
RoTaR: Efficient Row-Based Table Representation Learning via Teacher-Student Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose ROTAR, a row-based table representation learning method, to ad-
2 dress the efficiency and scalability issues faced by existing table representation
3 learning methods. The key idea of ROTAR is to generate query-agnostic row
4 representations that could be re-used via query-specific aggregation. In addition to
5 the row-based architecture, we introduce several techniques: cell-aware position
6 embedding, teacher-student training paradigm, and selective backward to improve
7 the performance of ROTAR model.

8 1 Introduction

9 Tabular data is one of the most widely used media for storing information. Table representation
10 learning has a wide range of downstream applications such as table question answering, table search,
11 table type detection, etc. It is thus vital to design an effective and practical solution for table
12 representation learning. However, despite its popularity and importance in modern data science, table
13 representation learning is not well addressed, compared to images, texts, or other media.

14 Most of the previous attempts [20, 11, 8, 7, 14, 19, 1] apply recent progress in natural language
15 processing (NLP), i.e., transformers and large language models (LMs). These works directly serialize
16 the entire table together with a query or related utterance into a sequence as the input to an LM, which
17 is pretrained on a sufficient amount of table corpus. However, as the most common and best practices
18 for table representation learning, this approach suffers from scalability and efficiency issues.

19 First, serializing a large table containing a large number of rows will result in a long sequence which
20 is hard to process by classical transformer-based models, because the complexity of such models
21 is quadratic to the length of the input sequence. To solve this problem, some works optimize the
22 transformer structure [15], while others use table specific solutions to reduce the complexity of
23 attention computation, such as restricting attention computation to the same row or column [6, 9] or
24 only between the schema and values [4]. However, these approaches do not eliminate the scalability
25 issue, because a pretrained LM is subject to a max sequence length constraint. For example, GPT-3
26 limits the input length to 2048 tokens, while BERT sets this limit as 512. A table with a small number
27 of rows can easily exceeds it, causing inevitable truncation and thus loss of information.

28 Second, the serialization process takes the given query as input, leading to query-specific encoding.
29 Because in real-world scenarios many queries concentrate on a few tables, repeatedly computing the
30 table representation for every new incoming query is inefficient.

31 To address the above problems, we proposed ROTAR which learns a *query-agnostic row-based* table
32 representation. Rather than serialize the whole table, ROTAR takes each row as input and efficiently
33 produces row level encodings which can be re-used by any queries. It then uses query-specific
34 aggregation to produce the table representation on top of these row encodings.

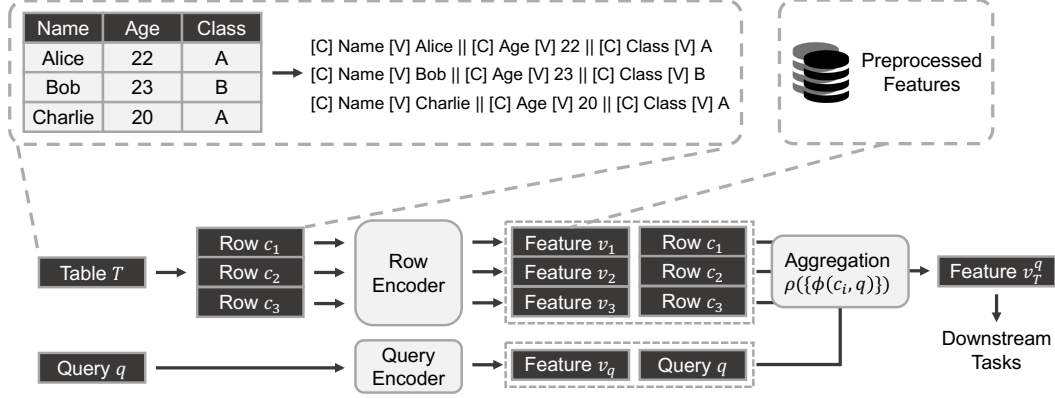


Figure 1: ROTAR model architecture.

35 2 ROTAR Methodology

36 2.1 Overall Architecture

37 **Row Independence Observation.** Interdependencies in the table structure are the key to reduce
 38 the computational complexity of transformer-based models: irrelevant attention can be saved in
 39 transformer-based models without harming their performance. We observe that although the informa-
 40 tion in different rows should be aggregated to form the representation of the whole table, there is no
 41 strong correlation among different rows.

42 Intuitively, a relational table can be viewed as a set of rows with homogeneous schema, because the
 43 order of the rows usually does not matter. The representation of a set is mathematically equivalent to
 44 an appropriate aggregation of the representations of each single item in the set [16, 21]. Therefore,
 45 table representation can be factorized into an aggregation of independent row representations.

46 Inspired by this observation, ROTAR uses a weight-shared row-based transformer model to encode
 47 each row in the table independently and ignores the inter-row correlation in this encoding (Fig. 1).

48 ROTAR consists of two components: query-agnostic row encoder M and query-specific aggregation.
 49 After training on a dataset, the learned row representations by the row encoder can be preprocessed
 50 and stored. Therefore, answering an upcoming query only requires computing the aggregation. It
 51 does not have to repeatedly run the row encoder, thus saving a massive amount of time.

52 **Row Encoder.** More specifically, ROTAR considers every cell $T_{i,j}$ in a table T as a textual cell, i.e.,
 53 $T_{i,j} = T_{i,j;1}T_{i,j;2} \cdots T_{i,j;n}$ where each $T_{i,j;k}$ is a token. Similarly, each attribute in the schema A
 54 is also viewed as textual, i.e., $A_j = A_{j;1}A_{j;2} \cdots A_{j;m}$, where each $A_{j;k}$ is a token. ROTAR thus
 55 serializes each cell $c_{i,j}$ by combining the attribute and the cell value $c_{i,j} = [\text{COL}]A_j[\text{VAL}]T_{i,j}$.
 56 Given a table with N rows and L columns, a shared encoder M encodes each concatenated row
 57 $c_i = c_{i,1} || c_{i,2} || \cdots || c_{i,L}$ into a fixed-dimension vector $v_i = M(c_i)$. Notice that the obtained set
 58 of row representations $\{v_1, v_2, \cdots, v_N\}$ is query-agnostic.

59 **Aggregation.** Then for each incoming query q , the resulted table representation can be computed by
 60 $v_T^q = \rho(\{\phi(c_i, q)\}_{i=1}^N)$, where ϕ is a learnable function that given a query q , extracts information
 61 from c_i . ρ is an appropriate aggregation function [16, 21]. The function ϕ and ρ together constitute
 62 a query-specific aggregation module. For instance, if a query q is encoded as a vector v_q of the
 63 same dimension with the row representations, setting $\phi(c_i, q) = v_i \odot v_q$ (\odot stands for point-wise
 64 multiplication) and $\rho(X) = \frac{1}{|X|} \sum_{x \in X} x$ yields the table representation as the average row vector
 65 projected onto the query vector v_q .

66 2.2 Query-agnostic Row Encoder

67 The ROTAR model tackles the first fore-mentioned issue in scalability, i.e., encoding a table with
 68 a large number of rows. Because the transformer model is run separately for each row and the
 69 aggregation module does not have to use transformer-based models, ROTAR is capable of handling

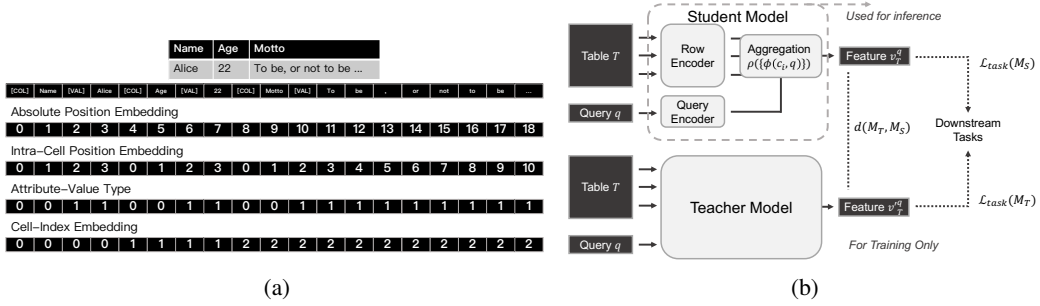


Figure 2: (a). Cell-aware position embedding. (b). Teacher-student paradigm.

70 any tables with any number of rows. Note the number of columns is not of concern, because the
 71 number of columns is usually much smaller than the number of rows.

72 The design of the query-agnostic row encoder M can be very flexible. For example, we can
 73 directly use a general-purpose pretrained LM like BERT [5] or RoBERTa [12], or pre-existing table
 74 representation models like TAPAS [8, 7] if we view each row as a table consisting of a single row. In
 75 addition to directly adopt the existing methods, we also propose two new techniques customized to
 76 row-based table representation learning.

77 When using a general-purpose pretrained LM as the row encoder M , instead of directly feeding the
 78 serialized row $c_i = c_{i,1} || c_{i,2} || \dots || c_{i,L}$ into the LM, we could take advantage of the structure of
 79 row to elaborate the position embedding design. Specifically, the position embedding $p_{i,j;k}^T$ of each
 80 token $T_{i,j;k}$ or $p_{j;k}^A$ of each token $A_{j;k}$ can be decomposed into inter-cell position embeddings (based
 81 on j or A_j) and intra-cell position embeddings (based on k) [3], which can then be customized for
 82 different purpose of use (Fig. 2a). For example, in many circumstances the order of attributes should
 83 be irrelevant to the representation of the row. By simply removing the absolute position embedding
 84 and the cell index embedding, the order of attributes is not perceived by the LM and thus the learned
 85 representation is robust against swapping order of columns.

86 2.3 Query-specific Aggregation

87 The ROTAR model also tackles the fore-mentioned issue in efficiency, i.e., learning representation
 88 re-usable for different queries. While the learned row-based representation is query-agnostic, a
 89 query-specific aggregation module is introduced to produce query-specific table representation.

90 The design of query-specific adaption function ϕ can be very straightforward, for example, $\phi(v_i, q) =$
 91 $v_i \odot v_q$, or $\phi(v_i, q) = \text{MLP}(v_i \oplus v_q)$ (\oplus stands for concatenation) or even $\phi(v_i, q) = \text{MLP}(v_i \oplus$
 92 $v_q \oplus |v_i - v_q| \oplus (v_i \odot v_q))$. Furthermore, notice that ϕ does not have to be differentiable or even
 93 numeric, traditional selective ϕ based on the textual input c_i could also be used. For example, the
 94 n-gram similarity weighted embedding $\phi(c_i, q) = \text{ns}(c_i, q) \cdot v_i = \frac{2 \cdot |\text{n-gram}(c_i) \cap \text{n-gram}(q)|}{|\text{n-gram}(c_i)| + |\text{n-gram}(q)|} \cdot v_i$, or a
 95 hard threshold $\phi_\alpha(c_i, q) = [\text{ns}(c_i, q) > \alpha] \cdot v_i$.

96 The choice of aggregation function ρ can be rather arbitrary, for example, Mean, Min, Max,
 97 LogSumExp, etc., are all feasible aggregation functions. However, the choice of aggregation naturally
 98 influences the table representation quality when handling different queries. Therefore, a learnable
 99 aggregation function ρ could make the model more flexible. For instance, we could learn a multi-
 100 head projection aggregation function, which has a set of learnable parameters $\Theta = \{\theta_l\}_{l=1}^d$, and
 101 $\rho_\Theta(X) = \frac{1}{|X|} \sum_{x \in X} \text{Nonlinear}(x \odot \theta_l)$.

102 2.4 Training Techniques

103 2.4.1 Teacher-Student Paradigm

104 ROTAR simplifies the table representation process by ignoring inter-row interactions. Therefore,
 105 it pursues efficiency and scalability at the expense of inevitable but acceptable performance decay.
 106 However, with the help of the previous table representation methods during training time, the ROTAR
 107 model is able to improve its performance while still being efficient during inference time.

Table 1: Experiment result.

Method	Test Acc. (%)	Inference Speed
TAPAS _{BASE} [†]	78.5% ± 0.3%	–
TAPAS _{BASE} (TabFact only)	69.9% ± 3.8%	–
Table-BERT	65.1%	–
TAPAS _{BASE} (no query)	50.3%	55ms/table
ROTAR	63.6%	15ms/table

108 Consider the teacher-student paradigm which is widely used in model distillation [10, 13, 17].
 109 Training the student model to mimic the features generated by the highly performant teacher model
 110 can provide more differentiable information and boost the student model’s learning process. The
 111 small and efficient student model is then used as a cost effective alternative to the original teacher
 112 model.

113 Technically, instead of only considering the loss function $\mathcal{L} = \mathcal{L}_{task}(M_S)$ of the downstream
 114 task during training, where M_S is the student model, the loss in teacher-student paradigm uses
 115 $\mathcal{L} = \alpha \cdot \mathcal{L}_{task}(M_T) + \beta \cdot \mathcal{L}_{task}(M_S) + \gamma \cdot d(M_T, M_S)$, where M_T is the teacher model, d is the
 116 mean squared error distance between features or logits generated by the teacher and student model,
 117 and α, β, γ are tunable hyperparameters (Fig. 2b).

118 2.4.2 Selective Backward

119 In practice, back-propagation through the aggregation of multiple rows could be unacceptable because
 120 of the GPU memory limit. However, since the row encoder M is shared, ROTAR is able to only
 121 sample some rows to back-propagate. The sampling process could be done either randomly or
 122 weighted according to a traditional selective ϕ , like n-gram similarity.

123 3 Experiments

124 We conduct preliminary experiments with the table fact verification task on the TabFact [2] dataset.
 125 In the table verification task, a query q is a statement to be evaluated based on a table T . The TabFact
 126 dataset consists of 16K tables obtained from Wikipedia and 118K labeled statements. A common
 127 public data split is provided along with the dataset. The results are shown in Tab. 1.

128 For fair comparison, we use the same model size to bert-base and google/tapas-base. Further,
 129 because the TAPAS_{BASE} best result[†] (78.5% ± 0.3%) is pretrained on 6.2M table-text examples
 130 obtained from Wikipedia, we compare against the results reported by TAPAS using only TabFact
 131 data (69.9% ± 3.8%) and Table-BERT (65.1%).¹

132 Note since the ROTAR model prioritizes efficiency and scalability over performance, it slightly
 133 sacrifices performance in exchange for big speed up. The performance sacrifice mainly comes from
 134 the query-agnostic property instead of the row-independency. In particular, with an accuracy drop of
 135 $\sim 6\%$, the ROTAR model with preprocessed feature vector is $\sim 3.7x$ faster than the TAPAS model,
 136 depending on the speed of the query encoder.

137 Further, we also compare against a TAPAS model that same to ROTAR, is not aware of the queries
 138 beforehand. We then finetuned it under the same setting to ROTAR. This TAPAS model achieves an
 139 accuracy of 50.3%, which is much lower than the accuracy of our ROTAR (63.6%). This confirms
 140 that ROTAR is indeed able to produce query-agnostic representation.

141 4 Conclusion

142 We propose ROTAR which uses a shared row-encoder to generate query-agnostic row representations
 143 and learns instance optimized aggregation function to produce query-specific table representation.
 144 Preliminary experiments on TabFact confirm that ROTAR significantly improves the scalability and

¹The TAPAS_{BASE} and TAPAS_{BASE} (TabFact only) result is reported by [7], while the Table-BERT result is reported by [2].

145 efficiency of table representation learning, with limited performance drop. In the further, we will
146 continue to optimize the ROTAR model to improve the speed-accuracy trade-off.

147 References

- 148 [1] Sercan Ö. Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In *AAAI*
149 *2021, IAAI 2021, EAAI 2021*. AAAI Press, 2021.
- 150 [2] Wenhua Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou
151 Zhou, and William Yang Wang. Tabfact: A large-scale dataset for table-based fact verification.
152 In *ICLR 2020*. OpenReview.net, 2020.
- 153 [3] Sarthak Dash, Sugato Bagchi, Nandana Mihindukulasooriya, and Alfio Gliozzo. Permutation
154 invariant strategy using transformer encoders for table understanding. In *Findings of the*
155 *Association for Computational Linguistics: NAACL 2022*. Association for Computational
156 Linguistics, 2022.
- 157 [4] Xiang Deng, Huan Sun, Alyssa Lees, You Wu, and Cong Yu. TURL: table understanding
158 through representation learning. *SIGMOD Rec.*, 2022.
- 159 [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of
160 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and
161 Tamar Solorio, editors, *NAACL-HLT 2019*. Association for Computational Linguistics, 2019.
- 162 [6] Julian Eisenschlos, Maharshi Gor, Thomas Müller, and William W. Cohen. MATE: multi-
163 view attention for table transformer efficiency. In Marie-Francine Moens, Xuanjing Huang,
164 Lucia Specia, and Scott Wen-tau Yih, editors, *EMNLP 2021*. Association for Computational
165 Linguistics, 2021.
- 166 [7] Julian Martin Eisenschlos, Syrine Krichene, and Thomas Müller. Understanding tables with
167 intermediate pre-training. In Trevor Cohn, Yulan He, and Yang Liu, editors, *Findings of the*
168 *Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*,
169 volume EMNLP 2020 of *Findings of ACL*, pages 281–296. Association for Computational
170 Linguistics, 2020.
- 171 [8] Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Mar-
172 tin Eisenschlos. Tapas: Weakly supervised table parsing via pre-training. In Dan Jurafsky, Joyce
173 Chai, Natalie Schluter, and Joel R. Tetreault, editors, *ACL 2020*. Association for Computational
174 Linguistics, 2020.
- 175 [9] Hiroshi Iida, Dung Thai, Varun Manjunatha, and Mohit Iyyer. TABBIE: pretrained repre-
176 sentations of tabular data. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek
177 Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao
178 Zhou, editors, *NAACL-HLT 2021*. Association for Computational Linguistics, 2021.
- 179 [10] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun
180 Liu. Tinybert: Distilling BERT for natural language understanding. In Trevor Cohn, Yulan He,
181 and Yang Liu, editors, *Findings of the Association for Computational Linguistics: EMNLP 2020,*
182 *Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*. Association
183 for Computational Linguistics, 2020.
- 184 [11] Qian Liu, Bei Chen, Jiaqi Guo, Morteza Ziyadi, Zeqi Lin, Weizhu Chen, and Jian-Guang Lou.
185 TAPEX: table pre-training via learning a neural SQL executor. In *ICLR 2022*. OpenReview.net,
186 2022.
- 187 [12] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,
188 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT
189 pretraining approach. *CoRR*, 2019.
- 190 [13] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version
191 of BERT: smaller, faster, cheaper and lighter. *CoRR*, 2019.
- 192 [14] Nan Tang, Ju Fan, Fangyi Li, Jianhong Tu, Xiaoyong Du, Guoliang Li, Samuel Madden,
193 and Mourad Ouzzani. RPT: relational pre-trained transformer is almost all you need towards
194 democratizing data preparation. *Proc. VLDB Endow.*, 2021.
- 195 [15] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey.
196 *CoRR*, 2020.

- 197 [16] Edward Wagstaff, Fabian Fuchs, Martin Engelcke, Ingmar Posner, and Michael A. Osborne.
 198 On the limitations of representing functions on sets. In Kamalika Chaudhuri and Ruslan
 199 Salakhutdinov, editors, *ICML 2019*, Proceedings of Machine Learning Research. PMLR, 2019.
- 200 [17] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep
 201 self-attention distillation for task-agnostic compression of pre-trained transformers. In Hugo
 202 Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin,
 203 editors, *NeurIPS 2020*, 2020.
- 204 [18] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony
 205 Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer,
 206 Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain
 207 Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-
 208 art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods
 209 in Natural Language Processing: System Demonstrations*. Association for Computational
 210 Linguistics, 2020.
- 211 [19] Tianbao Xie, Chen Henry Wu, Peng Shi, Ruiqi Zhong, Torsten Scholak, Michihiro Yasunaga,
 212 Chien-Sheng Wu, Ming Zhong, Pengcheng Yin, Sida I. Wang, Victor Zhong, Bailin Wang,
 213 Chengzu Li, Connor Boyle, Ansong Ni, Ziyu Yao, Dragomir R. Radev, Caiming Xiong,
 214 Lingpeng Kong, Rui Zhang, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. Unifiedskg:
 215 Unifying and multi-tasking structured knowledge grounding with text-to-text language models.
 216 *CoRR*, 2022.
- 217 [20] Pengcheng Yin, Graham Neubig, Wen-tau Yih, and Sebastian Riedel. Tabert: Pretraining for
 218 joint understanding of textual and tabular data. In Dan Jurafsky, Joyce Chai, Natalie Schluter,
 219 and Joel R. Tetreault, editors, *ACL 2020*. Association for Computational Linguistics, 2020.
- 220 [21] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov,
 221 and Alexander J Smola. Deep sets. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach,
 222 R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing
 223 Systems*. Curran Associates, Inc., 2017.

224 A Appendix: Experiment Settings

225 We use the huggingface [18] implementation of transformer models including BERT and TAPAS. All
 226 models share the parameter of virtual batch size 64, learning rate 2×10^{-5} , weight decay 10^{-5} , the
 227 AdamW optimizer, and cosine annealing with 2 warm-up epochs. The training process uses early
 228 stopping of patience 8. All experiments are run in half-precision on a cloud server with a single
 229 NVIDIA Tesla V100 Tensor Core GPU.

230 All features are extended to the same dimension of 2048 by a transformation module, which is a
 231 two-layer neural network with hidden size equal to the original feature dimension (768 in case of
 232 BASE models), LeakyReLU activation with negative slope 0.01 and dropout probability 0.1. All
 233 models share the same downstream binary classifier, which is a three-layer neural network with
 234 hidden size 2048, LeakyReLU activation with negative slope 0.01 and dropout probability 0.1. The
 235 query encoder is a separate transformer BASE model.