

---

# Virtual Task Selection in Meta-Learning for Domain Generalization in Semantic Parsing

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Recent work has explored applying meta-learning strategies to improve domain  
2       generalization in natural language processing tasks. We propose a simple strategy  
3       for improving such model training with meta-learning algorithms through strategic  
4       sampling of virtual tasks. We experiment with our task selection strategies on a  
5       challenging text-to-SQL semantic parsing dataset and show that it is possible to  
6       achieve performance comparable to a strong baseline training strategy while using  
7       half as many training steps. Furthermore, comparison with a strong multi-task  
8       learning baseline suggests that the benefits of our strategies are complementary to  
9       the benefits conferred by meta-learning itself.

## 10   1 Introduction

11   Recent studies have shown that meta-learning can be used to enhance domain generalization across nu-  
12   merous vision and language tasks (Balaji et al., 2018; Li et al., 2019). In particular, optimization-based  
13   Model-Agnostic Meta-Learning (MAML) seeks to optimize domain generalization by simulating  
14   zero-shot learning during training through virtual source and target tasks (Finn et al., 2017; Li et al.,  
15   2018). However, little work has explored the effects of virtual task selection for MAML. In this  
16   work, we explore how redefining domains for source/target tasks and utilizing different task selection  
17   methods affects the meta-training process with respect to model performance. In particular, we  
18   test our methods in the realm of semantic parsing, which is the task of translating natural language  
19   utterances into formal meaning representations (e.g., logical forms, database queries, programs).  
20   Domain generalization in semantic parsing involves producing new programs from previously unseen  
21   natural language utterances stemming from new domains and databases. This is a topic of interest  
22   since gathering data to train semantic parsers is difficult and data during inference can be largely  
23   different from data available during training (Chang et al., 2020; Huang et al., 2018; Conklin et al.,  
24   2021; Chen et al., 2020).

25   We experiment with different strategies for sampling virtual tasks during meta-learning for Spider, a  
26   zero shot text-to-SQL semantic parsing dataset<sup>1</sup> (Yu et al., 2019). The Spider task is a particularly  
27   interesting setting to study meta-learning strategies in for a number of reasons. There is value in  
28   studying the effects of meta-learning in a challenging semantic parsing task where adaptation to  
29   unseen domains is an explicit goal inherent to the task, while simultaneously, semantic parsing tasks  
30   are relatively under-explored compared to image classification tasks where meta-learning is most  
31   commonly applied. Meanwhile, the task is relatively tractable despite its difficulty; the domains

---

<sup>1</sup>Our repo is based on Wang et al. (2021)’s implementation in <https://github.com/berlino/tensor2struct-public> and will be available at <https://github.com/clarana/tensor2struct-meta-domain>

32 in the dataset are well-defined (different SQL databases correspond to different domains), with an  
33 underlying similarity in desired logical form structure (i.e., acceptable model outputs are well-formed  
34 SQL queries which produce correct outputs when executed).

35 Similar to Wang et al. (2021), we create virtual zero-shot tasks during meta-training by sampling  
36 examples from source and target domains pre-defined for Spider. We experiment with alternative  
37 domain definitions based on different similarity scores between natural language text and its cor-  
38 responding programs, and consider different source and target domain selection strategies during  
39 training in addition to a random baseline: namely choosing similar domains, different domains,  
40 increasingly different domains, or increasingly similar domains. We show that it is possible to achieve  
41 full benefits of DG-MAML while drastically reducing training time needed to achieve the results, or  
42 outperform baseline DG-MAML using our domain and task selection strategies. Other notable results  
43 include the capacity of "bad" virtual task sampling strategies to dramatically reduce task accuracy,  
44 and the stabilizing effect of DG-MAML versus a multitask learning baseline.

## 45 2 Methods

### 46 2.1 Redefining Domains

47 Text-to-SQL datasets contain pre-defined domains since queries tend to reference a single database  
48 (i.e., a database is a domain). In particular, the publicly available Spider dataset contains 166  
49 databases including 146 for the training set and 20 for the validation set. However, samples from the  
50 same domain may not be homogeneous when considering text representations (see Figure 5). Since  
51 meta-learning assumes that samples in source and target domains are cohesive, we re-define data  
52 domains in the following ways:

53 **Similarity of Question representation:** the text-SQL pairs in the Spider training set are clustered  
54 into 220 new domains using KNN based on cosine similarity between the natural language question  
55 representations.

56 **Similarity of SQL representation:** the training samples are clustered into 220 new domains based  
57 on SQL similarity using the same method.

58 **Similarity of Question-SQL representation similarity:** domains are redefined by the degree to  
59 which a natural language text is similar to its SQL representation. We calculate the cosine similarity  
60 between a text embedding and its corresponding SQL embedding. The new domains are the cosine  
61 similarities defined by increments of 0.1. The re-defined domains represent degrees of semantic and  
62 syntactic similarity and serve as a fine-grained proxy for training example *difficulty*.

63 The above ways of redefining domains purely rely on the question or SQL representations, which  
64 ignores some orthogonal features in the table schemes. We also propose another naive deconvolution  
65 method which we deconvolve domain-specific and domain-agnostic representation (see A.1), making  
66 use of both the similarity of representation and the pre-defined domain information. Hence we also  
67 propose an alternative method below.

68 **Similarity of Domain-Specific Question Representation:** The text-SQL pairs in the training set  
69 are clustered into 136 new domains using KNN and cosine similarity between the deconvolved  
70 domain-specific question representations.

### 71 2.2 Virtual Task Selection

72 We explore task selection strategies to potentially improve meta-learning performance by: 1) Consis-  
73 tently selecting similar domains; 2) Consistently selecting distinct domains; 3) Selecting domains  
74 with decreasing similarity; 4) Selecting domains with increasing similarity, and; 5) Random selection.

75 To quantify similarity between domains, we calculate the cosine similarity matrix for each pair of  
76 domains, where the mean embedding for each domain is used. In choosing a target domain for a  
77 source domain  $d \in D$  during meta-training, we retrieve the vector  $\vec{s}$  of similarity between that  $d$   
78 and all other domains  $D - d$  from the cosine similarity matrix, which we use to obtain a probability

79 distribution under which we pick the target domains. Formally, in each step  $i \in 1, \dots, n$  where  $n$  is the  
 80 maximum training steps, we calculate the target domain with the probability distribution,  $P(t)$ , as:

$$P(t) = \text{softmax}(\vec{s}^r) \quad \text{where} \quad (1)$$

$$r = \begin{cases} 1 & \text{if selection strategy 1 (constant rate)} \\ -1 & \text{if selection strategy 2 (constant rate)} \\ -\frac{2i}{n} + 1 & \text{if selection strategy 3 (linear rate)} \\ \frac{2i}{n} - 1 & \text{if selection strategy 4 (linear rate)} \end{cases}$$

### 81 3 Experiments

82 We experiment with RAT-SQL (Wang et al., 2019), a model with a relation-aware transformer encoder  
 83 and a LSTM-based decoder, as our baseline model. We use similar configurations as Wang et al.  
 84 (2021), including their batch size of  $B_s = B_t = 12$  and their inner and outer learning rates of  
 85  $\alpha = 5 \times 10^{-4}$  and  $6 \times 10^{-4}$ , where their outer learning rate schedule includes a 500 step warmup to  
 86  $6 \times 10^{-4}$  and a polynomial cooldown to their last step.

87 We run meta-training using the DG-MAML algorithm using combinations of our three domain  
 88 definitions and our five virtual task selection methods. In addition, we proposed an alternative,  
 89 stronger baseline of multi-task learning (Multi-task in Table 1), which allows us to implement task  
 90 selection on top of regular supervised training. This allows for an additional fairer comparison  
 91 between meta-learning and regular supervised training when comparing task selection strategies  
 92 across models. Details of this model can be found in Appendix A. Select results are depicted in  
 93 Table 1, with additional tables ( A.2, 2) containing more complete results. We run a subset of  
 94 experiments to 20k training steps and report those results in Table 3.

Model	Set Match Acc.	Exec. Acc.
<i>DG-MAML Wang et al. (2021), 20k train steps</i>	68.9	69.3
<i>DG-FMAML Wang et al. (2021), 20k train steps</i>	N/A	N/A
Regular	62.3	65.0
Multi-task + orig domain + .	64.5	66.6
DG-FMAML + orig domain + .	66.8	67.3
DG-MAML + orig domain + .	66.8	67.6
DG-FMAML + orig domain + different	67.5	68.8
DG-MAML + orig domain + similar	68.7	<b>69.2</b>
DG-MAML + question sim domains + different → similar	66.7	68.3
DG-MAML + question sim domains + similar → different	68.1	66.7
DG-MAML + domain-specific sim domains + dif- ferent	<b>68.8</b>	68.3
DG-MAML + domain-specific sim domains + sim- ilar → different	68.7 / 67.3	68.6 / 68.1

Table 1: Evaluation accuracy % on baseline and DG-MAML model with 10k training iterations, except the first section, which contains lines *in italics* which are reported dev accuracies from Wang et al. (2021) (DG-FMAML accuracy is unreported for English-only Spider). The second section contains our own baseline metrics, and the third section contains our most promising models. Each model name consists of three parts: The first component is the learning scheme, the second is the domain definition scheme, and the third is the task selection scheme, where ‘.’ indicates the default selection scheme – otherwise, similarity calculations among domains are based on example embeddings. Note that if task selection scheme is random, no similarity calculation is used.

## 95 4 Discussion

96 Our methods achieve varying levels of performance on the Spider task. Experiments that boost  
97 performance above the baseline at 10k steps in some cases approach the performance of Wang et al.  
98 (2021)’s reported results (68.9% and 69.3% accuracy for validation set exact set match and execution  
99 accuracy) at 20k steps, using the same batch sizes and other settings (Tables 1, 2). Additionally, we  
100 found that intentional task selection did not impact training time (the bottleneck is still by far the time  
101 required to compute first and second order gradients during meta-training), so halving the number of  
102 training steps also halves the wall clock time. Meanwhile, there are multiple methods that dramatically  
103 hurt accuracy compared to baselines, such as domain redefinition based on example-example text  
104 cosine similarity with a static sampling strategy (Table A.2 in Appendix), further confirming that task  
105 selection strategy does in fact affect model performance. Furthermore, models trained with multi-task  
106 learning or DG-FMAML often match or even exceed performance of DG-MAML models. However,  
107 meta-learning seems to encourage faster convergence (with respect to the number of steps) and may  
108 confer benefits in *stability* throughout individual training runs (Table A.2 in Appendix).

109 Our domain redefinition experiments are also instructive: When the domains are redefined by question  
110 embedding similarity, shifting from selecting similar to dissimilar domains or vice versa throughout  
111 training can improve the performance compared to using a consistent selection strategy. Similar-to-  
112 dissimilar sampling tends to outperform its inverse when the text-SQL similarity is high and worse  
113 when text-SQL similarity is low (Fig. 1 in Appendix). When using the question domain-specific  
114 representation as the redefined domains, selecting different samples constantly or shifting from  
115 similar to dissimilar domains boosts the performance.

116 When the domains are redefined by text-SQL similarity, selecting similar source and target domains  
117 during training results in higher accuracies when generating hard and extra hard queries compared to  
118 selecting dissimilar domains during training (Fig. 7, 8 in Appendix). Additionally, curriculum learning  
119 in this setting seems to equalize accuracies over the entire range of queries with different levels of  
120 text-SQL similarity, a proxy for example difficulty. Similar-to-dissimilar curriculum learning has a  
121 more prominent effect than dissimilar-to-similar curriculum learning (Fig. 3, Fig. 4 in Appendix).

### 122 4.1 Future work

123 **Further "tuning" of models trained with redefined domain selection** Despite meta-training  
124 being a strategy to aim domain generalization, there still may be inherently different distributions  
125 of data between training and evaluation or inference time. We can explore meta-training schemes  
126 that involve a final tuning of the model using originally defined domains, after meta-training using  
127 redefined domains.

128 **Experimentation with alternative curricula** In our experiments that involve curriculum learning,  
129 the rate in which domain similarity selection increases or decreases is linear with respect to the  
130 number of steps. However, it may be beneficial to explore non-linear rates or injecting outliers into  
131 the curriculum during training.

132 **Beyond text-SQL** Another potential future direction is to evaluate DG-MAML on other zero-  
133 shot tasks where domain is not as well-defined as the case where a domain is a SQL table. Our  
134 experimentation with redefined domains lays groundwork for domain selection in cases where  
135 domains may be defined in different ways – our results suggest that domain selection strategy does  
136 matter when meta-learning for domain generalization, and that different strategies may be beneficial  
137 depending on how domains are defined. A dynamic task selection strategy could be especially helpful  
138 to improve domain generalization in a model-agnostic, task-agnostic way.

### 139 4.2 Conclusion

140 Previous work suggests that a good "curriculum" in meta-learning is important (Zhan et al., 2021;  
141 Zhang et al., 2019). Having well-defined inherent domains for a challenging task, as they are in  
142 Spider (Yu et al., 2019), can help us ground experimentation and learn from and interpret results.  
143 This can in turn help us design human-interpretable curricula for meta-learning, potentially even for  
144 tasks beyond those specifically studied.

145 **References**

- 146 Yogesh Balaji, Swami Sankaranarayanan, and Rama Chellappa. 2018. Metareg: Towards domain  
147 generalization using meta-regularization. In *Advances in Neural Information Processing Systems*,  
148 volume 31. Curran Associates, Inc.
- 149 Shuaichen Chang, Pengfei Liu, Yun Tang, Jing Huang, Xiaodong He, and Bowen Zhou. 2020. Zero-  
150 shot text-to-sql learning with auxiliary task. In *In The Thirty-Fourth AAAI Conference on Artificial  
151 Intelligence, AAAI 2020*.
- 152 Xilun Chen, Asish Ghoshal, Yashar Mehdad, Luke Zettlemoyer, and Sonal Gupta. 2020. Low-resource  
153 domain adaptation for compositional task-oriented semantic parsing.
- 154 Henry Conklin, Bailin Wang, Kenny Smith, and Ivan Titov. 2021. Meta-learning to compositionally  
155 generalize. In *Proceedings of the 59th Annual Meeting of the Association for Computational  
156 Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume  
157 1: Long Papers)*, pages 3322–3335, Online. Association for Computational Linguistics.
- 158 Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast  
159 adaptation of deep networks.
- 160 Po-Sen Huang, Chenglong Wang, Rishabh Singh, Wen-tau Yih, and Xiaodong He. 2018. Natural  
161 language to structured query generation via meta-learning.
- 162 Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy Hospedales. 2018. Learning to generalize: Meta-  
163 learning for domain generalization. *Proceedings of the AAAI Conference on Artificial Intelligence*,  
164 32(1).
- 165 Yiyang Li, Yongxin Yang, Wei Zhou, and Timothy M. Hospedales. 2019. Feature-critic networks for  
166 heterogeneous domain generalization.
- 167 Bailin Wang, Mirella Lapata, and Ivan Titov. 2021. Meta-learning for domain generalization in  
168 semantic parsing. In *Proceedings of the 2021 Conference of the North American Chapter of  
169 the Association for Computational Linguistics: Human Language Technologies*, pages 366–379,  
170 Online. Association for Computational Linguistics.
- 171 Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2019.  
172 Rat-sql: Relation-aware schema encoding and linking for text-to-sql parsers.
- 173 Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li,  
174 Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2019. Spider: A large-scale  
175 human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task.
- 176 Runzhe Zhan, Xuebo Liu, Derek F Wong, and Lidia S Chao. 2021. Meta-curriculum learning for  
177 domain adaptation in neural machine translation. In *Proceedings of the AAAI Conference on  
178 Artificial Intelligence*, volume 35, pages 14310–14318.
- 179 Xuan Zhang, Pamela Shapiro, Gaurav Kumar, Paul McNamee, Marine Carpuat, and Kevin Duh.  
180 2019. Curriculum learning for domain adaptation in neural machine translation. In *Proceedings  
181 of the 2019 Conference of the North American Chapter of the Association for Computational  
182 Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1903–1915,  
183 Minneapolis, Minnesota. Association for Computational Linguistics.

184 **Checklist**

- 185 1. For all authors...
- 186 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
187 contributions and scope? [Yes]
- 188 (b) Did you describe the limitations of your work? [Yes] We mentioned what we can do  
189 for future work
- 190 (c) Did you discuss any potential negative societal impacts of your work? [N/A]

- 191 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
192 them? [Yes]
- 193 2. If you are including theoretical results...
- 194 (a) Did you state the full set of assumptions of all theoretical results? [N/A]  
195 (b) Did you include complete proofs of all theoretical results? [N/A]
- 196 3. If you ran experiments...
- 197 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
198 perimental results (either in the supplemental material or as a URL)? [No] We  
199 will release a public version of our repository along with camera-ready revisions.  
200 Our implementations are built on top of Wang et al. (2021)’s implementations at  
201 <https://github.com/berlino/tensor2struct-public>
- 202 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
203 were chosen)? [Yes] in §3
- 204 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
205 ments multiple times)? [Yes] When applicable (Table 3 in Appendix)
- 206 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
207 of GPUs, internal cluster, or cloud provider)? [No] We did not strictly control for the  
208 hardware we used across all experiments, but our wall clock times for training were  
209 generally in line with those reported by Wang et al. (2021). Furthermore, our strategies  
210 theoretically halve the training time required to reach previously reported baseline  
211 performance given any adequate compute resources.
- 212 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 213 (a) If your work uses existing assets, did you cite the creators? [Yes]  
214 (b) Did you mention the license of the assets? [No]  
215 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
216
- 217 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
218 using/curating? [N/A]
- 219 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
220 information or offensive content? [N/A]
- 221 5. If you used crowdsourcing or conducted research with human subjects...
- 222 (a) Did you include the full text of instructions given to participants and screenshots, if  
223 applicable? [N/A]
- 224 (b) Did you describe any potential participant risks, with links to Institutional Review  
225 Board (IRB) approvals, if applicable? [N/A]
- 226 (c) Did you include the estimated hourly wage paid to participants and the total amount  
227 spent on participant compensation? [N/A]

## 228 A Appendix

229 Optionally include extra information (complete proofs, additional experiments and plots) in the  
230 appendix. This section will often be part of the supplemental material.

### 231 A.1 Domain-specific and domain-agnostic representation

232 For datasets with domain categorical information, after obtaining an initial embedding for data  $X^{n \times k}$   
233 and domain label  $y \in \{1, \dots, C\}^n$ , we form the question as follows:

234 For domain-specific embedding, we want to find a transformation that minimizes the in-domain  
235 variance. Without loss of generality, we assume that we want to find a one-dimensional representation.  
236 We define the mean of  $X$  as  $\hat{X}$  such that

$$\hat{x}_i = \frac{1}{\sum_{j=1}^n \mathbf{1}_{\{y_j=y_i\}}} \sum_{j=1}^n \mathbf{1}_{\{y_j=y_i\}} x_j$$

237 We then define the objective as follows:

$$\min_p p^T (X - \hat{X})^T (X - \hat{X}) p, \text{ where } \|p\|_2 = 1$$

238 Then for domain-agnostic (syntax-specific) embeddings, we want to find a transformation that  
 239 maximizes the in-domain variance to account for the variations within the syntax for each domain.

240 We then define the objective as follows:

$$\max_q q^T (X - \hat{X})^T (X - \hat{X}) q, \text{ where } \|q\|_2 = 1$$

241 The dimension of the embedding can be more than  $r = 1$ .

242 Then the final transformation for domain-specific embeddings and domain-agnostic embeddings  
 243 would be

$$\begin{aligned} z^{spe} &= Xp \\ z^{agn} &= Xq \end{aligned}$$

244 In theory,  $p$  and  $q$  should be the last column and first column of the eigenvector matrix in the  
 245 eigendecomposition of  $X - \hat{X}$ .

## 246 A.2 Multi-task learning baseline

247 Our multi-task learning implementation is based on the DG-MAML implementation and less efficient  
 248 than it can be, but it nevertheless cuts training time in half (38 hours vs 19 hours). For multi-task  
 249 learning, we use:

$$\nabla_{\theta} \mathcal{L}_{\mathcal{T}}(\theta) = \nabla_{\theta} \mathcal{L}_{\mathcal{B}_t}(\theta) + \nabla_{\theta} \mathcal{L}_{\mathcal{B}_s}(\theta) \quad (2)$$

250 Furthermore, we implement DG-FMAML, the first order approximation of DG-MAML described  
 251 but not provided by Wang et al. (2021) and run additional experiments using it to explore a computa-  
 252 tionally less intensive alternative to DG-MAML. In DG-FMAML, we use:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\mathcal{T}}(\theta) &= \nabla_{\theta} \theta' \nabla_{\theta'} \mathcal{L}_{\mathcal{B}_t}(\theta') + \nabla_{\theta} \mathcal{L}_{\mathcal{B}_s}(\theta) \\ &\approx \nabla_{\theta'} \mathcal{L}_{\mathcal{B}_t}(\theta') + \nabla_{\theta} \mathcal{L}_{\mathcal{B}_s}(\theta) \end{aligned} \quad (3) \quad (4)$$

253 Because our implementation adds gradients manually from a copy to the original, our DG-FMAML  
 254 may require more memory at times, but training time is cut roughly in half from the DG-MAML  
 255 implementation. In Table 1, we highlight our best performing models and compare them with each  
 256 other as well as with Wang et al. (2021)’s reported results.

Learning algorithm	Domain Definition	Domain Selection Strategy	Set Match Acc.	Exec. Acc.	
Regular supervised	N/A	N/A	62.3	65.0	
Multi-task learning	SQL table	similar (table name)	64.5 / 64.9	66.6 / 66.7	
		similar (text)	66.5	66.3 / 67.2	
		different (text)	66.6	67.8	
		similar → different (text)	67.7	66.8 / 66.9	
		different → similar (text)	68.6 / 69.1	67.4	
DG-FMAML	SQL table	similar (table name)	66.8	67.3	
		random	67.2	68.3	
		similar (text)	65.0	67.3	
		different (text)	67.5 / 67.6	68.8	
		similar → different (text)	66.6 / 66.7	67.8 / 68.4	
		different → similar (text)	66.2	67.3.	
DG-MAML	SQL table	similar (table name)	66.8	67.6	
		random	67.1	67.1	
		similar (text)	<b>68.7</b>	<b>69.2</b>	
		different (text)	67.9	68.9	
		similar → different (text)	67.6 / 68.2	67.9 / 68.1	
		different → similar (text)	67.6	67.9 / 68.0	
	example-example text cosine similarity	random	66.8	66.8	
		similar (text)	64.0 / 65.4	64.6 / 66.2	
		different (text)	65.1	66.2	
		similar → different (text)	68.1	66.7	
		different → similar (text)	66.7	68.3	
		text-SQL cosine similarity	similar (text)	65.4	65.3
			different (text)	62.8	62.7 / 63.2
			similar → different (text)	65.1 / 64.9	64.9 / 65.2
different → similar (text)	62.5 / 63.7		64.8 / 66.3		



Learning algorithm	Domain Definition	Domain Selection Strategy	Set Match Acc.	Exec. Acc.
DG-MAML	example-example domain specific representation cosine similarity	random	67.3	67.4
		similar (text mean cosim)	66.3 / 66.3	67.0 / 66.5
		different (text mean cosim)	<b>68.8</b>	68.3
		similar → different (text mean cosim)	68.7 / 67.3	<b>68.6</b> / 68.1
		different → similar (text mean cosim)	66.7 / 65.7	66.8 / 66.7
	sql-sql domain specific representation cosine similarity	random	65.7 / 65.7	66.2 / 66.3
		similar (sql mean cosim)	65.7	66.2
		different (sql mean cosim)	65.7	66.7
		similar → different (sql mean cosim)	62.1	63.2
		different → similar (sql mean cosim)	61.3 / 61.9	62.3 / 62.7

Table 2: This description applies to both Table A.2 and Table 2. Validation set accuracy % on dev set with after 10k training iterations for all models trained. All models listed use a RAT-SQL base with BERT-base contextualized embeddings. Wang et al. (2021) define data domains by SQL table of origin (*SQL table*) and by default select domains weighted by cosine similarity of SQL table name embeddings so that domains with similar names are more likely to be selected for a task (*similar table name cosim*). We mark accuracy metrics at 'final checkpoint / best checkpoint' when the final checkpoint does not have the highest accuracy. Checkpoints were saved every 1000 steps starting from the 1000th step.

Model	Set Match Accuracy (10k $\rightarrow$ 20k steps)	Execution Accuracy (10k $\rightarrow$ 20k steps)
<i>Regular Wang et al. (2021)</i>	<i>N/A <math>\rightarrow</math> 66.8</i>	<i>N/A <math>\rightarrow</math> 66.8</i>
<i>DG-MAML Wang et al. (2021)</i>	<i>N/A <math>\rightarrow</math> 68.9</i>	<i>N/A <math>\rightarrow</math> 69.3</i>
Regular	N/A $\rightarrow$ 67.4 $\pm$ 0.8	N/A $\rightarrow$ 67.5 $\pm$ 1.1
Multi-task + orig domain + .	64.5 $\rightarrow$ 69.8 $\pm$ 1.6	66.6 $\rightarrow$ 69.3 $\pm$ 1.1
DG-FMAML + orig domain + .	66.8 $\rightarrow$ *69.9	67.3 $\rightarrow$ *69.2
DG-MAML + orig domain + .	66.8 $\rightarrow$ 68.2 $\pm$ 1.3	67.6 $\rightarrow$ 69.4 $\pm$ 1.3
DG-MAML + orig domain + similar	68.7 $\rightarrow$ 69.4 $\pm$ 1.0	69.2 $\rightarrow$ 69.7 $\pm$ 1.2
DG-MAML + orig domain + different	67.9 $\rightarrow$ 67.8 $\pm$ 0.4	68.9 $\rightarrow$ 68.4 $\pm$ 1.2

Table 3: Evaluation accuracy % on baseline and DG-MAML models trained for 20k iterations. Accuracies reached in (separate) training runs to 10k steps are included for comparison. Reported numbers for 20k step training runs are arithmetic means ( $n = 3$ )  $\pm$  standard deviations, unless \*otherwise specified. Comparison of 10k step accuracy with 20k step accuracy can hint at the rate of convergence using different training schemes. Although variability is high relative across random seeds, many of the differences we report, including the improvement in set match accuracy from selecting similar domains throughout training (vs baseline DG-MAML), are statistically significant (p-value  $<$  0.05).

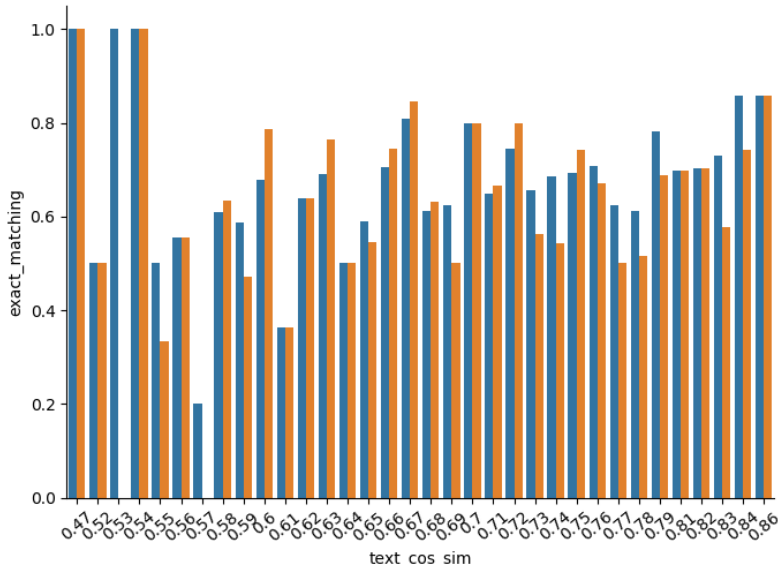


Figure 1: The accuracy among tasks with different cos text-SQL similarity using curriculum learning (blue) and reverse curriculum learning (orange) with redefined domains based on question embedding similarity

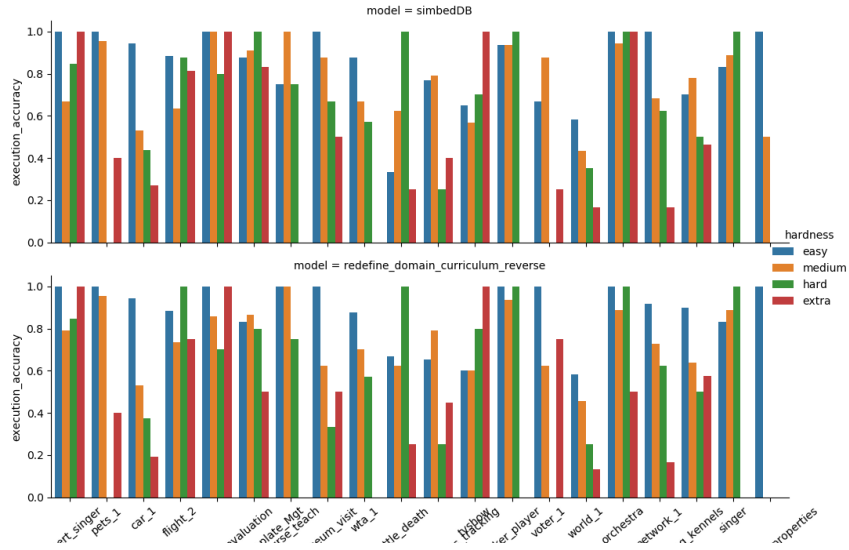


Figure 2: The execution accuracy of predicted SQL with consistently selecting similar domain based on mean question embedding similarity with original domain and reverse curriculum learning with redefined domains based on question embedding similarity

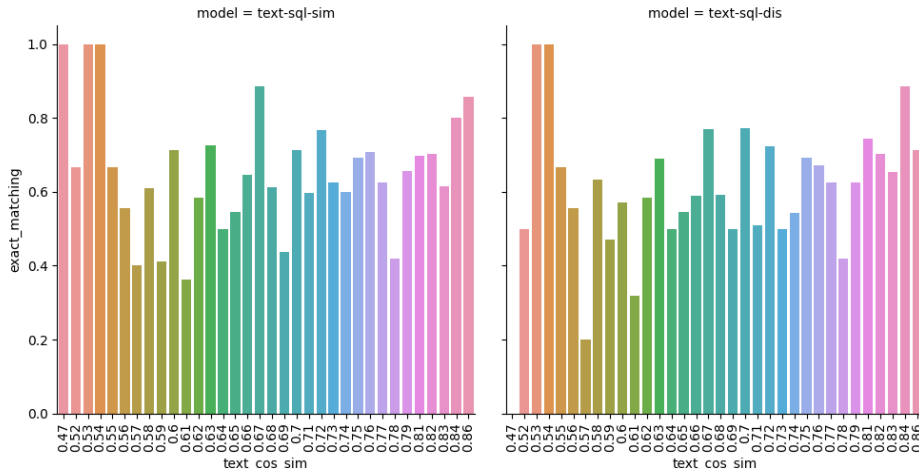


Figure 3: Accuracy among tasks with different cosine similarities between text and SQL. Source and target domains during training are redefined by cosine similarities between text and SQL. Left: similar domain are selected as source and target domains during training. Right: dissimilar domain are selected as source and target domains during training

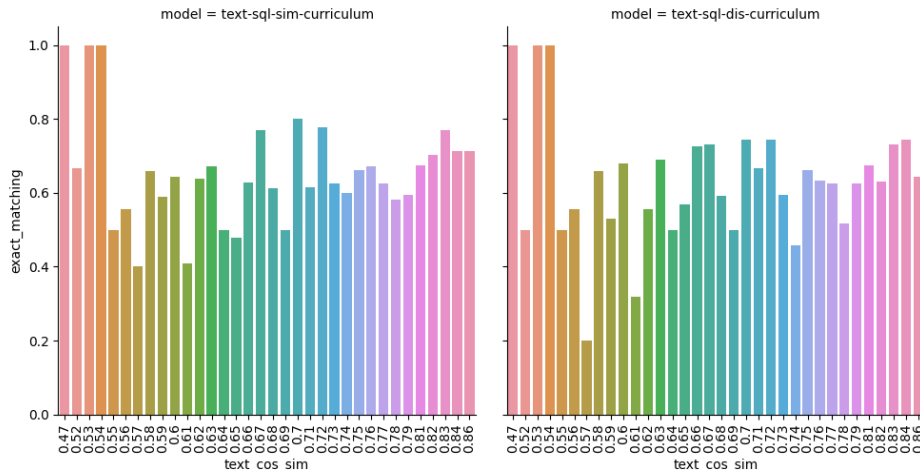


Figure 4: Accuracy among tasks with different cosine similarities between text and SQL. Source and target domains during training are redefined by cosine similarities between text and SQL. Left: curriculum learning is used (decreasing source and target domain similarity over time during training). Right: reverse curriculum learning is used (increasing source and target domain similarity over time during training)

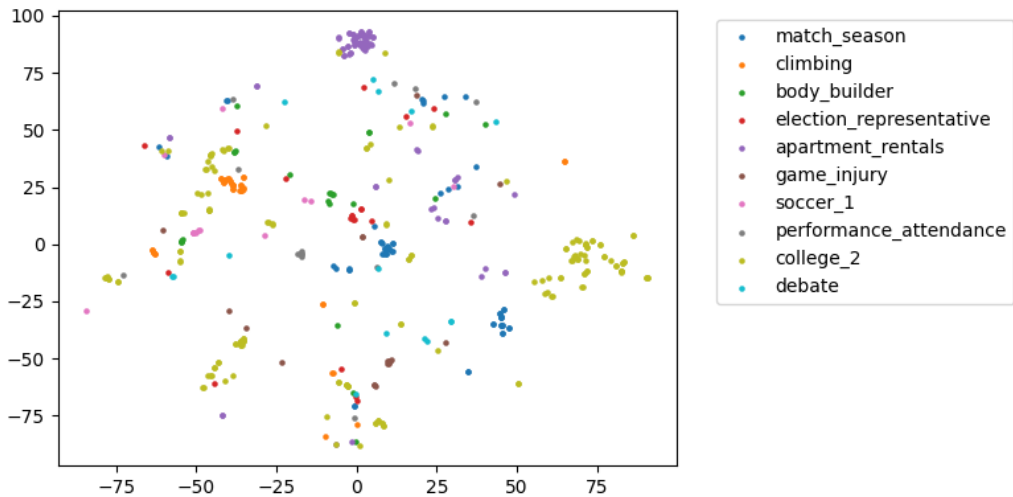


Figure 5: tSNE plot for sentence representation of questions from 10 example domains

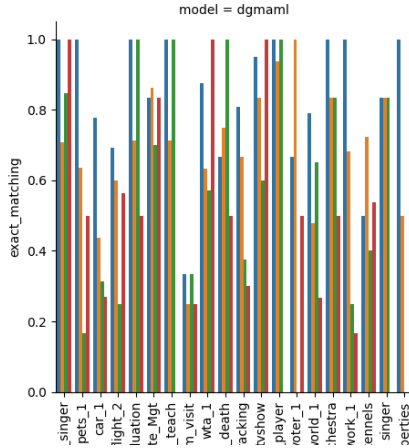


Figure 6: Accuracy among domains defined by hardness of SQL query using DGMAML

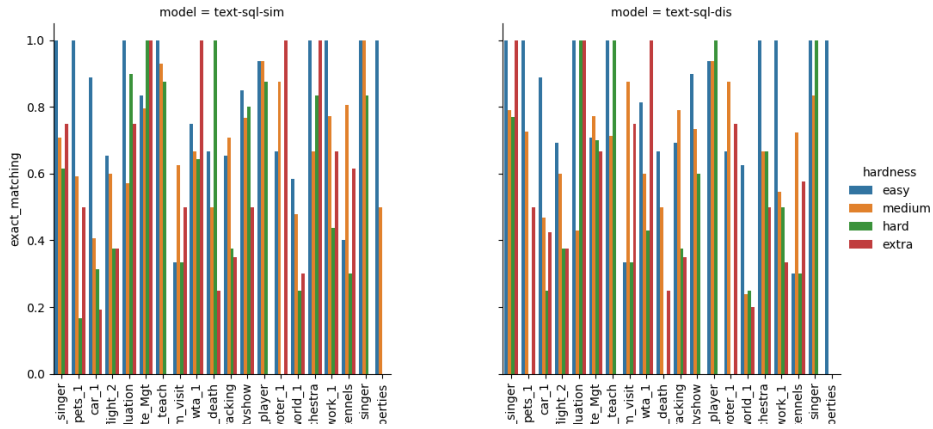


Figure 7: Accuracy among domains defined by hardness of SQL query. Source and target domains during training are redefined by cosine similarities between text and SQL. Left: similar domain are selected as source and target domains during training. Right: dissimilar domain are selected as source and target domains during training

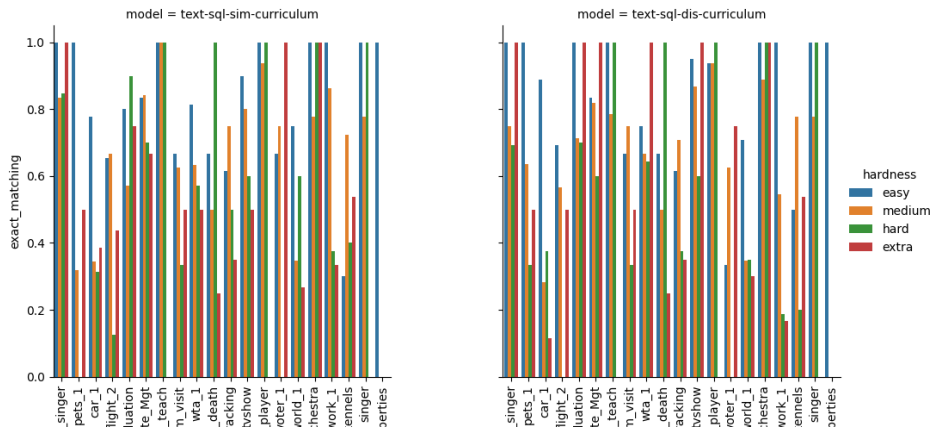


Figure 8: Accuracy among domains defined by hardness of SQL query. Source and target domains during training are redefined by cosine similarities between text and SQL. Left: curriculum learning is used (decreasing source and target domain similarity over time during training). Right: reverse curriculum learning is used (increasing source and target domain similarity over time during training)