# **CMTS: A Dataset and Benchmark for Text Summarization of Minority** Languages in China

**Anonymous ACL submission** 

#### Abstract

Minority languages in China, such as Tibetan, 002 Uyghur, and Traditional Mongolian, face significant challenges due to their unique writing systems, which differ from international standards. This discrepancy has led to a severe lack of relevant corpora, particularly for supervised tasks like text summarization. To address this gap, we introduce a novel dataset, Chinese Minority Text Summarization (CMTS), which includes 100,000 entries for Tibetan, and 50,000 entries each for Uyghur and Mongolian, specifically curated for text summarization tasks. Additionally, we propose a high-quality test set anno-013 tated by native speakers, designed to serve as a benchmark for future research in this domain. We hope this dataset will become a valuable resource for advancing text summarization in Chinese minority languages and contribute to the development of related benchmarks.

#### 1 Introduction

007

011

014

017

020

021

034

040

Recently, the rapid development of large language models (LLMs) has been fueled by the availability of high-quality pre-training data. However, these advancements have primarily benefitted highresource languages such as English and Chinese. In contrast, many languages with substantial user bases remain excluded due to the scarcity of suitable corpora, especially for specific tasks like text summarization. This exclusion poses challenges for both academic research and the practical application of AI technologies.

This paper focuses on underrepresented minority languages in China, including Tibetan, Uyghur, and Mongolian, which have rich linguistic and cultural significance but suffer from a lack of resources. Although these languages appear in multilingual datasets like OSCAR (Jansen et al., 2022) and CulturaX (Nguyen et al., 2024), quality issues limit their usefulness. As shown in Figure 1, there is a clear gap between the large speaker populations of these languages and the small amount of data



Figure 1: The relationship between population size and dataset size in OSCAR (y-axis, in MB) for various high-, middle-, and low-resource languages.

available in major corpora. Studies also reveal problems with data quality: Zhang et al. (2024) found that 34% of the Uyghur data in CulturaX contains Kazakh or Arabic texts, pointing to issues like language misidentification and noise. These challenges of data scarcity and quality undermine efforts to build effective natural language processing (NLP) systems for these communities.

042

044

045

046

047

048

051

054

056

059

060

061

062

063

064

065

Moreover, there is a complete lack of opensource datasets tailored for text summarization in these minority languages. This gap hinders the development of supervised methods and benchmarks for summarization tasks. To address this limitation, we introduce Chinese Minority Text Summarization (CMTS), a novel dataset specifically designed for text summarization in Tibetan, Uyghur, and Mongolian. CMTS consists of 100,000 Tibetan samples and 50,000 samples each for Uyghur and Mongolian. In addition to the main dataset, we collaborated with native speakers of these languages to further ensure data quality. From the existing data, we selected 3,000 samples for each language and conducted a detailed annotation process to evaluate the alignment and quality of the text summaries. These samples were reviewed by multiple annota067tors for each language, and only the data deemed068high-quality by consensus among native annota-069tors was retained. This subset of data provides a070reliable *benchmark* for future research, ensuring071consistency and reproducibility in evaluating sum-072marization models. By combining a large-scale073dataset with carefully curated benchmark samples,074CMTS bridges the gap in resources for text sum-075marization in Chinese minority languages.

In summary, this paper makes the following key contributions:

- We present **CMTS**, a novel and large-scale open-source dataset specifically designed for text summarization in three Chinese minority languages: Tibetan, Uyghur, and Mongolian. We release this dataset under MIT license.
  - We provide a carefully curated benchmark test set, annotated by native speakers, to ensure high-quality evaluation and support transparent, reproducible research in text summarization for Chinese minority languages.

By introducing CMTS, we aim to fill the resource gap and pave the way for advancing natural language processing research on underrepresented languages.

# 2 Data Sources

076

086

090

097

100

101

102

103

104

105

106

109

110

111

112

The Chinese Minority Text Summary (CMWS) dataset, proposed in this paper, is sourced from various online platforms in China, including government documents, broadcasts, and news articles (detailed list in Appendix A). We used web crawlers to collect the data, where the webpage title serves as the summary and the main text as the source content. To ensure data quality and reliability, we applied a thorough cleaning process, with the main methods outlined as follows:

- **Removal of Non-Textual Content**: We filtered out non-textual elements such as advertisements, pop-ups, navigation bars, and multimedia content (e.g., images, videos, and audio files). This ensured that only relevant text was retained.
- **Duplicate Detection and Removal**: We identified and removed duplicate entries to avoid redundancy in the dataset, which could potentially bias the summarization models.

• **Text Normalization**: We standardized the text by converting all characters to a uniform encoding format and removing any extraneous white spaces, special characters, or formatting inconsistencies.

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

• Language Purity Check: We conducted a language purity check to ensure that the collected texts were primarily in the target languages (Tibetan, Uyghur, and Mongolian). Non-target language content was filtered out to maintain linguistic integrity.

These steps were essential to ensure that the final dataset consisted of high-quality, relevant, and clean text that could be reliably used for text summarization tasks.

## 3 Annotation

After crawling the data, we further enhanced the quality of the three languages' evaluation set by native speaker annotation. For each language, we recruited two native speakers who helped assess the quality of the title-content matching. A total of 3000 randomly selected samples from the crawled data were annotated, with the task focused on evaluating how well the article titles matched the content of the articles.

# 3.1 Annotation Guidelines

The annotation process was designed to ensure the reliability and consistency of the evaluations. The annotators were provided with the following specific guidelines:

- Task Objective: Annotators were asked to assess the degree of relevance between the title and the content of the article and assign a score accordingly.
- **Title Evaluation:** Annotators were first asked to identify any issues between the title and the article content, such as:
  - *Incomplete Article:* whether the article content is incomplete, making it impossible for the title and content to match.
  - *Text Quality:* whether the title contains spelling, grammatical, or contextual errors that would significantly hinder its match with the article content.
  - Other Issues: whether there is any other noticeable discrepancies between the title and content, such as irrelevance.

Language	Size	Length In Token (Title/Content)	Length In Characters (Title/Content)
Tibetan	2901	12.3 / 376.7	74.0 / 1884.1
Mongolian	2931	27.2 / 429.8	134.3 / 2181.7
Uyghur	2950	30.2 / 815.7	91.5 / 3781.6

Table 1: Annotation Results for Valid Samples in the CMTS Dataset. The token length is counted using the CINO (Yang et al., 2022) tokenizer, while the character length is based on the raw character count for both title and content.

160	If no major issues were identified, the title-	• I
161	content match was considered "Normal."	ta
162	• Matching Score: Annotators were instructed	i
163	to rate the match based on how well the title	3.3
164	corresponded with the article content. The	т
165	scoring system was as follows:	work
166	- 1 point: Completely Mismatched (The	systen
167	title is entirely unrelated to the content).	2
168	– 2 points: Slightly Mismatched (The title	• 5
169	is related to the content but does not align	te
170	with the main theme).	
171	- 3 points: Slightly Inaccurate (There	
172	is some connection, but it is not fully	
173	aligned).	
174	– 4 points: Uncertain (The relationship	• A
175	between title and content is unclear or	n
176	ambiguous).	• A
177	– <b>5 points:</b> Slightly Matched (There is a	d
178	strong connection, but there are some	t
179	inconsistencies).	r
180	- 6 points: Well Matched (The title	<b>A</b>
181	matches the content with only minor dis-	Ani to the
182	crepancies).	consis
183	- 7 points: Fully Matched (The title per-	consis
184	fectly corresponds to the content).	3.4
185	3.2 Consistency and Quality Control	After
186	To ensure consistency and accuracy across anno-	or err
187	tations, multiple annotators evaluated each article.	sampl
188	The following steps were implemented to guarantee	numbe
189	the quality of the annotations:	eacn I
190	• Consistency Check: An annotation was con-	degree
101	sidered invalid if the score differed by more	had so

<sup>Consistency Check. An annotation was considered invalid if the score differed by more than 2 points from the majority of annotators. Additionally, if an annotator's judgment deviated significantly from the majority opinion (e.g., the majority rated the title as "matching," but the annotator rated it as "not matching"), the annotation would be discarded.</sup> 

193

194

195

196

197

• Handling Invalid Annotations: Invalid annotations were removed, and the annotators are incentivized to not produce such annotations. 198

199

201

202

203

204

205

206

207

208

210

211

212

213

214

215

216

217

218

219

221

222

223

224

225

226

227

228

229

230

231

232

#### 3.3 Incentive System

To encourage careful and consistent annotation work, we implemented a reward-based incentive system:

- Scores < 4 or > 4 are considered as different tendencies:
  - Scores < 4 indicate a *non-aligned* tendency.
  - Scores  $\geq$  4 indicate an *aligned* tendency.
- Annotation whose tendency aligns with the majority will receive **0.25 RMB**.
- Annotation that aligns with the majority tendency and furthermore deviates by no more than 1.5 points from the average score will receive an additional **0.25 RMB**.

Annotators were strongly encouraged to adhere to the guidelines to ensure the high quality and consistency of the dataset annotations.

#### 3.4 Annotation Results

After removing the data flagged as mismatched or erroneous by the annotators, we retained the samples with an average score above 4. The final number of valid samples and the average length for each language are shown in Table 1.

In general, the data we retained showed a high degree of quality. Most of the remaining samples had scores of 7, with a small number scoring 6. This indicates that the majority of the collected data is of high quality and well-suited for text summarization tasks. These results suggest that the annotation process, guided by native speakers, was effective in ensuring the reliability and relevance of the data. The average title and content lengths reflect the linguistic characteristics of these minority languages. For Tibetan, the average title is 12.3 tokens (74.0 characters) and the content is 376.7 tokens (1884.1 characters). Mongolian and Uyghur samples, however, show much longer lengths across both titles and contents.

#### 4 Experiment

234

235

236

240

241

242

243

244

245

246

247

248

249

251

252

258

259

260

261

262

263

264

265

267

269

271

272

273

277

278

279

In this section, we evaluate some of the most popular models available for Tibetan, Mongolian, and Uyghur on CMTS, including finetuning small encoder-decoder models and few-shot evaluation of LLMs.

#### 4.1 Experimental Settings

**Fine-tuned Models:** The small models, **cino-cum** (which uses the **cino** (Yang et al., 2022) encoder, based on the XLM-R model tailored for Chinese minority languages, and a transformer decoder in a seq2seq architecture) and **swcm** (Su et al., 2025) (which is based on the same structure as cino-cum, but incorporates shared weight optimization across the encoder and decoder for improved performance across languages), are fine-tuned on non-annotated data from the CMTS dataset. These models are then evaluated using high-quality annotated data to assess their summarization performance. The finetuning is conducted on raw, non-annotated data, while the evaluation is done using a set of annotated samples to measure the ROUGE-L scores.

Few-shot Models: The large models, Qwen2.5-72B (Yang et al., 2024) and LLaMA3.1-70B (Dubey et al., 2024) use a 2-shot learning paradigm, where two annotated samples are dynamically inserted as examples within the input of each annotated sample.

Detailed training configurations and hyperparameters are provided in Appendix B.

#### 4.2 High-Quality Small Sample Experiment

Given that evaluating large models like **Qwen2.5**-**72B** and **LLaMA3.1-70b** with nearly 3,000 annotated samples per language is resource-intensive, we also selected a high-quality subset for evaluation to facilitate future works. Specifically, we chose the top 500 annotated samples based on evaluation scores to create a high-quality small sample version, enabling more efficient performance assessment while maintaining data quality.

Model	Size	bo	mn	ug
cino-cum	411M	0.20	0.12	0.9
swcm	457M	0.23	0.18	0.15
Qwen2.5	72B	0.24	0.32	0.29
LLaMA3.1	70B	0.34	0.30	0.35

Table 2: Model Parameters and ROUGE-L F1 Scores across all annotated data

Model	bo	mn	ug
cino-cum	0.21	0.13	0.10
swcm	0.23	0.17	0.14
Qwen2.5	0.24	0.29	0.34
LLaMA3.1	0.34	0.31	0.34

Table 3: ROUGE-L F1 Score in High-Quality data

281

282

283

284

285

287

290

291

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

#### 4.3 Results and Discussion

The experimental results are summarized in Table 2, which presents the performance of the models on the CMTS dataset across the three languages: Tibetan (bo), Uyghur (ug), and Mongolian (mn). The fine-tuning results for the small models, **cino-cum** and **swcm**, show that both models achieved competitive ROUGE-L scores, demonstrating that finetuning with the CMTS dataset enables the models to effectively capture text summarization capabilities for all three languages. This indicates that the large amount of non-annotated data collected in the CMTS dataset plays a crucial role in enhancing model performance for these underrepresented languages.

For the large models, **Qwen2.5-72B** and **LLaMA3.1-70B**, the few-shot results, as shown in Table 2 and Table 3, reveal strong performance across both small and large sample tests. This demonstrates that the models exhibit high-quality summarization capabilities, regardless of the sample size, highlighting the effectiveness of using small annotated datasets for evaluating model performance. The ability of these models to perform well with just a few annotated samples supports the idea that the CMTS dataset, with its carefully curated annotated samples, can serve as a reliable benchmark for future research and evaluation in text summarization for minority languages.

Overall, both the fine-tuning and few-shot learning approaches contribute significantly to advancing text summarization for minority languages, and the CMTS dataset proves to be a valuable resource for further research in this area.

4

## 372 373 374 375 376 377 378 379 380 381 382 383 384 385 387 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424

425

369

370

371

### 5 Limitations

315

Despite the significant contributions of the CMTS 316 dataset in advancing text summarization for Chi-317 nese minority languages, several limitations remain. 318 First, while the CMTS dataset represents a substantial effort to address the data scarcity issue for Tibetan, Uyghur, and Mongolian, the availability 321 of high-quality linguistic resources for these lan-322 guages is still limited compared to high-resource languages like English and Chinese. The scarcity of large-scale annotated datasets for other minor-325 ity languages in China and beyond further highlights the need for continued efforts to expand the 327 scope of language resources. Additionally, the current dataset focuses primarily on text summariza-329 tion tasks, leaving other NLP applications underexplored. Future work will aim to address these lim-331 itations by expanding the dataset to include more minority languages and diversifying the types of 333 NLP tasks supported. We also plan to collaborate 334 335 with more native speakers and linguistic experts to enhance the quality and coverage of the dataset. By doing so, we hope to contribute to a more inclusive 337 and comprehensive development of NLP research for underrepresented languages. 339

#### References

341

343

346

351

356

357

361

362

364

366

367

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph

Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. The llama 3 herd of models. *CoRR*, abs/2407.21783.

- Tim Jansen, Yangling Tong, Victoria Zevallos, and Pedro Ortiz Suarez. 2022. Perplexed by Quality: A Perplexity-based Method for Adult and Harmful Content Detection in Multilingual Heterogeneous Web Data. *arXiv e-prints*, page arXiv:2212.10440.
- Thuat Nguyen, Chien Van Nguyen, Viet Dac Lai, Hieu Man, Nghia Trung Ngo, Franck Dernoncourt, Ryan A. Rossi, and Thien Huu Nguyen. 2024. CulturaX: A cleaned, enormous, and multilingual dataset for large language models in 167 languages. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 4226– 4237, Torino, Italia. ELRA and ICCL.
- Zeli Su, Ziyin Zhang, Guixian Xu, Jianing Liu, XU Han, Ting Zhang, and Yushuang Dong. 2025. Multilingual encoder knows more than you realize: Shared weights pretraining for extremely low-resource languages.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. arXiv preprint arXiv:2407.10671.
- Ziqing Yang, Zihang Xu, Yiming Cui, Baoxin Wang, Min Lin, Dayong Wu, and Zhigang Chen. 2022. CINO: A chinese minority pre-trained language model. In Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022, pages 3937–3949. International Committee on Computational Linguistics.
- Chen Zhang, Mingxu Tao, Quzhe Huang, Jiuheng Lin, Zhibin Chen, and Yansong Feng. 2024. Mc<sup>2</sup>: Towards transparent and culturally-aware NLP for minority languages in china. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 8832–8850. Association for Computational Linguistics.

426	A Dataset Details	Task:
427	A.1 1. Dataset Size	Content: Passage
428	Table 4 shows the size of annotated and non-	The model generates a concise and accurate sum-
429	annotated data for each language in the CMTS	mary based on the given passage and title.
430	dataset.	
431	A.2 2. List of Crawled Websites	
432	Table 5 lists the websites and URLs used for data	
433	crawling.	
434	<b>B</b> Training Details	
435	Fine-tune Training Details	
436	Hardware: NVIDIA A5000 GPU, 24 GB RAM,	
438	Intel i7 CPU.	
439	Software: Ubuntu 20.04, CUDA 11.7, PyTorch	
440	2.3	
441	Training Configurations	
443	Local Batch Size: 20	
444	<b>Gradient Accumulation Steps:</b> 4	
445	Global Batch Size: 80	
446	Epochs: 50	
447	<b>Optimizer:</b> AdamW with $\beta_1 = 0.9$ , $\beta_2 =$	
448	0.999	
449	Learning Rate: 1e-4	
450	Warm-up: Linear warm-up for the first epoch,	
451	gradually increasing the learning rate from 1e-5 to	
452	1e-4.	
453	Few-shot Training Details	
458	In the few-shot setting, the model is provided	
456	with a prompt and a few examples to generate a	
457	task-specific output. For this task, the prompt is	
458	designed to help the model generate concise and	
459	accurate summaries in Tibetan, Uygnur, or Mon-	
400	The examples are structured to guide the model's	
462	behavior in generating the expected output.	
463	Prompt	
465	Based on the provided passage with title and	
466	content, generate a concise and accurate summary	
467	in {Tibetan/Uyghur/Mongolian}:	
468	Example 1/2:	
469	Content: Passage	
470	Title: Title of the passage	
471	Example 2/2:	
472	Content: Passage	
473	Title: Title of the passage	

# 

Language	Annotation Size	Non-Annotation Size
Tibetan (bo)	2901	100,000
Mongolian (mn)	2931	50,000
Uyghur (ug)	2950	50,000

Table 4: Dataset size for each language in CMTS.

Website Name	URL	Language
Qinghai Lake Website (Tibetan Version)	https://www.amdotibet.cn	BO
China Tibet News Network	https://tb.xzxw.com	BO
Bon Religion Website	http://www.himalayabon.com	BO
Kamba Satellite TV Network	http://tb.kangbatv.com	BO
Qinghai Tibetan Language Radio and TV Station	http://www.qhtb.cn	BO
China Tibetan Calligraphy Website	http://www.zgzzsfw.com	BO
Inner Mongolia Government Website	https://mgl.nmg.gov.cn	MN
Hulunbuir City Government Website	http://mgl.hlbe.gov.cn	MN
Xilingol League Government Website	http://mgl.zlq.gov.cn	MN
Ula'gae Government Website	http://mgl.wlgglq.gov.cn	MN
Chifeng City Government Website	http://mgl.chifeng.gov.cn	MN
Tongliao City Government Website	http://mgl.tongliao.gov.cn	MN
Aksu News Network	https://uy.aksxw.com	UG
Nur Network	https://www.nur.cn	UG
Tianshan Net	http://uy.ts.cn	UG
Xinjiang Government Website	https://uygur.xinjiang.gov.cn	UG
Xinjiang Daily Website	http://xjrbuy.ts.cn	UG

Table 5: List of websites used for data crawling.