# CoLeM: A framework for semantic interpretation of Russian-language tables based on contrastive learning

Anonymous ACL submission

## Abstract

Tables are extensively utilized to represent and store data, however, they often lack explicit semantics necessary for machine interpretation of their contents. Semantic table interpretation is essential for integrating structured data with knowledge graphs, yet existing methods face challenges with Russian-language tables due to limited labeled data and linguistic peculiarities. This paper introduces a contrastive learning approach to minimize reliance on manual labeling and enhance the accuracy of column annotation for rare semantic types. The proposed method adapts contrastive learning for tabular data through augmentations and employs a distilled multilingual BERT model trained on the unlabeled RWT corpus (comprising 7.4 million columns). The resulting table representations are incorporated into the RuTaBERT pipeline, reducing computational overhead. Experimental results demonstrate a micro-F1 score of 97% and a macro-F1 score of 92%, surpassing several baseline approaches. These findings emphasize the efficiency of the proposed method in addressing data sparsity and handling unique features of the Russian language. The results further confirm that contrastive learning effectively captures semantic similarities among columns without explicit supervision, which is particularly vital for rare data types.

# 1 Introduction

005

007

011

017

018

019

028

Tabular data are one of the key formats for presenting structured information in various domains, ranging from scientific research to business analytics. It is widely used in relational databases, spreadsheets, web resources, and documents, making its processing critically important for automating data analysis. However, tables typically lack explicit semantics necessary for machine interpretation of their content. Therefore, the semantic interpretation of tables, especially in non-English languages, remains a challenging task (Gilbert Badaro, 2023; Jixiong Liu and Monnin, 2023). The primary challenges are associated with mapping individual table elements (columns, rows, cells) to concepts from knowledge graphs such as DBpedia or Wikidata, as well as handling the structural and linguistic diversity of data. 043

045

047

049

051

054

055

057

059

060

061

062

063

064

065

066

067

068

069

070

071

073

074

075

076

077

078

079

Russian-language tables pose a particular challenge due to the limited availability of specialized tools and annotated datasets. Most modern methods, particularly those based on pretrained language models like BERT (Xiang Deng and Yu, 2020; Jonathan Herzig and Eisenschlos, 2020; Pengcheng Yin and Riedel, 2020; Hiroshi Iida and Iyyer, 2021; Zhiruo Wang and Zhang, 2021; Yoshihiko Suhara and Tan, 2022), require vast amounts of labeled data, which are often unavailable or imbalanced for the Russian language. Moreover, existing solutions developed for English do not adapt well to other languages due to differences in tokenization and contextual semantics.

In this paper, we propose a novel approach for column type annotation in Russian-language tables based on contrastive learning. This approach effectively leverages unlabeled tabular data to train robust vector representations, reducing the reliance on manual annotation. Our contributions include:

- 1. Adaptation of contrastive learning for Russianlanguage tabular data using augmentations such as cell deletion and rearrangement.
- 2. Utilization of the distilled multilingual model DistilBERT, which balances performance and computational costs.
- 3. Integration of pre-trained tabular representations into an existing annotation pipeline based on the RuTaBERT framework, demonstrating the flexibility of the approach.
- 4. Experiments on the large Russian-language 080 dataset, RWT-RuTaBERT, showed that the 081

100

102

103

104

106

108

109

110

111

112

113

114

115

116

117

118

119

120

121

123

124

125

126

127

129

130

131

proposed approach outperforms certain baseline solutions, confirming its effectiveness under conditions of data sparsity and linguistic specificity.

The paper is organized as follows: Section 2 reviews the current state of research on semantic table interpretation. Section 3 describes the proposed approach for column type annotation in Russian-language tables, including data preparation, model architecture, and training algorithm. Section 4 presents experimental evaluations of the proposed approach's performance. Finally, Section 5 discusses the obtained results and outlines plans for future work.

# 2 Related works

Semantic table interpretation (STI) refers to the process of recognizing and linking tabular data to concepts from a target knowledge graph, ontology, or external vocabulary (e.g., DBpedia, Wikidata, Yago, Freebase, WordNet) (Jixiong Liu and Monnin, 2023; Zhang and Balog, 2020). One of the core tasks of STI is column type annotation, which involves mapping table columns to semantic types (classes and properties) from the target knowledge graph.

Over the past few years, existing methods and models have leveraged advances in deep machine learning, formulating the column type annotation task as a multi-class classification problem. For instance, (Madelon Hulsebos and Hidalgo, 2019) employed neural networks and various extracted feature groups, such as word and character embeddings, as well as global column statistics. The study by (Dan Zhang and Tan, 2020) incorporated analysis of local (intra-table) context (adjacent columns relative to the target column), while (Daheng Wang and Jiang, 2021) further added inter-table context to improve predictions. However, particular interest lies in works utilizing pre-trained language models based on the Transformer architecture. Transformer blocks employ an attention mechanism, enabling the model to generate useful contextualized embeddings for structural components of tabular data, such as cells, columns, or rows. Additionally, language models pre-trained on large-scale text corpora can encode semantics from the training text into model parameters, making fine-tuning on specific downstream tasks highly efficient. Examples of such works include models like TURL (Xiang Deng and Yu, 2020), TaPas (Jonathan Herzig



Figure 1: An example of data sparsity issue in the Viznet dataset.

and Eisenschlos, 2020), TaBERT (Pengcheng Yin and Riedel, 2020), TABBIE (Hiroshi Iida and Iyyer, 2021), TUTA (Zhiruo Wang and Zhang, 2021), and Doduo (Yoshihiko Suhara and Tan, 2022). 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

158

159

160

161

162

163

164

165

166

Existing solutions in this area achieve high performance due to the availability of large labeled training datasets. Specifically, English-language datasets may include hundreds of thousands of labeled columns (e.g., VizNet-Sato (Dan Zhang and Tan, 2020)  $\sim$  100,000, WikiTables-TURL (Xiang Deng and Yu, 2020)  $\sim$  600,000), while the Russian-language tabular dataset RWT-RuTaBERT contains over 1.4 million columns. Creating such datasets is a labor-intensive process requiring significant time and resources. Moreover, existing table datasets often suffer from data sparsity, manifested in a highly imbalanced distribution of semantic types (known as a "long-tail distribution"). For instance, some semantic types correspond to hundreds of thousands of columns, while others are associated with only a few dozen. As a result, models struggle to capture sufficient signals for minority (rare) semantic types (e.g., "athlete", "mountain range" or "insurance company"), even in supervised settings. Figure 1 illustrates this issue with a distribution chart of the 20 most frequent semantic types in the VizNet-Sato dataset. Figure 2 shows the same issue for the RWT-RuTaBERT dataset.

It should also be noted that current methods based on pre-trained language models are not universally applicable. There is a gap between the effectiveness of existing solutions on test cases and their practical applicability, particularly for tables in non-English languages and with varying structural layouts.



Figure 2: An example of data sparsity issue in the RWT-RuTaBERT dataset.

To enhance general table understanding and address various tabular tasks, recent works have employed large language models, which often outperform pre-trained models like BERT. These models are also more robust to unseen examples due to specific effects arising from their scale and training on vast text corpora. Examples include models such as Table-GPT (Peng Li and Chaudhuri, 2024), TableLlama (Tianshu Zhang and Sun, 2024), and approaches in (Korini and Bizer, 2024). However, a major drawback of such solutions is their requirement for substantial computational resources, hindering practical use.

To address the aforementioned challenges, we propose the use of self-supervised learning methods, specifically contrastive learning, to derive tabular representations from a large corpus of unlabeled tabular data. These representations can be used for determining relatedness between two tables (via cosine embedding similarity) and for fine-tuning with limited labeled data for specific downstream tasks.

**3** Proposed approach

167

168 169

170

171

172

173

174

175

176

177

178

180

181

182

184

188

189

190

191

192

194

195

196

198

#### 3.1 Problem statement

A table is a two-dimensional data structure composed of rows and columns. Table cells may contain textual data, numerical values, dates, times, etc.Tables can be categorized into three types based on the structure of information:

- 1. Highly structured (relational database tables);
- 2. Semi-structured (spreadsheets created in specialized software, e.g., MS Excel);

3. Unstructured (table images in PDF documents).

199

200

201

203

204

205

207

208

209

210

211

213

214

215

216

217

218

219

220

221

222

223

224

227

230

231

232

233

234

235

236

238

239

Tables can also be classified into three main groups based on orientation:

- 1. Vertical tables where data is arranged in vertical columns (i.e., top to bottom);
- 2. Horizontal tables where data is arranged in horizontal lines (i.e., left to right);
- 3. Matrix tables where each entry is indexed by row and column key(s).

This work focuses solely on vertical, highly structured, and semi-structured tables. The formal description of an input table can be represented as:

$$T = \{c_1, ..., c_n\}, c_i = \{v_1, ..., v_m\}, i \in \overline{1, n}$$
(1) 212

where T is a vertical table;  $c_i$  is an *i*-column;  $v_j$  is an *j*-cell of an *i*-column with  $j \in \overline{1, m}$ .

Our goal is to predict the column type, i.e., classify each column by its semantic type, such as "Book", "Writer", "Genre" or "Publication Date" rather than standard data types like string, integer, or datetime. The proposed approach involves using 170 distinct semantic types derived from selected classes and properties (value properties and object properties) from the general-purpose knowledge graph DBpedia<sup>1</sup>. Only Russian labels for these types (via language tags) were used, as the approach targets the annotation of Russian-language tables. Formally, this task can be described as:

$$P(c_i) \in KG_{st}, KG_{st} = \{st_1, \dots, st_{170}\}, \quad (2)$$

where  $P(c_i)$  is a predicted semantic type for a *i*-column;  $KG_{st}$  is a set of all semantic types with a cardinality of 170 in this case.

An example of solving the column annotation task for an input table is shown in Figure 3.

The core idea of the approach is to develop an encoder for robust tabular representations based on contrastive learning, which can then be applied to downstream tasks, specifically semantic annotation of columns in Russian-language tables. The general schema of the proposed approach is presented in Figure 4.

<sup>&</sup>lt;sup>1</sup>https://www.dbpedia.org/





Figure 3: An example of the CTA task.

#### 3.2 Dataset Description

240

241

242

244

245

246

247

249

251

257

258

259

260

261

262

264

The pre-trained table encoder is trained on a vast amount of tabular data that does not require manual annotation. The large-scale Russian Web Tables (RWT) corpus (Platon E. Fedorov and Chernishev, 2023) is used as the source dataset. This dataset represents a snapshot of tables from the Russian Wikipedia as of September 13, 2021. Key statistics for the RWT corpus are provided in Table 1.

Statistics	Value
Number of tables	1 266 731
Number of columns	7 419 771
Number of cells	99 638 194
Average number of cells per table	81.78
Set size	17 GB
Percentage of almost empty columns	6%
Average number of cells per column	13.42
Percentage of numeric columns	17%

Table 1: Statistics of the RWT table corpus.

During the initial data preprocessing stage, vertical tables were selected from the original RWT corpus. Each column from such a table is represented as a data string using the cell delimiter "«". Subsequent data cleaning was performed using

the following operations:

- Filtering out empty columns.
- Removing parser metadata wrapping text using regular expressions.
- Removing links to Wikipedia articles.
- Removing special characters (e.g., "@", "&", "?", and "!").
- Removing empty cells within columns.
- Removing columns with fewer than three cells, as such columns become unrepresentative after cell deletion augmentations.

As a result of these cleaning operations, an unlabeled dataset of Russian-language tabular data consisting of 4,656,668 columns was obtained. This preprocessing was automated using a specialized tool, LoReTA.

## 3.3 Training Algorithm

Contrastive learning is a self-supervised learning technique designed to obtain informative embeddings. It involves maximizing a consistency metric, in our case cosine similarity, between positive pairs (data instances) while minimizing this metric between negative pairs. Contrastive learning enables effective training on unlabeled data corpora.

In this work, we adapt the contrastive learning concept proposed in (Ting Chen and Hinton, 2020) for tabular data. The contrastive learning algorithm for tabular data is illustrated in Figure 5.

The main idea is to construct two augmentations for each column in a batch during training. Column embeddings are generated for the resulting augmentations using an encoder model. Representations of augmentations derived from the same column are considered a positive pair, and our goal is to maximize the cosine similarity metric for this pair. Conversely, representations of augmentations derived from different columns are considered negative pairs, for which the task is to minimize the cosine similarity metric.

#### 3.3.1 Data Augmentation

Data augmentation refers to a technique for artificially increasing the size of a training dataset by applying transformations to the original data. This technique is widely used in scenarios with limited or no labeled data to enhance the model's generalization ability. In contrastive learning, augmentations play a critical role in forming semantically consistent positive pairs.

Common augmentations for tabular data include:

- Random cell deletion.
  Deletion/rearrangement/replacement of tokens in a cell.
  Row sampling (e.g., 50% of rows).
  Cell rearrangement within a table row.
  Column deletion.
- Column rearrangement within a table. 309



Figure 4: General scheme of the proposed approach.

Currently, there is no research identifying the most effective augmentations for forming semantically consistent pairs in the context of tabular data processing. Therefore, in this work, we selected two augmentations deemed most promising: random cell deletion and cell rearrangement within a column. For random cell deletion, 10% of all cells in a column are removed.

## 3.3.2 Contrastive Loss

318

319

321

323

325

327

329

330

331

332

333

334

335

338

Contrastive loss functions are widely used in representation learning tasks, as they enable models to better distinguish internal data structures and, consequently, extract more useful representations. A contrastive loss function aims to maximize agreement between positive pairs and minimize agreement between negative pairs in the vector space.

There are several variations of contrastive loss functions. In this work, we adopt the NT-Xent loss (Normalized Temperature Cross-Entropy Loss) used in (Ting Chen and Hinton, 2020), defined as:

$$L = \frac{1}{2N} * \sum_{k=1}^{N} [l(2k-1,2k), (2k,2k-1)],$$
  
$$l(i,j) = -\log \frac{exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} 1_{k \neq i} * exp(s_{i,k}/\tau)}, \quad (3)$$
  
$$s_{i,j} = \frac{z_i * z_j}{||z_i|| * ||z_j||}$$

where  $1_{[k \neq i]}$  is 1 if  $k \neq i$ , otherwise 0;  $\tau$  is the temperature parameter; and *s* is cosine similarity.

### 3.4 Model Architecture

Currently, Transformer-based models are central to natural language processing tasks. These models are versatile tools for text processing due to their ability to capture contextual dependencies between words in sequences and to train on unlabeled or partially labeled data. They achieve this efficiently through high parallelism, making them preferable for training on large datasets. 339

340

341

343

344

345

347

348

350

351

352

353

354

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

According to (Ting Chen and Hinton, 2020), two critical hyperparameters in contrastive learning are batch size and the number of epochs. Larger batch sizes and more epochs result in more representative embeddings, leading to better performance on downstream tasks during fine-tuning.

Based on this, the distilled multilingual BERT model<sup>2</sup> was chosen as the base encoder. This model was trained on Wikipedia articles in 104 different languages. Unlike the base version<sup>3</sup>, it consists of only 6 layers (half the number of the base version) and 12 attention heads. It has 134 million parameters (compared to 177 million in the base version).

Model distillation is a technique in machine learning where knowledge is transferred from a more complex model (teacher) to a more compact one (student) while maintaining prediction quality.

This technique, combined with reducing the tokenizer's maximum sequence length to 256 tokens, enabled training with a batch size of 800, which is 25 times larger than that of a comparable state-ofthe-art English-language solution (Miao and Wang, 2023).

Research in (Ting Chen and Hinton, 2020) explored the use of projecting the encoder's output layer into a latent space for calculating the contrastive loss. Results indicate that applying a non-linear projection during training positively impacts representation quality. Thus, in this work, a two-layer perceptron (MLP) is used after the encoder's output layer to project into a 128-dimensional la-

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/distilbert/

distilbert-base-multilingual-cased

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/google-bert/ bert-base-multilingual-cased



Figure 5: Contrastive learning algorithm for tabular data.

tent space where the contrastive loss is computedusing the aforementioned formula.

# 4 Experimental Evaluation and Discussion

376

377

381

384

390

391

396

400

401

402

403

All experiments were conducted on a graphics cluster. The cluster configuration includes two 16-core Intel Xeon Gold 6326 "Ice Lake" 2.9 GHz processors, four NVIDIA A100 80 GB PCIe GPUs, and 2 TB of DDR4-3200 RAM.

#### 4.1 Contrastive Learning Setup

The approach was implemented in Python using the PyTorch and Transformers libraries. The AdamW optimizer (lr = 5e-5, eps = 1e-6) was chosen for gradient descent. To accelerate convergence, cosine annealing was applied to dynamically reduce the learning rate. The temperature parameter, a hyperparameter of the contrastive loss function, was set to 0.1, as this value was found to be optimal in (Ting Chen and Hinton, 2020). Under these settings, the pre-trained encoder model was trained for 100 epochs on 4 NVIDIA A100 GPUs using the Distributed-Data-Parallel technology of the PyTorch framework. Training lasted 9 days, 9 hours, and 53 minutes. GPU memory consumption amounted to 290 GB.

## 4.2 Setup for Semantic Column Annotation Model

In this work, semantic interpretation (annotation) of table columns was selected as the downstream task. Previously, the RuTaBERT framework was proposed for this task, based on fine-tuning a pre-404 trained multilingual BERT model using the spe-405 cially prepared RWT-RuTaBERT dataset. This 406 dataset contains approximately 1.56 million labeled 407 columns. The core idea is to utilize the existing 408 pipeline of this framework, replacing the standard 409 BERT model with a specialized pre-trained table 410 encoder. The RWT-RuTaBERT dataset, with all 411 standard settings, was used for training. The valida-412 tion set comprised 5% of the total training subset. 413 The technique of neighboring column serialization 414 was used to decompose column values into token 415 sequences. 416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

According to (Ting Chen and Hinton, 2020), the projection layer is trained to be invariant to data transformations, potentially losing information useful for downstream tasks. Therefore, for further fine-tuning of the table encoder, the output from the first linear layer of the projection with a LeakyReLU activation function was used. Standard training settings defined in the RuTaBERT framework were applied. The model was fine-tuned for 30 epochs with a batch size of 32 on the RWT-RuTaBERT dataset using 2 NVIDIA A100 GPUs. Training lasted 2 days, 20 hours, and 15 minutes, with GPU memory consumption of 9.9 GB. Additionally, a model with a batch size of 256 was trained with all other hyperparameters unchanged. Under these settings, training took 4 days, 3 hours, and 1 minute, with GPU memory consumption of 52 GB.

4.3

dataset.

follows:

**Evaluation Metrics** 

The primary metrics for evaluating the performance

of the proposed approach are averaged F1 scores, as

the task involves multiclass classification. Specif-

ically, micro F1, macro F1, and weighted F1 are

used due to the imbalance in the RWT-RuTaBERT

 $microF1 = 2 \frac{MicroPrecision * MicroRecall}{MicroPrecision + MicroRecall}$ 

Macro F1 is the average F1 score for each seman-

tic type (class), treating all classes equally without

accounting for class imbalance. It is calculated as

 $macroF1 = \frac{1}{N} * \sum_{i=1}^{N} F1_i$ 

where N is the number of semantic types (classes),

The weighted F-measure is calculated for each

class and then aggregated as a weighted average,

taking into account the number of instances for

each class. Unlike the micro F1, this metric consid-

ers the class imbalance. The weighted F-measure

 $weighted F1 = \sum_{i=1}^{C} [w_i * F1_i], w_i = \frac{n_i}{N}$ 

where C is the number of classes,  $n_i$  is the num-

ber of samples in the *i*-th class, N is the number of

samples and  $F1_i$  is the F1 score for the *i*-th class.

The results of the experimental evaluation are pre-

sented in Table 2. A comparison of the perfor-

mance of the proposed approach with several base-

which specializes in processing the Russian lan-

guage, was selected. One of the transfer learning

techniques was applied, where the weights of the

Firstly, a pre-trained language model, RuBERT<sup>4</sup>,

is computed using the following formula:

 $F1_i$  is the F1 for the i-th class.

4.4 Results and Discussion

line solutions is provided.

(5)

(6)

sion matrix and is defined as follows:

Micro F1 is calculated across the entire confu-

- 437 438
- 439
- 440 441
- 442 443
- 444
- 445 446
- 447 448

- 449
- 450
- 451
- 452 453

454

455 456

- 457
- 458
- 459 460

461

462

463

- 464

465

467

468

470 471

472

473

474

469

encoder layers remained unchanged during train-

ing. Thus, during fine-tuning of RuBERT on the RWT-RuTaBERT dataset, only the parameters of the classification layer were adjusted.

> <sup>4</sup>https://huggingface.co/DeepPavlov/ rubert-base-cased

Model	micro	macro	weighted
	F1	F1	F1
Doduo	0.140	0.040	_
RuBERT-ft	0.610	0.410	0.590
Doduo-ft	0.962	0.890	0.960
RuTaBERT	0.964	0.900	0.963
CoLeM-bs32	0.969	0.910	0.969
CoLeM-bs256	0.974	0.924	0.974

Table 2: Results of experimental evaluation on the RWT-RuTaBERT dataset.

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

Secondly, the state-of-the-art framework Doduo (Yoshihiko Suhara and Tan, 2022) was chosen, which is a leading solution for the task of semantic annotation of columns and relationships between them. In this case, transfer learning was also applied by freezing the transformer layers and finetuning only the final linear classifier layer. Additionally, a full fine-tuning of the multilingual BERT model was performed following the Doduo approach on the RWT-RuTaBERT dataset (Finetuned Doduo).

Thirdly, the original RuTaBERT approach was considered.

The obtained evaluation results demonstrated that the proposed approach outperformed all baseline solutions in both training configurations (batch sizes of 32 and 256). Specifically, the experiment showed that while the RuBERT model is tailored for processing the Russian language, it is not directly suited for tabular tasks, which proved challenging for this model. Consequently, existing Russian-language models cannot be effectively applied to the task of semantic column annotation.

The Doduo model, trained using transfer learning techniques, exhibited relatively low evaluation results. This is attributed to the fact that the model was trained on tabular data exclusively in English. Notably, the tokenizer of this model lacks sufficient Russian-language tokens. As a result, it can be concluded that a model trained on English data cannot be directly applied to another language, such as Russian, without modifying the base encoder to accommodate the target language.

Meanwhile, the fine-tuned multilingual encoder of the Doduo framework and the RuTaBERT approach demonstrated nearly comparable results in terms of evaluation metrics. However, it can be observed that the use of a pre-trained tabular encoder based on contrastive learning positively impacts

590

591

592

593

594

595

596

597

598

599

600

601

555

556

557

558

559

560

561

the performance. With a smaller model and iden-514 tical settings, the proposed approach achieved re-515 sults equivalent to those of the classical RuTaBERT 516 model or the fine-tuned Doduo. Additionally, the 517 model consumes approximately three times less GPU memory during training, requiring less than 519 10 GB (with a batch size of 32, consistent across all 520 three models), which enables training on a standard 521 home computer. Furthermore, with a larger batch size (e.g., 256), the proposed approach achieved 523 a performance gain of 1.5% compared to the clas-524 sical RuTaBERT model and nearly 3% compared 525 to the fine-tuned Doduo. The experimental results 526 highlight the potential of our approach for semantic 527 annotation of Russian-language tables.

531

536

539

540

541

542

543

544

546

547

548

551

552

553

554

To further evaluate CoLeM's performance, we conducted a statistical analysis on three aspects:

1) Datatype groups: We categorized all columns from the collected tables into 5 basic groups: Datetime, Numeric, Links, Short Text, and Long Text. Datetime columns included dates, years, or times. Numeric columns contained only numbers, for example, the results of measurements of length, weight or age. URL columns included different web addresses. Text columns were further divided into Short Text (tokens fewer than four) and Long Text (tokens four or more). We also identified a separate Persons datatype, given the prevalence of instances like "*employer*", "*screenwriter*", "*athlete*", and "*football player*". Table 3 summarizes the Micro F1 score and distribution for each datatype group.

Data type	F1 (CoLeM)	F1 (RuTaBERT)
Datetime	0.948	0.941
Long text	0.858	0.885
Numeric	0.760	0.749
Person	0.716	0.692
Short text	0.932	0.926
Links	0.611	0.699

Table 3: The performance for the six datatype groups.

2) Rare semantic types: Performance evaluations were also conducted for the 15 least frequently occurring semantic types. For comparison, checkpoints of the CoLeM-bs32 and RuTaBERT models, which achieved the highest macro F1 score on the training set, were used. The results are presented in Table 4.

The results demonstrate that, due to the robust tabular representations obtained, the CoLeM model

significantly outperforms the existing state-of-theart (SOTA) Russian-language solution, RuTaBERT, in terms of evaluation metrics for infrequently occurring semantic types.

**3) Model convergence:** To evaluate the convergence of the CoLeM model, experiments were conducted for checkpoints of CoLeM-bs32 and RuTaBERT models trained for 10 epochs. The performance results are summarized in Table 5.

It can be observed that the CoLeM model converges faster than the RuTaBERT model and has 1-3% better performance. This allows us to use a smaller number of epochs in training stage, while obtaining comparable or even superior performance to the RuTaBERT model.

## 5 Conclusion

This study proposes an approach for semantic annotation of columns in Russian-language tables based on contrastive learning. The experimental results demonstrate that the approach mitigates the dependency on large volumes of labeled data by leveraging self-supervised learning on unlabeled tables. Moreover, it outperforms existing baseline solutions (Doduo and RuTaBERT) in terms of evaluation metrics, particularly for rare semantic types. The approach also ensures computational efficiency through the use of a distilled model and optimized batch sizes, reducing memory requirements by 60% compared to analogous methods.

The results of the experimental evaluation confirm the effectiveness of the proposed solution. In the future, as part of a research project with the Ivannikov Institute for System Programming of the Russian Academy of Sciences (ISP RAS), it is planned to integrate these results into a specialized table processor within the Talisman platform<sup>5</sup>. Additionally, the approach will be extended to tables with horizontal and matrix layouts. Further investigation will also focus on the use of new data augmentations to enhance the robustness of tabular representations.

Overall, the proposed approach opens up opportunities for the development of universal systems for semantic interpretation of tables, which is relevant for tasks involving the integration of structured and semi-structured information, as well as business analytics.

<sup>&</sup>lt;sup>5</sup>http://talisman.ispras.ru

### References

602

610

611

612

613

614

615

618

623

635

641

643

646

647

650

651

- Colin Lockard Binxuan Huang Xin Luna Dong Daheng Wang, Prashant Shiralkar and Meng Jiang. 2021.
   Tcn: Table convolutional network for web table interpretation. In *Proceedings of the Web Conference (WWW'21)*, pages 4020–4032.
- Yoshihiko Suhara Çağatay Demiralp Jinfeng Li Dan Zhang, Madelon Hulsebos and Wang-Chiew Tan. 2020. Sato: Contextual semantic type detection in tables. *Proceedings of the VLDB Endowment*, 13(11):1835–1848.
- Paolo Papotti Gilbert Badaro, Mohammed Saeed. 2023. Transformers for tabular data representation: A survey of models and applications. *Transactions of the Association for Computational Linguistics*, 11:227–249.
  - Varun Manjunatha Hiroshi Iida, Dung Thai and Mohit Iyyer. 2021. Tabbie: Pretrained representations of tabular data. In *Proceedings of the 2021 Conference* of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 3446–3456.
  - Raphaël Troncy Viet-Phi Huynh Thomas Labbé Jixiong Liu, Yoan Chabot and Pierre Monnin. 2023. From tabular data to knowledge graphs: A survey of semantic table interpretation tasks and methods. *Journal of Web Semantics*, 76:100761.
  - Thomas Müller-Francesco Piccinno Jonathan Herzig, Pawel Krzysztof Nowak and Julian Eisenschlos. 2020. Tapas: Weakly supervised table parsing via pre-training. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL'2020), pages 4320–4333.
  - Keti Korini and Christian Bizer. 2024. Column property annotation using large language models. In *Proceedings of the Semantic Web: ESWC 2024 Satellite Events*, pages 61–70.
  - Michiel Bakker Emanuel Zgraggen Arvind Satyanarayan Tim Kraska Çagatay Demiralp Madelon Hulsebos, Kevin Hu and César Hidalgo. 2019. Sherlock: A deep learning approach to semantic data type detection. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery Data Mining (KDD'19)*, pages 1500– 1508.
- Zhengjie Miao and Jin Wang. 2023. Watchog: A lightweight contrastive learning based framework for column annotation. *Proceedings of the ACM on Management of Data*, 1(3):1–24.
- Dror Yashar Weiwei Cui Song Ge Haidong Zhang Danielle R. Fainman Dongmei Zhang Peng Li, Yeye He and Surajit Chaudhuri. 2024. Table-gpt: Table fine-tuned gpt for diverse table tasks. *Proceedings of the ACM on Management of Data*, 2(3):1–28.

Wen-tau Yih Pengcheng Yin, Graham Neubig and Sebastian Riedel. 2020. Tabert: Pretraining for joint understanding of textual and tabular data. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL'2020), pages 8413–8426. 656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

- Alexey V. Mironov Platon E. Fedorov and George A. Chernishev. 2023. Russian web tables: A public corpus of web tables for russian language based on wikipedia. *Lobachevskii Journal of Mathematics*, 44:111–122.
- Yifei Li Tianshu Zhang, Xiang Yue and Huan Sun. 2024. Tablellama: Towards open large generalist models for tables. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 6024–6044.
- Mohammad Norouzi Ting Chen, Simon Kornblith and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*, pages 1597–1607.
- Alyssa Lees You Wu Xiang Deng, Huan Sun and Cong Yu. 2020. Turl: Table understanding through representation learning. *Proceedings of the VLDB Endowment*, 14(3):307–319.
- Yuliang Li Dan Zhang Çağatay Demiralp Chen Chen Yoshihiko Suhara, Jinfeng Li and Wang-Chiew Tan. 2022. Annotating columns with pre-trained language models. In *Proceedings of the 2022 International Conference on Management of Data (SIGMOD'22)*, pages 1493–1503.
- Shuo Zhang and Krisztian Balog. 2020. Web table extraction, retrieval, and augmentation: A survey. *ACM Transactions on Intelligent Systems and Technology*, 11(2):1–35.
- Ran Jia Jia Li Zhiyi Fu Shi Han Zhiruo Wang, Haoyu Dong and Dongmei Zhang. 2021. Tuta: Treebased transformers for generally structured table pretraining. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery Data Mining* (*KDD*'21), pages 1780–1790.
- A Appendix: Evaluation for 15 least frequently occurring semantic types

Semantic type	Number of samples (test subset)	F1 (RuTaBERT)	F1 (CoLeM-bs32)
camera	102 (4)	0.250	0.750
employer	101 (10)	0.899	1.000
device	101 (8)	0.625	0.875
animal	93 (7)	0.857	1.000
magazine	93 (9)	0.440	0.440
continent	92 (8)	0.625	0.750
novel	89 (11)	0.818	0.909
law	89 (9)	1.000	1.000
wrestler	88 (5)	0.400	0.600
college	87 (5)	0.000	0.200
museum	86 (4)	0.500	0.750
firm	85 (6)	0.333	0.333
prefecture	83 (10)	0.600	0.699
road	83 (6)	0.500	0.666
quote	76 (7)	0.857	1.000

Table 4: Experimental evaluation results for the 15 least frequently occurring semantic types

# B Appendix: Model evaluation after 10 training epochs

Model	micro F1	Macro F1	Weighted F1
RuTaBERT (10 epochs)	0.952	0.856	0.952
CoLeM-bs32 (10 epochs)	0.966	0.888	0.966

Table 5: Results of model evaluation after 10 training epochs