, Mark, and Escher, 2022). This study 257 (Romberg and Escher, 2022), the variation in the ²²⁴ aims to address these gaps and analyze the potential ²⁵⁸ distribution of contributions is affected by factors 225 limitations of applying developed models to long 259 such as city size, local infrastructures, and 226 texts. As a potential solution, automatic text 260 participators' involvement. 227 summarization techniques, whose performances on 228 argumentative data in the German language are 261 4 229 largely unexplored (Alcantara et al., 2023), are

	Sin	gle La	ıbel	Multi-Label			
	Bonn	Ehrenfeld	Moers	Bonn	Ehrenfeld	Moers	
Bicycle parking	108	22	9	112	26	9	
Cycle path quality	449	58	111	519	71	118	
Cycling traffic management	1020	195	222	1056	204	229	
Lighting	37	1	10	47	2	15	
Misc	53	5	10	84	84	27	
Obstacles	319	35	31	364	45	33	
Signage	150	16	19	182	20	27	
Traffic lights	178	34	47	197	39	51	
Total	2314	366	459	2314	366	459	

Table 1: Distribution of the arguments across labels

231 is to assess how summarization impacts runtime 232 and overall performance.

233 3 Dataset Bonn Ehrenfeld Moers Mean 67.82 62.5 65.7 Wikipedia, 237 classification (Rohiberg and Escher, 2028). The summarization and is accessible through the 240 Bonn, Ehrenfeld, and Moers from September to

247 different Q2 ategoris6 for improving the citizens' 248 experien 23 were proposed and 2 each contribution 249 249 P90 d to 132 ... lor 119 cat _____ 223 ... n ²¹⁶ Many studies have sought to reduce the burden of ²⁵⁰ overview of the dataset is given in Table 1. In total,

Experimental Setup

230 tested and evaluated. The primary goal of this study ²⁶² In this section, we provide the applied techniques ²⁶³ and classification algorithms. As is well known the ²⁶⁴ uneven distribution of the contributions over ³¹² text comprehension tasks (Lewis et al., 2020). In 265 categories can potentially degrade 266 classification performance. Hence, weighting 314 specifically initialized for the German language 267 techniques are tested as a possible solution to 315 (gbert-base) is compared to a multilingual model, ²⁶⁸ improve the results. For the classification of the ³¹⁶ "facebook/mbart-large-50". 269 multi-label contributions, to improve 270 performance, in addition to the class weights, the 318 language. The summarized contributions are 271 Area Under Curve (AUC) score is used to evaluate 319 classified model's performance across 272 the 273 thresholds. The thresholds are applied to the output $_{274}$ of the model to determine the final decision on the $_{321}$ **4.4** 275 predicted class of a contribution. Experiments are 276 done using the Google Colab GPU (T4), Python 277 programming language, and the software provided 278 by Hugging Face, Inc. (Wolf et al., 2020). The 279 Transformers library of Hugging Face eases the 280 efforts for a variety of NLP tasks.

4.1 **BERT model** 281

282 The first applied classification model is based on ²⁸³ the BERT model initialized with the case-sensitive ³³¹ of training sequences can result in training 284 gbert-base. BERT is a pre-trained model ²⁸⁵ specialized for the tokenization and encoding of the ³³³ models, all the sequences of input texts need to be ²⁸⁶ German language. However, it is fundamentally an ³³⁴ equal in length. This necessitates padding of the ²⁸⁷ encoder-only model that cannot be directly used for ³³⁵ shorter documents while padding adds no classification tasks. 289 BertForSequenceClassification which provides a 290 classifier on top of the BERT model is used. ²⁹¹ BertForSequenceClassification is a class from the ³³⁹ can handle up to 512 or 1024 tokens (Hugging 292 HuggingFace Transformers library and lets us fine-²⁹³ tune the pre-trained BERT model on a specific $_{341}^{341}$ model supports is 512. Making decisions based on sequence classification task. The results obtained $\frac{1}{342}$ the maximum length can affect the classification ²⁹⁵ from this stage are considered a baseline for this ³⁴³ model performance. Different studies have 296 experiment.

German BERT2BERT Fine-Tuned on 297 4.2 **MLSUM DE model** 298

299 ³⁰⁰ initialized with the bert-base-german-cased model ₃₄₉ study on English texts found that keeping the head and fine-tuned on the MLSUM DE dataset. It is 350 and the tail of a text and dropping the middle part 302 specifically developed for German 303 summarization. After summarizing the citizens' 352 Mutasodirin and Prasojo (2021) compared the 304 proposals, the BertForSequenceClassification is 353 performance of truncation and summarization 305 used for classification.

306 4.3 **BART Model**

308 generalizes BERT (encoder) and GPT (decoder) 358 truncation is pre-trained and combines 359 (Mutasodirin and Prasojo, 2021). 309 models. It 310 Bidirectional and Auto-Regressive Transformers, 360

the 313 this study, the performance of the model This model is the 317 Sequence-to-Sequence and supports the German using the different 320 BertForSequenceClassification.

Implementation

322 Transformers, despite the notable improvements 323 brought to the NLP world, experience a significant 324 computational overhead due to their attention 325 mechanism. Indeed, from the autoregressive nature 326 of the attention mechanism, the complexity 327 depends on the quadratic of the sequence size $_{328}$ ($O(T^2)$) (Cukier, 2024). This dependency forces a 329 limit on the sequence length that can be passed to ³³⁰ the models. On the other hand, the extended length 332 inefficiency. From the nature of transformer Thus, 336 meaningful information, and their processing ³³⁷ wastes computational resources.

338 The majority of the transformer-based models 340 Face, 2024). The maximum length that the BERT 344 proposed different strategies to overcome this 345 limitation. One of the naïve solutions is truncation 346 of the longer sequences. However, studies have 347 shown that different parts of texts carry different This model is based on the German BERT, 348 amounts of information. (Chi et al., 2019) in a text 351 gives the best performance. In another study, 354 techniques in the classification of texts from the 355 Indonesian News Article dataset (IndoSum). They ³⁵⁶ stated that using the first block of the tokens results 307 BART is a transformer-based model that 357 in the best performance in comparison to the other and summarization techniques

The motivation for this study specifically lies in and is particularly effective in text generation or ³⁶¹ the tokenization stage of the contributions. Table 2 362 summarizes the statistics of the contribution

³⁶³ lengths after tokenization. The numbers show that ³⁶⁴ the distribution of the contribution lengths among 365 the cities is very close. Moreover, almost 90% of 366 the contributions in all of the cities have less than ³⁶⁷ 128 tokens. This number is half of the maximum ³⁶⁸ length applied by (Romberg and Escher, 2022) as 369 the hyper-parameters of the BERT model. 370 Choosing 256 as the maximum length while the 371 longest contribution has 247 tokens positively 372 avoids truncation and losing ³⁷³ However, more than 90% of the texts get padded to ³⁷⁴ at least twice their length. This can potentially ⁴¹⁶ $W_i = \frac{1}{N_i} / \sum_{i=1}^{C} \frac{1}{N_i}$ weighting strategy is applied. In $_{375}$ decrease the classification performance and $_{417}$ this formula, N_i is the number of samples under the ³⁷⁶ increase the run time unnecessarily. In the ⁴¹⁸ *i*-th category, w_i is the assigned weight to the *i*-th $_{377}$ following experiments, the performance of $_{419}$ category, and C is the total number of classes. The truncated and summarized texts is compared. 378

380 baseline for the rest of the experiments, the full-text 422 class. the 423 contributions are classified using 381 382 BertForSequenceClassification classifier, and the 424 of the baseline while in the tokenization stage of 383 BERT model is initialized with the case-sensitive 425 the BERT model, the maximum length is set to 128 384 385 BERT model, the maximum length is set to 256 and 427 The fifth experiment tests the abstractive 386 the rest of the hyper-parameters of the model are 428 summarized texts generated using the multilingual 387 kept unchanged.

388 ³⁸⁹ fold cross-validation is conducted. Thus, after the ⁴³¹ shortening results in the loss of information, only ³⁹⁰ preparation step, a new set of stratified 5 folds is ⁴³² the citizens' contributions with more than 128 391 generated to ensure the same label distribution in 433 tokens (based on white-space tokenization) are 392 both the training and testing data. However, for the 434 summarized. From the hyper-parameters, the ³⁹³ city of Ehrenfeld, since the number of contributions ⁴³⁵ minimum and maximum lengths are set to 50 and ³⁹⁴ for the category of "lighting" is only 1, stratified ⁴³⁶ 120, respectively. The decision on the minimum 395 sampling fails. To overcome this limitation, we 437 and maximum lengths is based on the statistics 396 applied synonym replacement, an NLP data 438 from Table 2 to balance the need for reducing the 397 augmentation technique, and generated 4 new 439 size while keeping the summary concise. From 398 contributions by randomly replacing some of the 440 Bonn, approximately 11% (250), from Ehrenfeld 399 words from the original contribution (Ma, 2019). 441 9% (34) of contributions, and from Moers 9% (43) 400 For synonym replacement, a BERT pre-trained 442 of the contributions are summarized. The average model initialized with 'deepset/gbert-base' is 443 length of the summarized contributions is 64 for 402 chosen. Reported results in the following sections 444 Bonn, 65 for Ehrenfeld, and 63 for Moers. 403 are averaged across all folds.

404 405 dataset was addressed in Section 3. For instance, in 447 experiment compares the performance of the 406 Bonn, the number of proposals under the category 448 German BERT2BERT summarization model in 407 of "cycling traffic management" is 27 times more 449 text classification against the multilingual mBART 408 than the lightning proposals. This condition exists 450 summarization model. The BERT2BERT model is 409 also in the other two cities and among the other 451 initialized 410 categories. Although such a distribution of the 452 german-finetuned-summarization, using the earlier 411 labels in the analysis of argument data is expected, 453 mentioned hyper-parameters, and is only applied 412 it can negatively affect the performance of the 454 on texts with more than 128 tokens. The average 413 model in the prediction of the minority classes or 455 length of the summarized texts for Bonn, 414 lower its generalization. One of the solutions is 456 Ehrenfeld, and Moers are 39, 53, and 55,

	Bonn	Ehrenfeld	Moers
Mean	67.82	62.5	65.7
Minimum	3	6	4
Q1	34	28	33
Q2	56	51	54
Q3	89	82	89
P90	132	119	123
Maximum	247	231	238

Table 2: Basic statistics of the tokens' distribution

information. 415 using class weights. In the second experiment, 420 weights are applied to the loss function to penalize In the first experiment, in order to have a 421 the model for the misclassification of the minority

The third and fourth experiments are a repetition 'gbert-base' model. In the tokenization stage of the 426 and the weighting strategy is applied, respectively. 429 BART (mBART) initialized with "facebook/mbart-In order to make the results more reliable, a 5- 430 large-50" model. Since any form of document

The sixth experiment is the repletion of the fifth 445 The issue of unbalanced label distribution in the 446 experiment with class weights applied. The seventh with mrm8488/bert2bert shared457 respectively. It should be noted that although the 458 minimum length hyperparameter of the models 459 was set to 50, the best summaries may still be 460 shorter depending on the influence of the other 461 hyperparameters. More details on the distribution 462 of the contributions' length after summarization are 463 provided in Appendix B. Finally, the eighth 464 experiment measures the performance of the sixth model when class weighting is applied. 465

At the end of the next section, all experiments 466 repeated to classify the multi-labeled 467 are contributions under two different scenarios. In the 469 first scenario, the applied threshold on the models' 470 output is set to 0. In the second, however, the AUC score is used to evaluate the model's performance 471 472 across different thresholds. The threshold with the 473 highest AUC score on the train data across all the 474 folds of the cross-validation is considered as the 475 best candidate.

Results and Evaluation 476 5

477 In this section obtained results for the earlier 478 described cases are presented. Table 3 shows the 479 category-wise performance of the trained models 480 using the F-score. Reported results are the average 481 of the outcome of testing the model on the test sets from the 5-fold cross-validation. The best model (with the highest micro F-score) for each city is 483 484 highlighted.

For the city of Bonn, the first model (the 485 486 baseline) successfully classified 78% of the 487 arguments. The best and worst performances were 488 obtained for categorizing "Bicycle Parking" and "Misc" labels. Experiment 2 shows the significant 489 effect of class weights in the classification of the 490 491 imbalanced data. While the overall performance 492 (Micro F-score) has decreased by less than 1%, the 493 classification performance of the least frequent 512 contribution texts' length and hence, loss of 494 labels "Lightning" and "Misc" have improved by 495 4% and 20%, respectively. Experiment 3 shows 496 closely the same performance as the first two 497 experiments. Nevertheless, since the contributions were truncated and only the first 128 tokens were 498 499 considered, the main achievement is the 50% 500 decrease in train time. Experiment 4 improved the performance of the third experiment, attained the 502 best result, and the micro and macro f-scores ⁵⁰³ surpassed the other experiments. Experiments 5 to ⁵⁰⁴ 8 show the performance of the summarized texts. 505 As it can be seen from the results, summarized texts ⁵⁰⁶ using the multilingual model (mBART) were more ⁵⁰⁷ successful. The main reason is in the retained

	Experiments	Bicycle Parking	Cycle Path Quality	Cycling Traffic Management	Lighting	Misc	Obstacles	Signage	Traffic Lights	Micro F-score	Macro F-score	Train run-time(second) GPU	Train run-time (second) CPU
	1 _B	0.94	0.77	0.81	0.77	0.21	0.75	0.61	0.80	0.78	0.71	438	18480
	2 _B	0.93	0.76	0.80	0.81	0.42	0.73	0.59	0.80	0.77	0.73	435	18490
	3 _T	0.91	0.78	0.80	0.82	0.27	0.74	0.59	0.77	0.77	0.71	252	9185
onn	4 T	0.92	0.77	0.81	0.79	0.52	0.75	0.63	0.81	0.78	0.75	236	9385
В	5s	0.94	0.78	0.81	0.76	0.34	0.75	0.59	0.81	0.78	0.72	207	9140
	6s	0.92	0.76	0.80	0.81	0.47	0.75	0.61	0.79	0.77	0.74	213	9010
	7s	0.91	0.77	0.80	0.85	0.15	0.73	0.54	0.78	0.76	0.69	213	9010
	8 s	0.90	0.77	0.79	0.80	0.40	0.72	0.60	0.78	0.76	0.72	221	9025
	1 _B	0.87	0.45	0.78	0.80	0.00	0.56	0.00	0.10	0.68	0.45	71	2882
	2 _B	0.68	0.31	0.63	0.80	0.00	0.43	0.22	0.39	0.54	0.44	77	2941
pi	3 _T	0.83	0.51	0.77	0.00	0.00	0.52	0.00	0.00	0.67	0.33	52	1492
enfe	4 _T	0.73	0.36	0.62	0.93	0.00	0.43	0.29	0.34	0.52	0.46	49	1497
Ehr	5s	0.80	0.56	0.80	0.40	0.00	0.56	0.00	0.16	0.70	0.41	51	1508
	6s	0.78	0.35	0.61	0.93	0.00	0.47	0.17	0.35	0.52	0.46	63	1502
	7s	0.87	0.53	0.78	0.40	0.00	0.59	0.00	0.05	0.69	0.40	47	1542
	8s	0.75	0.37	0.65	0.93	0.00	0.47	0.30	0.37	0.56	0.47	59	1545
	1B	0.63	0.77	0.82	0.00	0.00	0.29	0.21	0.84	0.75	0.45	416	4004
	2B	0.63	0.66	0.70	0.00	0.00	0.36	0.33	0.78	0.65	0.43	433	4248
s	3T	0.53	0.79	0.85	0.00	0.00	0.31	0.26	0.84	0.77	0.45	264	2000
loer	4T	0.92	0.75	0.82	0.13	0.00	0.47	0.43	0.80	0.76	0.54	233	1949
Σ	5s	0.47	0.78	0.83	0.00	0.00	0.30	0.00	0.83	0.75	0.40	249	2010
	6s	0.83	0.75	0.78	0.26	0.00	0.34	0.30	0.76	0.72	0.50	244	1965
	7s	0.47	0.76	0.83	0.00	0.00	0.20	0.24	0.80	0.74	0.41	248	2070
	8 s	0.81	0.75	0.80	0.00	0.00	0.31	0.44	0.73	0.73	0.48	255	2030

Table 3: Performance of the developed classification models on single-labeled contributions. Subscripts of B, T, and S represent Baseline, Truncation, and Summarization, respectively.

⁵⁰⁸ information. As noted earlier the average length of 509 the summarized contributions using the 510 BERT2BERT model is 39 for the city of Bonn. 511 Undoubtedly this is a significant reduction in the 513 information. Moreover, experiments 6 and 8 514 despite better performance in the classification of 515 rare labels, because of decrements in the 516 classification of frequent labels, show almost ⁵¹⁷ similar performance to experiments 5 and 7.

518 For Ehrenfeld, the results are different. In the 519 first experiment, the model failed to classify 520 contributions under the "Signage" and "Misc" 521 labels. The model in the third experiment which 522 was trained on the truncated contributions failed to ⁵²³ classify the contributions of half of the categories. 524 Nevertheless, it showed a close performance to that 525 of the first model. The main reason for the

526 insignificant effect of classification failures on the 527 micro f-score is due to the low number of samples ⁵²⁸ under the failed categories. Furthermore, reducing 529 the length of contributions by half has significantly 530 decreased the train time. The best performance is achieved for the fifth experiment where the model 531 532 is trained on the summarized texts from the ⁵³³ multilingual model. The rest of the trained models 534 show that applying weighting does not improve the 535 overall performance of the models. In fact, ⁵³⁶ weighting led to model confusion by affecting the 537 decision boundaries and only enhanced the model's 538 performance in classifying "Signage" and "Traffic 539 Lights" labels. Finally, it should be noted that the 540 applied data augmentation technique in the "lighting" class was successful. However, the 542 result should be considered carefully as all the 543 samples have the same characteristics.

544 Lastly, for Moers, the results show the same 545 pattern as for the other two cities. Truncation and 546 summarization significantly decreased train time 547 while the former performed slightly better in the 548 classification task. Weighting improved macro f-549 scores by enhancing the model's performance in the 550 classification of rare categories, while micro fscores decreased. 551

Table 4 provides detailed information on the 552 performance of the train models on the multi-553 554 labeled contributions. The general behavior of the 555 models is similar to the trained ones on the single-556 labeled contributions, and the changes in the 557 training time of the models are in line with the 558 former findings. Shortening reduced the run time 577 enhance the models' performance in most 559 by almost 47% for experiments conducted using a 578 experiments. 560 CPU compared to 37% for on GPU experiments. 579

561 ⁵⁶² multi-label classification tasks, simultaneous ⁵⁸¹ raw values from the final layer of a neural network. ⁵⁶³ prediction of multiple labels is probable; the ⁵⁸² Logits are converted into probabilities or binary 564 Hamming Loss metric is typically used instead of 583 decisions according to a threshold. Determination 565 accuracy. The hamming loss shows the percentage 584 of the most appropriate threshold for every class 566 of the mislabeled arguments. ⁵⁶⁷ experiments have labeled approximately 6% of the ⁵⁸⁶ directly on it. Thus, in the following part, the Area ⁵⁶⁸ public contributions wrongly. For Ehrenfeld and ⁵⁸⁷ Under the Curve (AUC) is used to find the optimal 569 Moers, however, the difference in the percentage of 588 threshold and is computed for every training fold. 570 mislabeled contributions between the models with 589 The optimum solution is the threshold that yields 571 and without adjusted weights is significant. 590 the best performance over the folds. The results 572 Weighting has lowered the models' performance 591 after detecting and applying the best thresholds are 573 under frequent classes in favor of less frequent 592 given in Table 4. 574 ones. The results in Table 4 indicate that, aside 593 575 from the failure of no-weighting models, the 594 majority of the models have improved while micro 576 considered weighting technique also failed to 595 F-scores have declined. This outcome is an

		М	ulti-la AU	bele C ap	Multi-Labeled with AUC applied						
	iments	F-score	F-score	iing (%)	T ti (s	Train time (sec.)		F-score	iing (%)	Train time (sec.)	
	Experi	Micro	Macro	Hamm	GPU	CPU	Micro	Macro	Hamm	GPU	CPU
	1_{B}	0.77	0.69	6	384	17482	0.76	0.73	7	430	22606
	$2_{\mathbf{B}}$	0.75	0.71	6	383	18378	0.75	0.76	7	427	22126
	3 _T	0.77	0.69	6	219	9955	0.76	0.73	7	336	11529
uu	4 _T	0.75	0.71	6	221	8940	0.75	0.72	7	301	12300
Bo	$5_{\mathbf{S}}$	0.77	0.69	6	220	9935	0.75	0.72	8	256	13814
	6 s	0.75	0.71	7	225	9702	0.74	0.69	8	269	12743
	$7_{\rm S}$	0.76	0.67	6	234	9120	0.74	0.71	8	329	13791
	$8_{\mathbf{S}}$	0.74	0.69	7	233	10128	0.73	0.71	8	252	12164
	1_{B}	0.72	0.50	7	94	2782	0.63	0.44	11	71	3877
	2 _B	0.57	0.24	12	93	2930	0.40	0.30	28	71	4100
Ыd	3_{T}	0.73	0.48	7	63	1635	0.64	0.44	10	45	2219
nfe	4 _T	0.55	0.18	12	65	1669	0.36	0.30	29	48	2290
hre	$5_{\mathbf{S}}$	0.73	0.47	7	47	1750	0.63	0.41	11	44	1801
£	$6_{\mathbf{S}}$	0.55	0.19	12	46	1534	0.40	0.31	26	46	1765
	$7_{\rm S}$	0.70	0.44	8	48	1517	0.63	0.43	10	50	2056
	$8_{\mathbf{S}}$	0.55	0.18	12	46	1619	0.41	0.33	28	52	1985
	1_{B}	0.71	0.31	7	85	4406	0.67	0.48	9	85	5449
	2 _B	0.39	0.18	12	85	3730	0.50	0.39	16	87	4820
	3_{T}	0.73	0.31	7	57	1987	0.67	0.49	1	54	2800
ers	4_{T}	0.40	0.18	12	77	1837	0.47	0.42	20	55	2605
Mo	5s	0.72	0.31	7	54	2099	0.66	0.52	10	60	2770
L	6 _s	0.40	0.17	12	55	1730	0.56	0.41	14	65	2526
	7_{s}	0.07	0.31	7	68	2170	0.68	0.48	9	58	2526
	8 s	0.40	0.17	12	70	1870	0.50	0.38	18	50	2560

Table 4: Results of the developed classification models on the multi-labeled contributions

In multi-label classification tasks, the outputs of In addition to the discussed metrics, as in the 580 models are usually referred to as logits which are For Bonn, 585 will be critical since class predictions depend

> As the results show, macro F-scores in the ⁵⁹⁶ indicator of the models' better performance in the

⁵⁹⁷ classification of rare classes. This is an expected ⁶⁴⁸ 598 result in the classification of unbalanced datasets. 649 two text summarization models for the German 599 Moreover, contrary to the developed models on 650 language on argument data which was left 600 single-labeled texts, the base models have shown 651 undiscovered. The summarized texts using the sol slightly better performance over multi-labeled data. 652 multilingual mBART model achieved around 2% This slightly lower performance of shortened texts 653 higher performance in text classification against 602 603 is insignificant when compared to the reduction in 654 the BERT2BERT model which is specially fine-⁶⁰⁴ train time. For Bonn, the variation in the F-scores ⁶⁵⁵ tuned for summarization of the German language. 605 is less significant. However, for Ehrenfeld and 656 One of the potential reasons for the lower 606 Moers, due to the presence of very rare classes, the 657 performance of the BERT2BERT model is the 607 combination of weighting and AUC has shown a 658 shorter length of the summarized documents. 608 significant improvement in both the micro and 659 Shorter texts mean more loss of information which 609 macro F-scores. More detailed information on the 660 could be informative for the model. Although the 610 class-wise performance of the models can be found 661 quality of the summarized argument texts is not 611 in Appendix 1.

Discussion 612 6

613 Based on the discussions in the previous section, 614 we observed that the maximum length in the BERT 615 models is an effective parameter. However, to 616 provide the practitioners with a feasible argument 617 data analytic system, limitations such as the ability 668 The objective of this study was to explore the effect 618 to handle long texts, the requirement for less 669 of abstractive summarization models in analyzing 619 powerful systems, and the need for faster data 670 the German language argument data. Limited 620 analysis need to be addressed. Summarization is an 674 studies have worked on this topic and most of the 621 effective technique for shortening the contributions 672 existing ones are focused on the English language. ⁶²² and reducing the processing time while saving the ⁶⁷³ Hence, this study was an attempt to fill this gap. 623 key information. Based on the run-time values in 674 624 Table 3, summarization reduces the train time using 675 the German language (German BERT2BERT) and 625 CPU by 50% on average without significantly 676 a BART model were considered. The German e26 affecting the performance. Reduction in time eases 677 BERT2BERT fine-tuned on the MLSUM DE 627 the use of models on systems without special 678 dataset was a reasonable candidate as it was 628 specifications such as having a GPU. Moreover, the 679 specifically 629 results from the trained models on the data of the 680 summarization. The performance of the generated 630 city of Bonn showed that if a sufficient number of 684 summarized texts using this model in the 631 contributions is available under each label, 682 classification task was studied against the 632 truncation as a shortening technique performs close 683 summarized texts from the BART model. Pre-633 to the summarization technique. Thus, despite the 684 trained on a large number of texts from different 634 removed parts of the texts, the model still has 685 languages made the BART model a good training instances to learn 635 enough 636 distinguishing features of each class. On the 687 637 contrary, if the number of contributions is limited, 688 methods reduce the run-time of the model while ess especially under each category, like in the case of 689 keeping the classification quality almost the same. 639 Ehrenfeld, summarization preserves the text's key 690 Additionally, our results suggest that the generated 640 features and hence surpasses the performance of 691 summarized texts using the German BERT2BERT truncation in classification. 641

642 643 are valid in the classification of multi-labeled 694 summarization is still open to be discovered. 644 contributions. However, like other multi-label 695 Moreover, in following the long-term goal of 645 classification problems, deciding on the optimal 696 constructing a system for public contributions 646 value of the thresholds plays a significant role in 697 analysis, the need for developing a pre-trained 647 the performance of the models.

This study also compared the performance of 662 measured in this study, the results follow the findings of (Alcantara et al., 2023) on Wikipedia articles which confirmed that the mBART model 665 outperforms the BERT2BERT model in abstractive 666 German text summarization.

667 7 **Conclusion and Future Work**

For this purpose, a BERT model fine-tuned in developed for German text the 686 alternative to the German BERT2BERT model.

The results proved that text summarization ⁶⁹² model are briefer than the ones generated from the Our findings confirmed that similar conclusions 603 mBART model. Hence, the quality of the generated ⁶⁹⁸ model on the argument data is important.

699 Limitations

700 One of the obstacles in the analysis of German texts 751 701 despite being a high-resource language is the 752 702 shortage of high-quality labeled data. This 753 ⁷⁰³ limitation can either result in lower accuracy of the ⁷⁵⁴ 704 trained models or over-fitting due to the fine-tuning 755 Arana-Catania, M., van Lier, F.-A., Procter, R., 705 of the model on a specific in-reach dataset. 756 706 Regarding the analysis of argument data, only a 757 707 few labeled datasets exist. In this study, one of the 758 708 important steps is the evaluation of the generated 759 709 summarized texts. The first German Text 760 R 710 Summarization Challenge was held in 2019 at the 761 711 SwissText conference in which 100,000 texts 762 712 collected from Wikipedia along with their 763 713 reference summaries were analyzed (SwissText, 765 714 2019). However, to the best of our knowledge, this 715 is the only available and suitable collection for the ⁷⁶⁶ Romberg, J., and Escher, T. 2022. Automated Topic 716 evaluation of summarization systems, in which 768 ⁷¹⁷ each data is equipped with a referenced summary. ⁷⁶⁹ 718 Hence, for the summarization of argument texts, 770 719 there is a lack of a developed dataset to evaluate the $_{720}$ quality of models and generated summaries. This $_{772}^{772}$ $_{721}$ limitation also gets highlighted as there is a lack of $_{773}$ 722 benchmark studies in the summarization of 774 723 German argument data. Finally, the size of 775 Romberg, J. 2023. Mind the User! Measures to More $_{724}$ available datasets is limited. The used dataset in $_{776}^{777}$ 725 this experiment consisted of 3139 contributions, 777 726 while for some of the classes, the number of 778 727 contributions is very low. A limited number of 779 728 samples in a category compared to the other classes 780 729 not only raises the challenges of unbalanced 781 Balta, D., Kuhn, P., Sellami, M., Kulus, D., Lieven, C., 730 datasets but also negatively affects the model 782 731 performance due to ineffective training. 783

732 References

733 Rowe, G., and Frewer, L. 2004. Evaluating Public 787 Participation Exercises: A Research Agenda. 734 788 Science, Technology, and Human Values. 29, 4. 735 https://doi.org/10.1177/0162243903259197. 736 790 737 Romberg, J., and Escher, T. 2023. Making Sense of 791 Citizens' Input through Artificial Intelligence: A 738

792 Review of Methods for Computational Text 739 793 Analysis to Support the Evaluation of Contributions 740 794 in Public Participation. Digital Government: 741 795 Practice. 5. Research and 742 https://doi.org/10.1145/3603254. 743

797 744 Simonofski, A., Fink, J., and Burnay, C. 2021. 798 Supporting policy-making with social media and e-745 799 participation platforms data: A policy analytics 746 800 framework. Government Information Quarterly. 38, 747 801 3, 101590. 748

- 749 Lee, J., and Kim, S. 2014. Active Citizen E-Participation in Local Governance: Do Individual 750 Social Capital and E-Participation Management Matter? Proceedings of the Annual Hawaii International Conference on System Sciences. 2044-2053. https://doi.org/10.1109/HICSS.2014.259.
 - Tkachenko, N., Zubiaga, A., and Liakata, M. 2021. Citizen Participation and Machine Learning for a Better Democracy. Digital Government: Research and Practice. 2. https://doi.org/10.1145/3452118.
 - omberg, J., Mark, L., and Escher, T. 2022. A Corpus of German Citizen Contributions in Mobility Planning: Supporting Evaluation Through Multidimensional Classification. Proceedings of the Thirteenth Language Resources and Evaluation Conference, 2874–2883. Marseille, France.
 - Categorization of Citizens' Contributions: Reducing Manual Labelling Efforts Through Active Learning. In Proceedings of the Conference, 369–385. https://doi.org/10.1007/978-3-031-15086-9 24.
- 771 Egger, R., and Gokce, E. 2022. Natural Language Processing (NLP): An Introduction: Making Sense of Textual Data. https://doi.org/10.1007/978-3-030-88389-8 15.
 - Accurately Evaluate the Practical Value of Active Learning Strategies. In Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing, 996-1006. https://aclanthology.org/2023.ranlp-1.107.
 - and Krcmar, H. 2019. How to Streamline AI Application in Government? A Case Study on Citizen Participation in Germany. In Lindgren, I., et al. (Eds.), Electronic Government. EGOV 2019, Lecture Notes in Computer Science, vol. 11685. Springer, Cham. https://doi.org/10.1007/978-3-030-27325-5_18.
- 789 Romberg, J., and Conrad, S. 2021. Citizen Involvement in Urban Planning - How Can Municipalities Be Supported in Evaluating Public Participation Processes for Mobility Transitions? In Proceedings of the 8th Workshop on Argument Mining, ACL, Punta Cana, 88-99. https://aclanthology.org/2021.argmining-1.9.
- 796 Romberg, J. 2022. Is Your Perspective Also My Perspective? Enriching Prediction with Subjectivity. In Proceedings of the 9th Workshop on Argument Mining, International Conference on Computational Linguistics, Gyeongju, Republic of Korea, 115-125. https://aclanthology.org/2022.argmining-1.11.

802 Aruneshwari, R. R., Anandkumar, K. M., and Kavitha, D. 2024. A Comprehensive Review of Text 803

784

785

786

Summarization. AIP Conference Proceedings. 857 804 2802.1, 140003. 858 805

https://doi.org/10.1063/5.0186988. 806

- lcantara, T. H. M., Krütli, D., Ravada, R., and Hanne, 860 807 A T. 2023. Multilingual Text Summarization for 861 808
- Using Transformer German Texts Models. 862 809
- Information. 14. 6, 303. 810
- https://doi.org/10.3390/info14060303. 811
- Verma, P., and Verma, A. 2020. A Review on Text 865 812

864

- Summarization Techniques. Journal of Scientific 866 813
- 64, 251-257. 867 Research. 814
- https://doi.org/10.37398/JSR.2020.640148. 815
- 869 Talib, R. 2021. Multilingual Text Summarization using 816
- Deep Learning. International Journal 817
- Engineering Research and Advanced Technology. 7, 871 818
- 29-39.https://doi.org/10.31695/IJERAT.2021.3712. 872 819
- 820 Anguiano, E., Villaseñor-Pineda, L., Montes, M., and 874
- Rosso, P. 2010. Summarization as Feature Selection 821
- 822
- Proceedings of the Conference. 39-44. 876 823
- https://doi.org/10.1007/978-3-642-14770-8 6. 824
- 878 825 Dewi, K., and Sagala, R. 2018. Using Summarization
- to Optimize Text Classification. IOP Conference 826 880
- Series: Materials Science and Engineering. 407, 827
- https://doi.org/10.1088/1757-012157. 828
- 899X/407/1/012157. 829
- 883 830 Khan, B., Shah, Z., Usman, M., Khan, I., and Niazi, B.
- 2023. Exploring the Landscape of Automatic Text 831
- Summarization: A Comprehensive Survey. IEEE 832
- PP. Access. 833
- https://doi.org/10.1109/ACCESS.2023.3322188. 834
- 835 Scialom, T., Dray, P.-A., Lamprier, S., Piwowarski, B.,
- 889 and Staiano, J. 2020. MLSUM: The Multilingual 836
- Summarization Corpus. In Proceedings of the 2020 837 891
- Conference on Empirical Methods in Natural 838
- Language Processing (EMNLP), 839
- Computational Association for Linguistics. 893 840
- https://doi.org/10.18653/v1/2020.emnlp-main.647. 894 841
- Schroeter, R., Scheel, O., Renn, O., and Schweizer, P.-842 J. 2016. Testing the Value of Public Participation in 843
- Germany: Theory, Operationalization and a Case 844
- 898 Study on the Evaluation of Participation. Energy 845
- 899 Research æ Social Science. 13. 846
- https://doi.org/10.1016/j.erss.2015.12.013. 847
- 848 Karic, S., Althaus, M.-C., and Heissler, J. 2024. Digital and Multi-Channel Citizen Participation in 849
- 903 Germany: A Comprehensive Overview of Patterns, 850
- 904 Methods and Determinants. Raumforschung und 851 905
- Raumordnung | Spatial Research and Planning. 82. 852
- https://doi.org/10.14512/rur.2170. 853
- 854 Kwon, N., Shulman, S., and Hovy, E. 2006. 908 Multidimensional Text Analysis for eRulemaking. 855 In ACM International Conference Proceeding 856

Series, 151. 157-166. https://doi.org/10.1145/1146598.1146649.

- 859 Arguello, J., Callan, J., and Shulman, S. 2008. Recognizing Citations in Public Comments. Journal of Information Technology & Politics. 5, 49-71. https://doi.org/10.1080/19331680802153683.
- 863 Habernal, I., and Gurevych, I. 2016. Which Argument is More Convincing? Analyzing and Predicting of Web Arguments Using Convincingness Bidirectional LSTM. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, 1589-1599. Berlin, Germany. https://doi.org/10.18653/v1/P16-1150.
- of 870 Cardie, C., Farina, C., Aijaz, A., Rawding, M., and Purpura, S. 2008. A Study in Rule-Specific Issue Categorization for E-Rulemaking. In Proceedings of the Conference, 244-253. https://doi.org/10.1145/1367832.1367874.
- for Document Categorization on Small Datasets. In 875 Konat, B., Lawrence, J., Park, J., Budzynska, K., and Reed, C. 2016. A Corpus of Argument Networks: Using Graph Properties to Analyse Divisive Issues. In Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16), 3899-3906. https://aclanthology.org/L16-1617.
 - 882 Fierro, C., Pérez, J., Quezada, M., and Fuentes Bravo, C. 2017. 200K+ Crowdsourced Political Arguments for New Chilean Constitution. а https://doi.org/10.18653/v1/W17-5101.
 - 1-1. 886 Liebeck, M., Esau, K., and Conrad, S. 2016. What to Do with an Airport? Mining Arguments in the 887 German Online Participation Project Tempelhofer Feld. In Proceedings of the Workshop on Argument Mining, 144-153. https://doi.org/10.18653/v1/W16-2817.
 - 8051-8067. 892 Ker, S. J., and Chen, J. N. 2000. A Text Categorization Summarization Based on Technique. In Proceedings of the ACL-2000 Workshop on Recent Advances in Natural Language Processing and Information Retrieval: Held in Conjunction with the 38th Annual Meeting of the Association for Computational Linguistics - Volume 11 (RANLPIR '00). 79-83. Association for Computational Linguistics, USA.
 - https://doi.org/10.3115/1117755.1117766.
 - 902 Zhao, J., and Gui, X. 2017. Comparison Research on Text Pre-processing Methods on Twitter Sentiment IEEE PP. Analysis. Access. 1-1. https://doi.org/10.1109/ACCESS.2017.2672677.
 - 906 Hovy, E., and Lin, C. Y. 1998. Automated Text Summarization and the Summarist System. In 907 of Proceedings Workshop, 197-214. а https://doi.org/10.3115/1119089.1119121.

900

901

910 Syed, A., Gaol, F., and Matsuo, T. 2021. A Survey of 964

the State-of-the-Art Models in Neural Abstractive 965 911

Text Summarization. IEEE Access, 9, 13248-966 912 967

13265. 913

https://doi.org/10.1109/ACCESS.2021.3052783. 914

969 915 Chen, C., Yin, Y., Shang, L., Jiang, X., Qin, Y., Wang,

970 F., Wang, Z., Chen, X., Liu, Z., and Liu, Q. 2021. 916

- 971 bert2BERT: Towards Reusable Pretrained Language 917 972
- Models. arXiv preprint. 918

https://arxiv.org/abs/2110.07143. 919

974 920 GitHub 2022. repository. Available online: https://github.com/google-research/bert/. Accessed 921 on 9 September 2024. 922

923 Scialom, Thomas, Dray, Paul-Alexis, Lamprier, 977 A Additional Literature Review (2020). MLSUM: The Multilingual Summarization ⁹⁷⁸ NLP in Public Contributions Over the past 50 924 925 926 Natural Empirical Methods in

927 Processing (EMNLP) (pp. 928

929

930 931

- preprint 932
- https://doi.org/10.48550/arXiv.2408.10277. 933

935 2024. 936 course/chapter2/5?fw=tf. 937

Chi, Sun & Qiu, Xipeng & Xu, Yige & Huang, 991 technology. 938

939

940 v.1905.05583. 941

943 944 945 946 947 364. 948

949 950 951 952 953 954 955 956 2020). 957 https://doi.org/10.18653/v1/2020.acl-main.703. 958 Ma, Edward. (2019). NLP Augmentation. GitHub¹⁰¹¹ methods using pre-trained models (Khan et al., 959 960

962

963

Remi & Funtowicz, Morgan & Davison, Joe & Shleifer, Sam & Platen, Patrick & Ma, Clara & Jernite, Yacine & Plu, Julien & Xu, Canwen & Scao, Teven & Gugger, Sylvain & Rush, Alexander. (2020). Transformers: State-of-the-Art Natural Language Processing. Proceedings of the 2020 Conference on Empirical Methods in Natural 38-45. Language Processing (EMNLP), https://doi.org/10.18653/v1/2020.emnlp-demos.6.

973 SwissText. (2019). German Text Summarization Challenge. Retrieved from https://www.swisstext.org/2019/shared-task/germantext-summarization-challenge.html

968

Corpus. Proceedings of the 2020 Conference on 979 years, there have been growing discussions Language 980 regarding the public's role in determining policies. 8051-8067). 981 However, the majority of the existing research is https://doi.org/10.18653/v1/2020.emnlp-main.647. 982 limited to some basic analyses (Simonofski, Fink, Cukier, R. (2024). Increasing transformer token length 983 and Burnay, 2021). "Citizen Participation and with a Maximum Entropy Principle Method. arXiv 984 Machine Learning for a Better Democracy" is one arXiv: 2408.10277. 985 of the recent studies that highlight how effective 986 NLP is in getting over the difficulties of analyzing 934 Hugging Face. (2024). "Chapter 2: Transformer 987 public participation platforms (Arana-Catania et Models" NLP Course. Accessed September 10, 988 al., 2021). Automatic Text Summarization is one of https://huggingface.co/learn/nlp- 989 the branches of NLP and its importance has been ⁹⁹⁰ highlighted in recent years due to the advances in

Xuanjing. (2019). How to Fine-Tune BERT for Text 992 Text Summarization The Vast generation of Classification?.arXiv.https://doi.org/10.48550/arXi 993 information in the form of documents requires a ⁹⁹⁴ searching system to access usable information 942 Mutasodirin, Mirza & Prasojo, Radityo. (2021). 995 efficiently, save time, and convey the main and Investigating Text Shortening Strategy in BERT: 996 important information to the reader (Alcantara et Truncation vs Summarization. Proceedings of the 997 al., 2023). Text summarization is one of the International Conference on Advanced Computer 998 solutions. Early experiments applied the traditional Science and Information Systems (ICACSIS), 1-5. 999 text summarization technique such as weighting https://doi.org/10.1109/ICACSIS53237.2021.9631 1000 terms according to their frequencies (Ker & Chen, 1001 2000), as a feature selection strategy to improve Lewis, Mike & Liu, Yinhan & Goyal, Naman & 1002 classification performance (Anguiano et al., 2010, Ghazvininejad, Marjan & Mohamed, Abdelrahman₁₀₀₃ Dewi and Sagala, 2018). (Zhao & Gui, 2017) & Levy, Omer & Stoyanov, Veselin & Zettlemoyer, 1004 emphasized the preprocessing stage as a method of Luke. (2020). BART: Denoising Sequence-to-Sequence Pre-training for Natural Language mostly based on the new proportion of the torte Generation, Translation, and Comprehension.¹⁰⁰⁶ mostly based on the naïve properties of the texts Proceedings of the 58th Annual Meeting of the 1007 such as sentence position and word frequency Association for Computational Linguistics (ACL 1008 counts (Khan et al., 2023). With the advances in 7871-7880.1009 neural networks and deep learning, techniques 1010 focused more on abstractive summarization repository. https://github.com/makcedward/nlpaug. 1012 2023). The work of Hovy and Lin (1998) is one the ⁹⁶¹ Wolf, Thomas & Debut, Lysandre & Sanh, Victor & ¹⁰¹³ first ones in Automatic Text Summarization (ATS) Chaumond, Julien & Delangue, Clement & Moi,¹⁰¹⁴ (Hovy & Lin, 1998). They developed the Anthony & Cistac, Pierric & Rault, Tim & Louf, 1015 SUMMARIST system which benefited from

1016 language-specific techniques of parsing and
1017 semantic analysis combined with Information
1018 Retrieval (IR) and statistical methods.

1020 B Dataset Distribution

In deciding on the best document length, Figure 1
provides a clearer depiction of the contributions
length distribution. It can be observed that the
distributions are skewed to the left and the length
of the most of contributions is less than 128 tokens.
Figure 2 shows the distribution of the
summarized texts using the mBART model. It can
be seen that the lengths are more evenly
distributed.

Figure 3 illustrates the texts' length distribution after summarization using the BERT2BERT model. In comparison with the mBART model, the lengths are more skewed to the left. However, a similar pattern is observed in the summarized texts' distribution.

C Class-wise performance of the multilabeled models

Table 5 provides detailed information on the performance of the trained models on the multilabeled contributions. Table 6 shows the results that after detecting and applying the best thresholds. As optimal threshold has significantly improved the performance of models, especially in classifying tota contributions of less frequent classes. In Table 6 For Ehrenfeld, except for the lighting category, we have an improvement in the F-scores, especially when the weighting is applied.



Figure 1: Distribution of contributions length



Figure 2: Distribution of the length of the contributions after summarizing texts using the multilingual mBART model



Figure 3: Distribution of the length of the contributions after summarizing texts using the BERT2BERT model

	Experiments	Bicycle Parking	Cycle Path Quality	Cycling Traffic Management	Lighting	Misc	Obstacles	Signage	Traffic Lights	Micro F-score	Macro F-score	Hamming	Train run-time (second) GPU	Train run-time (second) CPU
	1 _B	0.921	0.777	0.809	0.833	0.044	0.769	0.56	0.824	0.774	0.692	0.059	383.8	17482
	2в	0.919	0.739	0.798	0.833	0.329	0.663	0.596	0.820	0.753	0.712	0.063	383.2	18378
	3 T	0.909	0.771	0.814	0.804	0.084	0.767	0.571	0.821	0.774	0.693	0.059	218.8	9955
nn	4 T	0.923	0.729	0.795	0.856	0.332	0.665	0.588	0.809	0.750	0.712	0.064	221.0	8940
\mathbf{B}_0	5s	0.904	0.769	0.816	0.829	0.063	0.765	0.553	0.817	0.773	0.689	0.059	219.8	9935
	6s	0.909	0.728	0.794	0.862	0.341	0.665	0.563	0.803	0.747	0.708	0.065	224.6	9702
	7s	0.898	0.768	0.799	0.800	0.043	0.746	0.515	0.806	0.759	0.672	0.063	233.8	9120
	8 s	0.898	0.717	0.784	0.839	0.270	0.642	0.540	0.805	0.735	0.687	0.067	233.0	10128
	1 _B	0.762	0.631	0.854	0	0	0.573	0.194	0.657	0.722	0.459	0.073	93.8	2782
	2 _B	0.561	0	0.751	0	0	0	0.279	0.295	0.571	0.236	0.116	93.2	2930
bla	3 T	0.779	0.629	0.859	0	0	0.568	0.313	0.697	0.734	0.481	0.070	62.6	1635
nfe	4 T	0.364	0	0.740	0	0	0	0.08	0.239	0.546	0.177	0.121	65.2	1669
hre	5s	0.780	0.632	0.855	0	0	0.542	0.247	0.685	0.727	0.467	0.072	47.2	1750
Е	6s	0.466	0	0.735	0	0	0	0.08	0.222	0.545	0.188	0.121	45.6	1534
	7s	0.669	0.546	0.847	0	0	0.512	0.293	0.666	0.701	0.441	0.079	47.6	1517
	8 s	0.350	0	0.736	0	0	0	0.160	0.233	0.545	0.184	0.121	45.6	1619
	1 _B	0	0.810	0.813	0	0	0	0	0.829	0.713	0.306	0.068	85.2	4406
	2 _B	0.533	0.062	0.559	0	0	0	0	0.250	0.388	0.175	0.123	85.4	3730
	3 _T	0	0.828	0.831	0	0	0	0	0.832	0.727	0.311	0.066	56.8	1987
)er9	4 T	0.533	0.048	0.582	0	0	0	0	0.258	0.404	0.177	0.121	76.8	1837
Mo	5s	0	0.799	0.815	0	0	0	0	0.861	0.716	0.309	0.069	53.8	2099
	6s	0.533	0.016	0.586	0	0	0	0	0.240	0.398	0.172	0.123	54.8	1730
	7s	0	0.796	0.806	0	0	0	0	0.849	0.071	0.306	0.070	67.6	2170
	8 s	0.533	0.016	0.583	0	0	0	0	0.258	0.399	0.173	0.121	70.0	1870

Table 5: Performance of the developed classification models on the multi-labeled contributions.

	Experiments	Bicycle Parking	Cycle Path Quality	Cycling Traffic Management	Lighting	Misc	Obstacles	Signage	Traffic Lights	Micro F-score	Macro F-score	Hamming	Train run-time (second) GPU	Train run-time (second) CPU
	1 _B	0.901	0.779	0.825	0.855	0.349	0.733	0.607	0.780	0.761	0.729	0.071	430.0	22606
	2 _B	0.908	0.741	0.808	0.723	0.506	0.712	0.599	0.799	0.754	0.725	0.071	426.8	22126
	3 T	0.913	0.778	0.823	0.879	0.383	0.718	0.581	0.754	0.760	0.729	0.072	336.0	11529
uu	4 _T	0.920	0.740	0.799	0.738	0.464	0.707	0.622	0.806	0.749	0.724	0.073	301.4	12300
BC	5s	0.917	0.766	0.818	0.829	0.354	0.725	0.550	0.780	0.750	0.717	0.075	256.2	13814
	6s	0.908	0.735	0.794	0.586	0.436	0.715	0.568	0.804	0.741	0.693	0.076	269.4	12743
	7s	0.900	0.762	0.812	0.856	0.306	0.710	0.537	0.777	0.739	0.708	0.080	329.4	13791
	8 s	0.902	0.730	0.790	0.783	0.419	0.689	0.567	0.762	0.732	0.705	0.078	251.6	12164
	1 _B	0.654	0.537	0.814	0	0.363	0.286	0.391	0.442	0.633	0.436	0.106	71.0	3877
	2 _B	0.426	0.319	0.641	0	0.449	0.36	0.094	0.148	0.403	0.304	0.279	71.4	4100
eld	3 T	0.599	0.573	0.816	0	0.387	0.38	0.316	0.442	0.642	0.439	0.100	45.0	2219
enfe	4 T	0.463	0.266	0.552	0	0.398	0.338	0.219	0.153	0.362	0.299	0.291	47.6	2290
hre	5s	0.495	0.554	0.822	0	0.277	0.287	0.358	0.453	0.633	0.406	0.107	44.4	1801
Ξ	6s	0.355	0.246	0.626	0	0.59	0.306	0.199	0.116	0.398	0.305	0.257	46.4	1765
	7s	0.483	0.492	0.796	0	0.345	0.445	0.363	0.552	0.628	0.431	0.104	49.6	2056
	8 s	0.473	0.307	0.646	0	0.44	0.299	0.147	0.290	0.409	0.325	0.283	52.0	1985
	1 _B	0.333	0.790	0.820	0.406	0.260	0.044	0.325	0.829	0.670	0.476	0.093	84.8	5449
	2 _B	0.533	0.469	0.699	0.180	0.230	0.165	0.217	0.597	0.504	0.386	0.165	87.2	4820
s	3 _T	0.223	0.804	0.823	0.311	0.249	0.292	0.345	0.856	0.670	0.488	0.095	54.2	2800
oer	4 T	0.333	0.612	0.636	0.179	0.179	0.190	0.337	0.875	0.468	0.418	0.197	54.6	2605
M	5s	0.466	0.822	0.837	0.381	0.252	0.160	0.406	0.822	0.663	0.518	0.100	59.8	2770
	6s	0.333	0.488	0.708	0.196	0.338	0.140	0.233	0.820	0.555	0.407	0.142	65.4	2526
	7s	0.333	0.809	0.834	0.530	0.138	0.121	0.363	0.737	0.682	0.483	0.087	58.4	2526
	8 s	0.333	0.490	0.683	0.215	0.298	0.134	0.243	0.627	0.495	0.378	0.183	50.4	2560

 Table 6: Performance of the developed classification models on the multi-labeled contributions.

 Thresholds are optimized using AUC.