# StRuCom: A Novel Dataset of Structured Code Comments in Russian

Anonymous ACL submission

#### Abstract

002 Structured code comments in *docstring* format are essential for code comprehension and maintenance, but existing machine learning models for their generation per-006 form poorly for Russian compared to English. To bridge this gap, we present StRu-Com - the first large-scale dataset (153K)examples) specifically designed for Russian code documentation. Unlike machinetranslated English datasets that distort terminology (e.g., technical loanwords vs. 012 literal translations) and docstring struc-013 tures, StRuCom combines human-written 014 comments from Russian GitHub reposi-016 tories with synthetically generated ones, ensuring compliance with Python, Java, 017 018 JavaScript, C#, and Go standards through automated validation.

#### 1 Introduction

007

The automated generation of structured code comments in *docstring* format, including detailed descriptions of functionality, parameters, return values, exceptions, and usage ex-024 amples, greatly improves codebase mainte-026 nance. Structured code comments provide developers with quick and easy access to the 027 required information, and can also be used to automatically generate project documentation, for instance, in HTML format. However, modern language models, such as Qwen2.5-Coder (Hui et al., 2024) and DeepSeek-Coder (Guo et al., 2024), primarily focus on Englishlanguage data and therefore perform poorly for Russian-language comment, neglecting the needs of Russian-speaking developers. These developers, working on localized projects, who 037 often encounter linguistic barriers, which can lead to code misunderstanding and a waste of time. In view of this, there is a strong need for a specialized model for this task, which requires curated training data.

041

042

043

044

045

046

047

050

051

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

Unfortunately, existing datasets (Englishcentric CodeSearchNet (Husain et al., 2019) or multilingual MCoNaLa (Wang et al., 2023b)) mostly focus on code summarization and retrieval tasks, not on function-level documentation generation. The datasets that contain both simple comments and docstrings in English (for example, the Vault (Nguyen et al., 2023)), firstly, require a tool for structurebased filtration to check comments for existence of detailed functionality descriptions, covering all function parameters, exceptions and its return value. Secondly, machine translation of English comments cannot be straightforwardly used, as it introduces distortions (e.g., translating "endpoint" as "конечная точка" instead of the established loanword "эндпоинт") (Wang et al., 2023b) and disrupts *docstring* structure.

In this work, we present StRuCom, the first specialized dataset for generating structured Russian-language code comments. To create it, we developed a tool for filtering and validating comment structures, supporting five popular documentation styles: Python - GoogleDoc<sup>1</sup>, JavaScript - JSDoc<sup>2</sup>, Java -JavaDoc<sup>3</sup>, C# - XML<sup>4</sup>, and Go - GoDoc<sup>5</sup>. The dataset combines real-world comments from Russian repositories with synthetically generated examples. Using this data, we finetuned the Qwen2.5-Coder model family (0.5B, 1.5B,

<sup>&</sup>lt;sup>1</sup>https://google.github.io/styleguide/ pyguide.html

<sup>&</sup>lt;sup>2</sup>https://jsdoc.app

<sup>&</sup>lt;sup>3</sup>https://docs.oracle.com/javase/8/docs/ technotes/tools/windows/javadoc.html <sup>4</sup>https://learn.microsoft.com/en-us/

dotnet/csharp/language-reference/xmldoc/ recommended-tags

<sup>&</sup>lt;sup>5</sup>https://tip.golang.org/doc/comment

3B, and 7B parameters), demonstrating statistically significant improvements in generation quality via chrf++ (Popović, 2017) and BERTScore (Zhang et al.) metrics compared to baseline versions.

Our contributions: Filtering tool for structured comments. We developed an automated tool to validate comment structures across five documentation standards (Python, Java, Go, C#, JavaScript). Dataset. We compiled a dataset of 153K Russian-language code-comment pairs, combining real-world examples from GitHub repositories with synthetically generated annotations for five programming languages.

#### 2 Related Work

074

075

076

077

087

880

090

096

100

101

102

103

104

105

106

108

109

110

111

112

113

114

115

116

117

118

119

120

121

The existing datasets for code-to-text tasks are mainly focused on English-language content. **The Stack** (Kocetkov et al., 2022) combines multilingual code from 658 programming languages (67 TB in version 2x), collected from a variety of sources: Software Heritage Archive, GitHub Issues, Stack Overflow, etc. Despite its scale, the set is not adapted for supervised finetuning (SFT) tasks and requires significant preprocessing. The Vault (Nguyen et al., 2023), derived from The Stack v1, includes 43 million English-language code-text pairs from 10 programming languages. The data was obtained by extracting docstrings and inline comments using the *Code-Text* parser  $^{6}$ . However, structured comments (with parameters and usage examples) remain rare, which is partly explained by the predominance of short functions in the source data. CodeSearchNet (Husain et al., 2019), part of the CodeXGLUE benchmark (Lu et al., 2021), contains 1 million English-language code-text pairs for 6 languages. The set is focused on code search: text descriptions are limited to the first paragraphs of the documentation, which simplifies comparison, but excludes complex descriptions. MCoNaLa (Wang et al., 2023b) offers limited multilingual support: 345 Russian, 341 Spanish, and 210 Japanese intent-snippet pairs for Python. The focus on narrow "how-to" scenarios and a small size limit the applicability of this dataset for structured documentation

tasks.

#### 3 StRuCom Dataset

Collection **Process.** To construct our dataset, we crawled all existing Russianlanguage repositories on GitHub for the selected programming languages (Python, Java, JavaScript (JS), C#, and Go). Since the GitHub API does not provide a direct query to identify the natural language used by repository authors, we developed a novel approach to address this limitation. Our program retrieved repositories with Russian-language descriptions and permissive licenses (allowing commercial use or lacking licensing restrictions). The crawled repositories contained comments written in various languages. For details on comment extraction see Appendix A.

Filtration Process. At the initial stage of filtering, all comments were standardized to follow a uniform style based on the conventions established for each programming language: Python - GoogleDoc, JavaScript - JS-Doc, Java - JavaDoc, C# - XML, and Go -GoDoc. Examples of these standardized formats can be seen on Fig. 1. To further divide comments into types by structure, we suggest the following terminology: A structured com*ment* is a comment that can be parsed by the docstring\_parser library<sup>7</sup> and contains either parameter lists, return value descriptions, or exception descriptions. A complete com*ment* is a structured comment that provides a comprehensive description of all its component parts, including types (if needed). Anincomplete comment is a structured comment that lacks a description of any of its component parts, which is why it cannot be called complete. Unstructured comments are those that do not correspond to a specific format used in a given programming language. For more information about filtration by structure see Appendix D. Only structured and complete comments were included in the final version of the dataset.

**Enhancement with LLM.** Based on the statistics on the structuredness of the collected data from GitHub, many code comments are incomplete or unstructured and gen-

123

124 125

126

127

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

159

160

161

162

163

164

165

166

167

168

169

<sup>&</sup>lt;sup>6</sup>https://github.com/FSoft-AI4Code/ CodeText-parser/tree/main

<sup>&</sup>lt;sup>7</sup>https://github.com/nmd2k/docstring\_parser





(c) C # XML comment style



// NameOfFunction description



Figure 1: Comparison of documentation styles in different programming languages

3

erally of poor quality. For some program-170 ming languages (for example, JavaScript and 171 Python), there is very little data and this is 172 not enough to finetune neural networks. To 173 solve these problems, we used large language 174 models (LLM), generating synthetic data using them in two ways: generating comments from 176 scratch and improving existing comments. For additional information about comment's enhancement see Appendix E. 179

Dataset Overview Tab. 1 presents the fi-180 nal statistical data of the final set, com-181 bining synthetic improved by the Miqu-70B model comments and generated from scratch by Qwen2.5-Coder-32B-Instruct ones with real comments from more than 150,000 Russianlanguage GitHub repositories of five program-187 ming languages: Python, Java, Go, C# and The total amount of data is JavaScript. 153,181 examples, of which 79,548 are improved, 65,914 are synthetic, and 7,719 are real comments. 191

Prog. lang.	Enhanced	From scratch	Real
Python	$14,\!625$	10,078	359
Java	$16,\!283$	$10,\!536$	$2,\!619$
Go	7,278	20,339	232
$\mathrm{C}\#$	39,715	$5,\!617$	$4,\!435$
JavaScript	$1,\!647$	$19,\!344$	100
$\sum$	79,548	65,914	7,719

Table 1: Statistics of the collected Russianlanguage data on programming languages and methods of obtaining them. The table shows the amount of improved (modification of existing comments by the Miqu-70B model), generated from scratch (synthetic data from Qwen2.5-Coder-32B-Instruct) and real comments.

The uniqueness of the proposed dataset is determined by several factors (see Table 2). Firstly, this is the first large corpus with Russian-language documentation for functions. The only existing dataset with comments in Russian, MCoNaLa, is designed to solve a different problem - searching for a code snippet based on the user's intent and, there-

196

197

199

Feature	CSN	Vault	MCoNaLa	Our dataset
#Pairs «code-text»	$6.5\mathrm{M}$	43K	341 - es, 210 - ja, 345 - ru	153K
Code format	Functions	Functions, classes, snippets	Code snippets	Functions
Text	Unstr.,	Mixed (unstr. and str. $w/o$	Unstr.,	Str. complete
format	1-2 sent.	filtration by structure)	(1-2  sent.)	(>5  sent.)
Progr. lang.	Go, Java, PHP, JavaScript, Python, Ruby	Java, JavaScript, Python, Ruby, Rust, Golang, C#, C++, C, PHP	Python, Java, JavaScript	Java, Python, C#, Go, JavaScript
Nat. lang.	en	en	ru, ja, es	ru
Data source	GitHub	The Stack	Stack Overflow	GitHub

Table 2: Comparison of the characteristics of the proposed dataset with existing analogues (CSN, Vault, MCoNaLa) by key parameters. The table shows the amount of data, the formats of code and text representation, the coverage of programming languages, linguistic features and data sources. The dataset we propose stands out with a strict focus on Russian-language structured comments on functions (153 thousand pairs), which contrasts with English-language counterparts operating with unstructured or mixed comments.

fore, is not suitable for generating structured 201 comments in the *docstring* style. Secondly, our dataset was strictly checked for structure and completeness: all comments were modified to one of the formats used in the industry for each specific programming language. In 205 other datasets, either there are no structured comments at all (MCoNaLa, CodeSearchNet), or they have not been filtered by structure (the Vault). Thirdly, as a result of the addition of synthetic data, the proposed set, unlike 210 MCoNaLa, has a sufficient size to train large 211 language models for all five selected program-212 ming languages. 213

#### 4 Experimental Evaluation

214

215

216

217

218

219

221

222

We conducted experiments, where we first benchmark existing open-source code-specific LLMs of different size (Qwen2.5-Coder (0.5B -7B) and DeepSeek-Coder (1.3B - 6.7B)), then finetune Qwen2.5-Coder (0.5B - 7B) on 7,500 comments, sampled from a synthetic part of our dataset and evaluate all models on our test set, 500 comments, sampled from real comments.

Evaluation We evaluated the models using standard natural language generation metrics, including ChrF++ (Popović, 2017) and a modified BERTScore (Zhang et al.). Instead of the traditional BERT (Kenton and Toutanova, 2019), we employed E5-Mistral 7B (Wang et al., 2022, 2023a), which offers superior performance for Russian, outperforming BERT models. 230

231

233

234

235

236

237

239

240

241

242

243

244

246

247

248

249

250

251

252

253

254

255

256

**Training and Results** The additional information about training setup, hyperparameters, etc. is located in Appendix F. Finetuning on the proposed dataset significantly improves the quality of comment generation using the BERTScore metric for all model sizes and most languages. For chrf++, significant improvements are observed in small number of cases. The results confirm that the proposed approach is effective for adapting language models to the task of generating Russian-language comments, especially in terms of semantic correctness (BERTScore).

#### 5 Conclusion

In this paper, we have developed a tool for filtering structured comments, collected a dataset of 153 thousand Russian-language code-comment pairs (real and synthetic data for 5 programming languages). We plan to expand the dataset by adding other programming languages, and develop and implement a quality criterion for structured code comments to automatically filter data and therefore improve the quality of the dataset.

## 257

6

Limitations

always true.

References

the dataset's quality.

Deepseek-coder:

arXiv:2106.09685.

arXiv:2409.12186.

Mikolov. 2016.

Minnesota.

arXiv:1612.03651.

of semantic code search.

text classification models.

The study has several limitations, including a

specific commenting style limitation, an imbal-

anced test dataset, and the assumption that

code comments always contain useful infor-

mation about code functionality, which is not

from GitHub may be redundant, uninforma-

tive, or contain errors, negatively impacting

Daya Guo, Qihao Zhu, Dejian Yang, Zhenda

Xie, Kai Dong, Wentao Zhang, Guanting

Chen, Xiao Bi, Yu Wu, YK Li, et al. 2024.

model meets programming-the rise of code in-

telligence. arXiv preprint arXiv:2401.14196.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan

Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang,

Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun

Zhang, Bowen Yu, Keming Lu, et al. 2024.

Qwen2. 5-coder technical report. arXiv preprint

Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Mil-

tiadis Allamanis, and Marc Brockschmidt. 2019.

Codesearchnet challenge: Evaluating the state

Armand Joulin, Edouard Grave, Piotr Bojanowski,

Armand Joulin, Édouard Grave, Piotr Bojanowski,

and Tomáš Mikolov. 2017. Bag of tricks for ef-

ficient text classification. In Proceedings of the

15th Conference of the European Chapter of the

Association for Computational Linguistics: Vol-

Jacob Devlin Ming-Wei Chang Kenton and

training of deep bidirectional transformers for

language understanding. In Proceedings of

naacL-HLT, volume 1, page 2. Minneapolis,

Denis Kocetkov, Raymond Li, Loubna Allal, Jia Li,

Chenghao Mou, Carlos Ferrandis, Yacine Jer-

nite, Margaret Mitchell, Sean Hughes, Thomas

Wolf, Dzmitry Bahdanau, Leandro Werra, and

ume 2, Short Papers, pages 427–431.

Lee Kristina Toutanova. 2019.

Matthijs Douze, Hérve Jégou, and Tomas

Fasttext.zip:

Compressing

Bert: Pre-

arXiv preprint

Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,

and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. arXiv preprint

Additionally, code comments

When the large language

## 2

2

26

- 262
- 26
- 265
- 26
- 267
- 2
- 2
- 2

273

- 27
- 2<sup>.</sup> 2<sup>.</sup>
- 2
- 2

281 282

- 2
- 286 287

288 289

291 292 293

294 295

- 297 298
- 299

302 303

3(

305 306 307

307 308 Harm Vries. 2022. The stack: 3 tb of permissively licensed source code.

- Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, MING GONG, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie LIU. 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation. In *Proceedings* of the Neural Information Processing Systems Track on Datasets and Benchmarks, volume 1.
- Dung Nguyen, Le Nam, Anh Dau, Anh Nguyen, Khanh Nghiem, Jin Guo, and Nghi Bui. 2023. The vault: A comprehensive multilingual dataset for advancing code understanding and generation. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, pages 4763–4788, Singapore. Association for Computational Linguistics.
- Maja Popović. 2017. chrf++: words helping character n-grams. In Proceedings of the second conference on machine translation, pages 612–618.
- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pages 3505–3506.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*.
- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023a. Improving text embeddings with large language models. arXiv preprint arXiv:2401.00368.
- Zhiruo Wang, Grace Cuenca, Shuyan Zhou, Frank F. Xu, and Graham Neubig. 2023b. MCoNaLa: A benchmark for code generation from multiple natural languages. In *Findings of* the Association for Computational Linguistics: EACL 2023, pages 265–273, Dubrovnik, Croatia. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

## A Comment Extraction

To extract comments, we used the *func* $tion_parser^8$  tool for Python, Java, and Go. 309 310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

336

337

338

339

340

341

348

346

349 350 351

352

353

354

355

357

358

360

361

<sup>&</sup>lt;sup>8</sup>https://github.com/ncoop57/function\_parser

For JavaScript and C#, we employed *Code*-Text. The GitHub data collection process con-364 sisted of several steps. First, code snippets from Python and JavaScript libraries with very few non-English comments were excluded. The formatting of comments in Java, JavaScript, and C# was then standardized. In C#, XML tags such as <summary> were corrected. For Java and JavaScript, redundant whitespaces, 371 line breaks in block comments (delimited by 372 /\*\* and \*/), and HTML tags were removed. 373 Next, automatically generated comments in 374 C# and JavaScript were filtered out. Dupli-375 cate comments in the function and docstring columns were eliminated, along with dupli-377 cates based on function and docstring independently. The language of each comment was then identified using Lingua<sup>9</sup>. More information about language identification methods 381 that we used is in Appendix B. If Lingua failed to determine the language, the corresponding comments were excluded from the dataset. To improve language identification accuracy, Lingua was provided with short descriptions of 387 comments, ensuring tags and identifier names that could degrade identification quality were removed. This process was applied to all programming languages except Go, which has a relatively simple comment structure.

The final dataset, after filtering, is summarized in Table 3. The results show that JavaScript and Go are characterized by a similar trend: a high proportion of commented repositories (70.8% and 55.9%) and functions (70.2% and 25.8%) are combined with a low percentage of Russian-language comments (24.0% and 16.4%), which may indicate the predominance of English-language documentation in their ecosystems. On the contrary, Python and C# show an increased proportion of Russian–language comments (49.2%) and 36.4%), which is probably due to regional development practices - the active participation of Russian-speaking communities in projects in these languages, where comments are often written in their native language for the local context.

397

398

400

401

402

403

404

405

406

407

408 409

## **B** Language identification

We applied two language identification methods to determine the language of the comments: FastText (Joulin et al., 2017, 2016) and Lingua. FastText uses a bag-of-n-grams approach to capture partial word order information, enabling efficient processing of large datasets on consumer hardware. Its pretrained models can classify text into one of 217 supported languages with high speed and efficiency. Lingua, on the other hand, employs a probabilistic n-gram model combined with rule-based heuristics, focusing on achieving high detection accuracy across 75 supported languages. While FastText offers broad language coverage and high efficiency, it demonstrated high precision but low recall for identifying Russian comments, frequently misclassifying them as less popular languages. Lingua, although slower and more memory-intensive, excels at handling short text and mixedlanguage inputs, which are common in code comments where natural language often intermixes with programming-specific syntax (e.g., tags and identifier names). Lingua's robustness in these scenarios makes it a preferable choice for detecting natural language within code comments.

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

## C Comment Structure

The examples of comment structure for five selected programming languages are shown in Figure 1. Notably, Python's GoogleDoc and JavaScript's JSDoc are the only styles among the selected ones that require explicit descriptions of parameter types and return types, reflecting the dynamically-typed nature of these languages. JSDoc shares stylistic similarities with JavaDoc, emphasizing structured documentation. By contrast, C# utilizes XML for comment formatting, providing a more tagbased approach. GoDoc stands apart with its flexible and descriptive style, as it imposes no strict format requirements, allowing developers to use a nearly free-form commentary approach.

## D Filtration by structure

For filtration-by-structure stage, we utilized the fork of  $docstring\_parser$  library  $^{10}$  and

<sup>&</sup>lt;sup>9</sup>https://github.com/pemistahl/lingua-py

<sup>&</sup>lt;sup>10</sup>https://github.com/rr-/docstring\_parser

Programming	#Repositories			#Functions			# Comments			
languago	With	Total	07	With	Total	07	in Pussion	Total	% in	
language	comments	Total	/0	comments	10tai	/0	in Russian	Total	Russian	
Python	18,535	64,440	28.8%	305,187	1,627,726	18.7%	150,255	$305,\!187$	49.2%	
Java	13,525	42,271	32.0%	409,506	$2,\!684,\!650$	15.3%	$98,\!622$	409,506	24.1%	
Go	2,592	$4,\!639$	55.9%	$117,\!691$	456, 347	25.8%	19,276	$117,\!691$	16.4%	
$\mathrm{C}\#$	8,858	26,329	33.6%	291,142	596,905	48.8%	106,058	291,142	36.4%	
JavaScript	15,073	$21,\!291$	70.8%	129,767	$184,\!871$	70.2%	31,084	129,767	24.0%	

Table 3: Statistics on data collection from GitHub, including analysis of repositories, functions, and comments on programming languages, grouped into three categories: **repositories** (the total number of repositories for each programming language, the number of at least one comment, and the percentage of the latter), **functions** (the total number of functions, the number of functions with comments and their relative proportion) and **comments** (the total number of comments, the number of Russian-language comments and their percentage).

*javalang*<sup>11</sup> tools to extract information about 458 comment structure and *Code-Text* to gather 459 information about code structure. We also 460 added missing types in Python comments 461 where possible using *Code-Text*. The dataset's 462 463 collection showed significant differences in 464 structured comments' availability and completeness across programming languages, as 465 summarized in Table 4. The results demon-466 strate an inverse relationship between the 467 complexity of the commenting standard and 468 the proportion of complete structured com-469 ments. Go, with minimal requirements (only 470 471 the function name at the beginning of the comment), shows the maximum percentage 472 of full comments (56.4%, 10,880). On the 473 contrary, Python and JavaScript, where stan-474 dards require specifying types and complex 475 476 annotations, have an extremely low proportion of complete comments (1.5% and 1.4%), 477 with unstructured ones dominating (94,968 478 and 14,091). Java and C++ with moderately 479 complex standards occupy an intermediate po-480 sition: 29.8% and 22.7% of full comments, re-481 spectively, but a significant number of unstruc-482 tured (48,347 and 30,188). The table con-483 firms that the simpler the syntax of a struc-484 tured comment, the higher the proportion of 485 486 its compliance. The extremely high Go score is explained by the simplified standard, and 487 the low Python/JavaScript values are due to 488 489 the excessive complexity of the requirements, which leads to a preference for unstructured 490 491 comments.

## E Enhancement of comments via LLM

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

The final dataset includes only those data with the length of both the code and the comment ranging from 250 to 1,000 characters. Very short comments and functions were excluded, as the goal was to create a dataset with detailed and comprehensive documentation. Very long comments or features are outliers and therefore were not considered. Comments were generated from scratch using the Qwen2.5-Coder-32B-Instruct model for functions without comments (see Table 3) and for functions, which comments were not successfully enhanced. To improve the dataset, the MIQU 70B<sup>12</sup> model was used, which was further trained in Russian. The goal of the improvement is to generate a complete and detailed comment of the best quality based on the function and the existing comment on it. An example is illustrated in figure 2. Candidates for improvement were selected from all the structuredness groups that were not included in the dataset in the "real" group. Comment is considered improved if it has become complete as a result of the improvement. Table 5 shows statistics on improving the dataset. Go stands out for the maximum efficiency of improvements (avg = 84.3%), especially for complete comments (91.5%), which is explained by a simple commenting standard, where it is enough to specify the function name. Python and JavaScript show the lowest averages (31.9% and 33.5%), which is due to the complexity of their standards, which re-

<sup>&</sup>lt;sup>11</sup>https://github.com/c2nes/javalang

<sup>&</sup>lt;sup>12</sup>https://huggingface.co/miqudev/miqu-1-70b

Programming	St				
languago	% complete out	Complete	Incomplete	Non-structured	
language	of all Russian	Complete	mcomplete		
Python	1.5%	$2,\!176$	$30,\!115$	94,968	
Java	29.8%	29,367	12,221	$48,\!347$	
Go	56.4%	$10,\!880$	-	8,396	
$\mathrm{C}\#$	22.7%	24,017	$41,\!898$	$30,\!188$	
JavaScript	1.4%	431	$1,\!484$	14,091	

Table 4: The structure of Russian-language comments on programming languages. For each language, the following are indicated: the percentage of complete structured comments out of the total number of Russian-language comments (% of the total number), the absolute values of complete and incomplete structured comments, as well as the number of unstructured ones. In Go, the dash in the "Incomplete" column is due to a feature of the commenting standard: comments are considered complete if they begin with the function name, which excludes the "incomplete" category.



Figure 2: An example of improving a comment. On the left is a function and a comment on it before improvement, which, firstly, has a typo, and secondly, contains a minimum of information about the code. The comment after the improvement is devoid of these shortcomings.

quire specifying data types, which makes automatic modification difficult. C# and Java occupy an intermediate position: C# shows a high average percentage of improvements (80.1%) with a peak in the full comments category (92.4%), while Java shows moderate results (avg = 48.2%).

#### F Training and Results

The models were trained for 5 epochs with a context length of 2000, a learning rate of 1e-4, and a cosine scheduler with a weight decay of 0.1 and a warmup ratio of 0.01. We used LORA (Hu et al., 2021) adapters with a rank of 8, alpha of 16, and a dropout rate of 0.05 for finetuning. From the synthetic part of the dataset, we sampled 1,500 examples for each programming language, resulting in 7,500 examples. For calculating metrics on real data, we sampled 100 examples for each programming language. The comparison is made with the base models to determine the extent to which training on our synthetic dataset improves the quality. Notably, with a batch size of 1, the model takes approximately 20 hours to train on 5 programming languages using DeepSpeed Zero2 (Rasley et al., 2020) on a single A100 GPU. The results are shown in Table 6.

544

545

546

547

548

549

550

551

552

553

554

539

540

541

542

543

Programming		Non structured	Incomplete	Complete	
language		Non-structured	mcomplete	Complete	
Deathar	#Enhanced comments	10 775	$3\ 455$	395	$\sum = 14\ 625$
i yünon	% out of the original quantity	24.2%	23.2%	48.1%	$\mathrm{avg}=31.9\%$
Iarra	#Enhanced comments	7066	3 810	5  407	$\sum = 16~283$
Java	% out of the original quantity	32.0%	57.6%	55.1%	$\mathrm{avg}=48.2\%$
C	#Enhanced comments	$3 \ 018$	-	$4\ 260$	$\sum = 7 \ 278$
GO	% out of the original quantity	77.1%	-	91.5%	$\mathrm{avg}=84.3\%$
$\mathrm{C}\#$	#Enhanced comments	$12 \ 467$	18  148	9 100	$\sum = 39~715$
	% % out of the original quantity	74.8%	73.1%	92.4%	$\mathrm{avg}=80.1\%$
JS	#Enhanced comments	1 386	164	97	$\sum = 1~647$
	% % out of the original quantity	20.4%	20.4%	59.5%	$\mathrm{avg}=33.5\%$

Table 5: Statistics on the improvement of Russian-language comments on programming languages, divided into categories: unstructured, incomplete and complete structured comments. For each language, the absolute number of improved comments, the percentage of improvements relative to the initial number in the category (from the Table 4), the total number of improvements ( $\Sigma$ ) and the average percentage of improvements (avg) are indicated. The dash in the category of incomplete comments for Go reflects their absence in the source data due to the simplified standard for documenting functions.

Madal	Python Java		Go		C#		JavaScript			
Woder	BERTScore	e chrf++	BERTScore	ore chrf++ BERTS		BERTScore chrf++		BERTScore chrf++		e chrf++
Baselines										
DeepSeek-Coder 1.3B	0.837	18.3	0.827	19.2	0.811	10.4	0.812	18.4	0.839	24.7
	$\pm 0.041$	$\pm 9.8$	$\pm 0.040$	$\pm 7.2$	$\pm 0.042$	$\pm 4.5$	$\pm 0.044$	$\pm 16.9$	$\pm 0.038$	$\pm 8.7$
DeepSeek-Coder 6.7B	0.878	34.1	0.873	36.9	0.838	21.0	0.844	36.3	0.876	38.4
	$\pm 0.043$	$\pm 10.5$	$\pm 0.044$	$\pm 14.2$	$\pm 0.047$	$\pm 11.1$	$\pm 0.052$	$\pm 18.2$	$\pm 0.033$	$\pm 10.9$
Qwen2.5-Coder 0.5B	0.863	26.6	0.839	20.7	0.816	10.9	0.815	14.1	0.799	9.6
	$\pm 0.052$	$\pm 9.8$	$\pm 0.056$	$\pm 9.3$	$\pm 0.052$	$\pm 5.6$	$\pm 0.052$	$\pm 8.5$	$\pm 0.035$	$\pm 6.1$
Qwen2.5-Coder 1.5B	0.841	22.8	0.838	21.2	0.815	11.5	0.821	31.5	0.841	23.8
	$\pm 0.045$	$\pm 10.8$	$\pm 0.045$	$\pm 10.5$	$\pm 0.039$	$\pm 5.0$	$\pm 0.051$	$\pm 14.9$	$\pm 0.035$	$\pm 7.9$
Qwen2.5-Coder 3B	0.784	14.2	0.829	17.2	0.819	11.0	0.817	25.7	0.841	23.7
	$\pm 0.061$	$\pm 8.4$	$\pm 0.039$	$\pm 6.0$	$\pm 0.041$	$\pm 4.4$	$\pm 0.046$	$\pm 15.5$	$\pm 0.033$	$\pm 6.2$
Qwen2.5-Coder 7B	0.880	34.3	0.873	35.0	0.854	23.5	0.847	24.3	0.872	33.5
	$\pm 0.040$	$\pm 7.7$	$\pm 0.039$	$\pm 9.8$	$\pm 0.039$	$\pm 9.1$	$\pm 0.037$	$\pm 12.2$	$\pm 0.031$	$\pm 7.9$
Finetuned Models										
Qwen2.5-Coder 0.5B	0.873	35.3	0.872	39.7	0.859	28.7	0.849	44.4	0.871	40.3
	$\pm 0.042$	$\pm 9.0$	$\pm 0.040$	$\pm 9.8$	$\pm 0.038$	$\pm 6.8$	$\pm 0.041$	$\pm 10.2$	$\pm 0.035$	$\pm 0.03$
Qwen2.5-Coder 1.5B	0.877	34.4	0.880	41.6	0.863	32.1	0.857	45.7	0.877	40.3
	$\pm 0.040$	$\pm 7.5$	$\pm 0.036$	$\pm 8.8$	$\pm 0.035$	$\pm 6.3$	$\pm 0.038$	$\pm 9.3$	$\pm 0.031$	$\pm 0.03$
Qwen2.5-Coder 3B	0.880	34.9	0.881	40.6	0.864	32.5	0.859	46.4	0.878	41.3
	$\pm 0.040$	$\pm 7.5$	$\pm 0.035$	$\pm 8.3$	$\pm 0.035$	$\pm 6.2$	$\pm 0.037$	$\pm 9.7$	$\pm 0.031$	$\pm 8.5$
Qwen2.5-Coder 7B	0.878	35.5	0.882	42.0	0.867	32.9	0.859	45.9	0.879	41.4
	$\pm 0.039$	$\pm 7.3$	$\pm 0.036$	$\pm 8.9$	$\pm 0.035$	$\pm 6.2$	$\pm 0.034$	$\pm 9.5$	$\pm 0.032$	$\pm 7.6$

Table 6: Comparison of base and finetuned models using BERTScore and chrf++ metrics with statistical significance testing (Mann-Whitney criterion). Statistically significant improvements (p < 0.05) are highlighted in **bold** when comparing the finetuned model with the corresponding sized base version. The values are presented as the average  $\pm$  standard deviation.