

Measuring and Improving Compositional Generalization in Text-to-SQL via Component Alignment

Anonymous ACL submission

Abstract

001 Recently, the challenge of *compositional generalization* in NLP has attracted more and
002 more attention. Specifically, many prior works
003 show that neural networks struggle with com-
004 positional generalization where training and
005 test distributions differ. However, most of
006 these works are based on word-level synthetic
007 data or a specific data split method to gener-
008 ate compositional biases. In this work, we
009 propose a clause-level compositional example
010 generation method, and we focus on text-to-
011 SQL tasks. We first split the sentences in the
012 Spider text-to-SQL dataset into several sub-
013 sentences, then we annotate each sub-sentence
014 with its corresponding SQL clause, resulting
015 in our new dataset Spider-SS. Building upon
016 Spider-SS, we further construct a new dataset
017 named Spider-CG, by composing Spider-SS
018 sub-sentences to test the ability of models
019 to generalize compositionally. Experiments
020 show that previous models suffer significant
021 performance degradation when evaluated on
022 Spider-CG, even though every sub-sentence is
023 seen during training. To deal with this prob-
024 lem, we modify the RATSQ+GAP model to
025 fit the segmented data of Spider-SS, and we
026 show that this method improves the generaliza-
027 tion performance.¹
028

029 1 Introduction

030 Neural models in supervised learning settings show
031 good performance on data drawn from the training
032 distribution. However, the generalization perfor-
033 mance can be poor on out-of-distribution (OOD)
034 samples (Finegan-Dollak et al., 2018; Suhr et al.,
035 2020; Kaushik et al., 2020; Sagawa et al., 2020).
036 This observation might hold even when the new
037 samples are composed of known constituents; e.g.,
038 on SCAN dataset (Lake and Baroni, 2018), many
039 models provide the wrong predictions for the input
040 “jump twice and walk”, even when “jump twice”,

“walk”, and “walk twice” are seen during training.
This (often lacking) ability to generalize to novel
combinations of elements observed during training
is referred to as *compositional generalization*.

Previous work on compositional generalization
in text-to-SQL focuses on data split and word sub-
stitution (Finegan-Dollak et al., 2018; Shaw et al.,
2021). However, data split methods are limited by
the dataset content, making it difficult to construct
a challenging benchmark while ensuring that every
compound appears in the training set. Ensuring a
reasonable data split may also lead to a reduction
in the size of the dataset, e.g., the training set drop
from 7000 to 3282 in the Spider TCMD split (Yu
et al., 2018b; Shaw et al., 2021).

Word substitution is often used to generate com-
positional variations for data augmentation (Yu
et al., 2021; Andreas, 2020). However, this method
cannot generate more complex sentences, i.e., can-
not generate “jump twice and walk” from the
data only containing “jump twice”, “walk”, and
“walk twice”. In addition, in the *cross-domain*
text-to-SQL setting, examples generated by word
substitution are supposed to be considered as in-
distribution data, since cross-domain models are
expected to generalize to unseen domains with un-
seen utterances, including examples generated by
word substitution.

In this work, we first introduce our Spider-SS
dataset (SS stands for sub-sentence) derived from
the Spider benchmark (Yu et al., 2018b). Figure 1
presents a comparison between Spider and Spider-
SS. To build Spider-SS, we design a sentence split
algorithm to split every Spider sentence into sev-
eral sub-sentences until indivisible. Next, we an-
notate every sub-sentence with its corresponding
SQL clause. In order to reduce the difficulty of an-
notation, we annotate the query in an intermediate
representation called NatSQL (Gan et al., 2021b),
because it is simpler and syntactically aligns better
with natural language (NL).

¹We will release code and dataset upon publication.

Spider Example:	
<i>Sentence:</i>	What type of pet is the youngest animal, and how much does it weigh?
<i>SQL:</i>	<code>SELECT PetType , Weight FROM Pets ORDER BY Pet_Age LIMIT 1</code>
Spider-SS Example:	
<i>SubSentence:</i>	What type of pet
<i>NatSQL:</i>	<code>SELECT Pets.Pettytype</code>
<i>SubSentence:</i>	is the youngest animal
<i>NatSQL:</i>	<code>ORDER BY Pets.Pet_Age LIMIT 1</code>
<i>SubSentence:</i>	, and how much does it weigh?
<i>NatSQL:</i>	<code>SELECT Pets.Weight</code>

Figure 1: A natural language sentence in the original Spider benchmark is split into three sub-sentences in Spider-SS, where each sub-sentence has a corresponding NatSQL clause.

Our annotated Spider-SS provides us with sub-sentences paired with NatSQL clauses, which serve as our compounds. Based on Spider-SS, we further construct the dataset Spider-CG (CG stands for compositional generalization), by substituting the sub-sentences with those from other samples, or composing two sub-sentences to form a more complicated sample. Spider-CG contains two subsets, and we present one example for each subset in Figure 2. The first subset includes 24,134 examples generated by substituting sub-sentences, and we consider most data in this subset as in-distribution. The second subset contains 22,531 examples generated by appending sub-sentences, which increases the length and complexity of the sentence and the SQL query compared to the original samples, and we consider most examples in this subset as OOD. We demonstrate that when models are only trained on the original Spider dataset, they suffer a significant performance drop on the second OOD subset of Spider-CG, even though the domain appears in the training set.

To improve the generalization performance of text-to-SQL models, we modify the RATSQL+GAP+NatSQL (Wang et al., 2020; Shi et al., 2021; Gan et al., 2021b) model so that it can be applied to the Spider-SS dataset, with the model trained sub-sentence by sub-sentence. This modification obtains more than 7.8% accuracy improvement evaluated on the second subset of Spider-CG. To our knowledge, this is the first sub-sentence-based text-to-SQL model.

In short, we make the following contributions:

- Besides the sentence split algorithm, we

Spider-SS :	
<i>SubSentence:</i>	Example-1: What is the name and nation of the singer
<i>NatSQL:</i>	<code>SELECT Singer.Name SELECT Singer.Country</code>
<i>SubSentence:</i>	who have a song having 'Hey' in its name?
<i>NatSQL:</i>	<code>WHERE Concert.Song_Name like '%Hey%'</code>
Example-2:	
<i>SubSentence:</i>	What are the names of the singers
<i>NatSQL:</i>	<code>SELECT Singer.Name</code>
<i>SubSentence:</i>	who performed in a concert in 2014?
<i>NatSQL:</i>	<code>WHERE Concert.Year = 2014</code>
Spider-CG :	
Subset-1: sub-sentence substitution in Example 1 and 2	
<i>Sentence:</i>	What is the name and nation of the singer who performed in a concert in 2014?
<i>NatSQL:</i>	<code>SELECT Singer.Name, Singer.Country WHERE Concert.Year = 2014</code>
Subset-2: Example-1 append a sub-sentence from Example-2	
<i>Sentence:</i>	What is the name and nation of the singer who have a song having 'Hey' in its name and who performed in a concert in 2014?
<i>NatSQL:</i>	<code>SELECT Singer.Name, Singer.Country WHERE Concert.Song_Name like '%Hey%' AND Concert.Year = 2014</code>

Figure 2: Two Spider-CG samples generated by: (1) substituting the sub-sentence with one from another example; or (2) composing sub-sentences from 2 examples in Spider-SS.

introduce Spider-SS, a human-curated sub-sentence-based text-to-SQL dataset built upon the Spider benchmark, by splitting its NL questions into sub-sentences.

- We construct the Spider-CG benchmark for measuring the compositional generalization performance of text-to-SQL models.
- We show that the RATSQL+GAP+NatSQL model can be adapted to sub-sentence-based text-to-SQL data, and that this improves its generalization performance.

2 Spider-SS

2.1 Overview

Figure 1 presents a comparison between Spider and Spider-SS. Unlike Spider, which annotates a whole SQL query to an entire sentence, Spider-SS annotates the SQL clauses to sub-sentences. Spider-SS uses NatSQL (Gan et al., 2021b) instead of SQL because some examples are difficult to annotate using SQL. The Spider-SS provides a combination algorithm that collects all NatSQL clauses and then generates the NatSQL query, where the NatSQL query can be converted into an SQL query.

The purpose of building Spider-SS is to attain clause-level text-to-SQL data to generate more

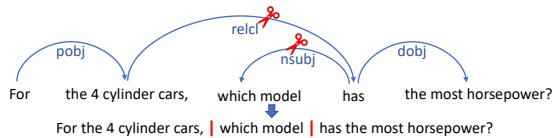


Figure 3: Dependency structure of a sentence and how to split this sentence into three sub-sentences.

complex examples through the combination of clauses. Besides, it is easier to build models with better performance based on more finely labeled data, especially in areas with relatively few resources, i.e., there are only 7000 training samples in the Spider dataset. Consistent with Spider, Spider-SS contains 7000 training and 1034 development examples, but Spider-SS does not contain a test set since the Spider test set is not public. There are two steps to build Spider-SS. First, design a sentence split algorithm to cut the sentence into sub-sentences, and then manually annotate the NatSQL clause corresponding to each sub-sentence.

2.2 Sentence Split Algorithm

We build our sentence split algorithm upon the NL dependency parser spaCy², which provides the grammatical structure of a sentence. Basically, we split the sentence with the following dependencies: *prep*, *relcl*, *advcl*, *acl*, *nsubj*, *npadvmod*, *csubj*, *nsubjpass* and *conj*. According to (de Marnee and Manning, 2016), these dependencies help us separate the main clause, subordinate clauses, and modifiers. Figure 3 shows the dependency structure of a sentence and how to split this sentence into three sub-sentences. However, not every sentence would be split since there are some non-splittable sentences, such as the third example in Figure 4, with the same annotation as the Spider dataset. Although this method can separate sentences well in most cases, due to the variability of natural language, some examples cannot be perfectly split.

To address the remaining issues in sentence split, we design some refinement steps tailored to text-to-SQL applications. For example, when the phrase of a schema column or table is accidentally divided into two sub-sentences, these two sub-sentences are automatically concatenated. Besides, when there is only one word in a sub-sentence, the corresponding split should also be undone.

We sampled 500 examples from the Spider-SS development set to evaluate the acceptability

²<https://github.com/explosion/spaCy>

Spider-SS:

Example-1: Use the "extra" keyword to compensate for split errors

SubSentence: Find the emails and phone numbers of all the customers.

NatSQL: `SELECT Customers.Email_Address
SELECT Customers.Phone_Number`

SubSentence: ordered by email address

NatSQL: `ORDER BY Customers.Email_Address ASC`

SubSentence: and phone numbers.

NatSQL: `EXTRA Customers.Phone_Number`

Example-2: Columns that are not mentioned in the sub-sentence are specifically annotated

SubSentence: List the total number of horses on farms

NatSQL: `SELECT Farm.Total_Horses`

SubSentence: in ascending order.

NatSQL: `ORDER BY Farm.Total_Horses ASC`

NO MENTIONED

Example-3: Some sentences cannot be split

SubSentence: Who advises student 1004?

NatSQL: `SELECT Student.Advisor
WHERE Student.StuID = 2014`

Figure 4: Spider-SS examples in three special cases.

of splitting results manually, and only < 3% of the splitting results are unsatisfactory. For example, in the splitting results of the first example in Figure 4, the last two sub-sentence should be combined to correspond to “**ORDER BY Customer.Email_Address, Customer.Phone_Number ASC**”. In this example, we did not simply give an “**ORDER BY Customer.Phone_Number ASC**” to the last sub-sentence, because it does not mention anything related to “**ORDER BY**”. Here, we introduce “*extra*”, a new NatSQL keyword designed for the Spider-SS dataset, indicating that this sub-sentence mentions a column that temporarily does not fit in any other NatSQL clauses. When combining NatSQL clauses into the final NatSQL query, the combining algorithm determines the final position for the “*extra*” column based on the clauses before and after. Note that even if there is a small proportion of unsatisfactory splitting results, as long as the model trained on Spider-SS can give the correct output according to the input sub-sentence, the quality of the sub-sentences itself does not strongly affect the model utility.

2.3 Data Annotation

When we get the split results from the last step, we can start data annotation. We give precise annotations based on the sub-sentence content, such as the “*extra*” column annotation discussed in the last subsection. Besides, if the description of the schema column is missing in the sub-sentence, we will give the schema column an additional “*NO*”

MENTIONED” mark. For example, in the second example of Figure 4, the “in ascending order” sub-sentence does not mention the “Farm.Total_Horses” column. Therefore, we add a “NO MENTIONED” mark for it. For those sub-sentences that do not mention anything related to the query, we give a “NONE” mark, representing there are no NatSQL clauses.

Since the annotation is carried out according to the sub-sentence content, the equivalent SQL that is more consistent with the sub-sentence will be preferred to the original SQL. Similarly, if the original SQL annotation is wrong, we correct it according to the content.

We annotate the sub-sentence using NatSQL instead of SQL, where NatSQL is an intermediate representation of SQL, only keeping the *SELECT*, *WHERE*, and *ORDER BY* clauses from SQL. Since some sub-sentences need to be annotated with *GROUP BY* clause, we choose the version of NatSQL augmented with *GROUP BY*. We did not use SQL directly because it is difficult to annotate in some cases, such as the SQL in Figure 5. The difficulty is that there are two *SELECT* clauses in this SQL query, but none of the sub-sentences seem to correspond to two *SELECT* clauses. In addition, considering that the two *WHERE* conditions correspond to different *SELECT* clauses, the annotation work based on SQL is far more difficult to complete. As shown in Figure 5, we can use NatSQL to complete the annotation quickly, while the NatSQL can be converted back to the target SQL.

3 Spider-CG

3.1 Overview

Spider-CG is a synthetic dataset, which is generated by recombining the sub-sentences of Spider-SS. There are two recombination methods. The first is sub-sentence substitution between different examples, and the second is to append a sub-sentence into another sentence. To facilitate the follow-up discussion, we named the Spider-CG subset generated by the sub-sentence substitution method **CG-SUB**, and the other named **CG-APP**.

In CG-SUB, there are 21,168 examples generated from the Spider-SS training set, while 2,966 examples are generated from the development set. In CG-APP, examples generated from training and development sets are 19,241 and 3,290, respectively. Therefore, the whole Spider-CG contains 46,665 examples, which is about six times the Spi-

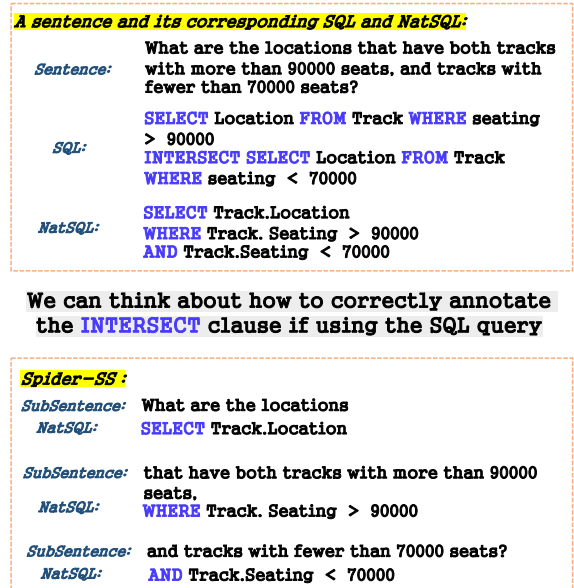


Figure 5: It is difficult to annotate if using the SQL instead of NatSQL.

der dataset. If need more data, we can append sub-sentences to the CG-SUB examples.

3.2 Generation Algorithm

Before generating the Spider-CG, we need to generate the compositional element. Considering that there are some unsatisfactory splitting results in spider-SS, in order to ensure the quality of generated sentences, compositional elements abandoned the examples with “NO MENTIONED” and “extra” column, such as the first and second examples in Figure 4. Each element contains one or more sub-sentences, but these sub-sentences must be annotated as *WHERE* or *ORDER BY* clauses. We collect the sub-sentences for compositional elements by scanning all sub-sentence from top to bottom or from bottom to top and stopping when encountering clauses except *WHERE* and *ORDER BY*. For example, we generate a compositional element containing the last two sub-sentences of the Spider-SS example in Figure 5. In contrast, we do not generate any element from the example in Figure 1.

According to Algorithm 1, we can generate the CG-SUB and CG-APP based on compositional elements. It should be noted that elements in a domain cannot be used in another because the schema items are different. So as many domains as there are, it needs to run Algorithm 1 how many times.

Algorithm 1 Generate CG-SUB and CG-APP dataset in a certain domain

Input: e_list ▷ All compositional elements in a domain
Output: cg_sub and cg_app ▷ CG-SUB and CG-APP dataset in a certain domain

- 1: **for** Every $element_1$ in e_list **do**
- 2: **for** Every $element_2$ in e_list **do**
- 3: **if** $element_1 \neq element_2$ **then**
- 4: **if** $element_1.can_be_substituted_by(element_2)$ **then**
- 5: $cg_sub.append(element_1.generate_substitution_example(element_2))$
- 6: **if** $element_1.can_append(element_2)$ **then**
- 7: $cg_app.append(element_1.generate_appending_example(element_2))$
- 8: **return** cg_sub, cg_app

Ques	List the name and age of the heads of departments are older than 56 ?
SQL	SELECT name, age FROM head WHERE age > 56
Ques	Show the name of employees named Mark Young ?
SQL	SELECT name FROM employee WHERE name = 'Mark Young'

Table 1: Two acceptable but not perfect examples in the Spider-CG.

3.3 Quality Assurance

We randomly sampled 2000 examples from the Spider-CG dataset, around 99% of which are acceptable. Acceptable does not mean that there are no grammatical errors; it means that the sentence has been clearly expressed and does not cause ambiguity. Besides, acceptable examples may be meaningless. We give two acceptable but not perfect examples in Table 1, where the first sentence contains the grammatical error. At the same time, the other is meaningless because the content it wants to query is the condition it gave. We define an acceptable example as an example where you may not find a problem without studying it carefully.

The ‘*can_be_substituted_by*’ and ‘*can_append*’ function in Algorithm 1 prevent mistakes while generating the Spider-CG. These two methods check whether the sub-sentences can be appropriately connected, e.g., whether the previous word is the same and whether the part of speech of the first word in the sub-sentence is the same. Besides, these two methods also prevent generating overly complex or contradictory SQL, e.g., ‘give me the name of the youngest student who is the oldest’.

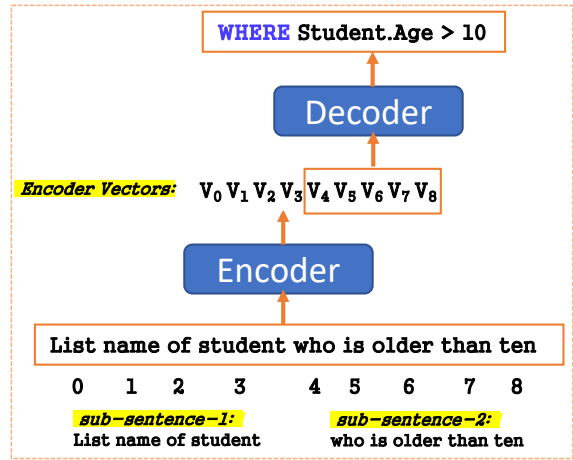


Figure 6: A example of encoding the whole sentence but decoding only the sub-sentence.

4 Model

Existing text-to-SQL models input a sentence and output the corresponding SQL query. So the easiest way to think of using the Spider-SS dataset is to train the model where inputting sub-sentence and outputting the corresponding NatSQL clauses. However, this method is not workable because it will lose some essential schema information. For example, if you only look at the third sub-sentence in Figure 1, you do not know whether he ask about the weight of pets or people.

In order to take into account the context and the sub-sentence data of Spider-SS, we propose that a seq2seq model can encode the whole sentence but decode only the sub-sentence. Figure 6 presents the workflow of encoding the whole sentence but only decoding the sub-sentence of ‘*who is older than ten*’ and outputting the corresponding NatSQL clause. Based on this modification, a seq2seq text-to-SQL model can be adapted to the Spider-SS.

We modify the RATSQL+GAP+NatSQL (Wang

et al., 2020; Shi et al., 2021; Gan et al., 2021b) model to adapt the Spider-SS. We choose this model because its performance is high enough, and it has the best performance among the text-to-SQL models that support NatSQL. We keep the same hyperparameters. Although re-search hyperparameters may improve the performance, it requires considerable computing resources (Wang et al., 2020).

5 Experiment

5.1 Experimental Setup

We evaluate the previous state-of-the-art models on the Spider-CG and Spider (Yu et al., 2018b) datasets. Since the Spider test set is not publicly accessible, Spider-CG does not contain a test set. As discussed in Section 3.1, we divide the Spider-CG into two subsets: CG-SUB and CG-APP. Therefore, there are five evaluation sets:

- **Spider_D**: the original Spider development set with 1,034 examples for *cross-domain in-distribution* text-to-SQL evaluation.
- **CG-SUB_T**: the CG-SUB training set, containing 21,168 examples generated from Spider-SS training set by substituting sub-sentence method. CG-SUB_T is used for *in-domain in-distribution* text-to-SQL evaluation.
- **CG-SUB_D**: the CG-SUB development set containing 2,966 examples for *cross-domain in-distribution* text-to-SQL evaluation.
- **CG-APP_T**: the CG-APP training set, containing 19,241 examples generated from Spider-SS training set by appending sub-sentence method. CG-APP_T is used for *in-domain out-of-distribution* text-to-SQL evaluation.
- **CG-APP_D**: the CG-APP development set containing 3,290 examples for *cross-domain out-of-distribution* text-to-SQL evaluation.

Our evaluation is based on the exact match and execution match metric defined in the original Spider benchmark. The exact match metric measures whether the syntax tree of the predicted query without condition values is the same as that of the gold query. The execution match metric measures whether the query results from the predicted query are the same as the gold query results. All models are only trained on 7000 Spider or Spider-SS training examples. We evaluate the following open-source models that reach competitive performance on Spider:

- **RATSQL**: The RATSQL+BERT model trained on Spider (Wang et al., 2020; Devlin et al., 2019).

Dataset	Exact Match	Execution Match
Training Set	89.4%	94.0%
Development Set	90.0%	94.3%

Table 2: Use exact match and execution match metrics to evaluate the difference between the SQL in Spider and the SQL generated by NatSQL in Spider-SS.

Dataset	easy	medium	hard	extra
Spider _D	24.1%	43.1%	16.8%	16.1%
CG-SUB _T	28.3%	38.4%	20.8%	12.5%
CG-SUB _D	33.8%	37.4%	13.6%	12.6%
CG-APP _T	3.2%	30.3%	27.3%	39.1%
CG-APP _D	2.3%	41.9%	22.9%	32.8%

Table 3: The difficulty distribution of five different evaluation sets.

- **RATSQL_G**: The RATSQL+GAP model trained on Spider (Shi et al., 2021).
- **RATSQL_{GN}**: The RATSQL+GAP+NatSQL model trained on NatSQL (Gan et al., 2021b).
- **RATSQL_{GNS}**: The modified RATSQL+GAP+NatSQL model trained on Spider-SS, as discussed in Section 4.

5.2 Dataset Analysis

Spider-SS. Table 2 presents the difference between the SQL in Spider and the SQL generated by NatSQL in Spider-SS. Our evaluation results are lower than the original NatSQL dataset (Gan et al., 2021b) because the Spider-SS uses equivalent SQL and corrects some errors, as discussed in Section 2.3. Most equivalent and corrected SQL cannot get positive results in exact match metric, while a small part can not either in execution match. Therefore, the model trained on Spider-SS may not be ideal for chasing the Spider benchmark, especially based on the exact match metric. Similarly, the RATSQL_G extending NatSQL had achieved a previous SOTA result in the execution match of the Spider test set but get a worse result than the original in the exact match (Gan et al., 2021b). Thus, we recommend using NatSQL-based datasets to evaluate models trained on NatSQL.

Spider-CG. Table 3 presents the difficulty distribution of five different evaluation sets. The difficulty criteria are defined by Spider benchmark, including *easy*, *medium*, *hard* and *extra hard*. Experiments show that the more difficult the SQL is, the more difficult it is to predict correctly (Wang et al., 2020; Shi et al., 2021; Gan et al., 2021b).

It can be found from Table 3 that the difficulty distribution of CG-SUB_T and CG-SUB_D is similar to that of Spider_D. The similar distributions among CG-SUB_T, CG-SUB_D, and Spider_D support the view discussed in Section 1 that the examples generated by the substitution method are in-distribution.

On the other hand, the difficulty distributions of CG-APP_T and CG-APP_D are obviously different from that of Spider_D. Due to appending the sub-sentence, the NL and SQL in CG-APP become more complex, where the proportion of SQL in *extra hard* increased significantly, while *easy* was the opposite.

5.3 Sentence Split Algorithm Evaluation

We generate the Spider-CG based on the combination of Spider-SS sub-sentences split by the algorithm introduced in Section 2.2. We can reuse this algorithm to split the sentence in Spider-CG and then compare the splitting results with the Spider-SS sub-sentences to evaluate the stability of the splitting algorithm. We consider that a deviation of one or two words in the splitting result is acceptable. For example, in Figure 1, we consider that putting the comma of the third sub-sentence into the second sub-sentence does not change the meaning of sub-sentences, same for moving both the comma and the word ‘and’.

Table 4 presents the similarity between sub-sentences in Spider-SS and Spider-CG, which are generated by the same split algorithm under the deviation of one or two words. The similarity exceeds 90% in all evaluation set when two deviation words are allowed. Considering that the model trained on the Spider-SS does not require consistent split results, as discussed in Section 2.2, the similarity results of the splitting algorithm are good enough. The similarity of CG-SUB is higher than that of CG-APP, which means the more complex the sentence, the greater the challenge to the algorithm. Although the algorithm has been refined on the training set, the similarity between training and development in CG-SUB and CG-APP is close, showing that the algorithm performs consistently for unseen datasets. In summary, we consider that as long as the sentences are not more complex than CG-APP, the algorithm can be used stably in other text-to-SQL datasets.

Dataset	Deviation <= 1	Deviation <= 2
CG-SUB _T	93.2%	94.4%
CG-SUB _D	92.9%	94.1%
CG-APP _T	86.0%	90.4%
CG-APP _D	88.9%	92.6%

Table 4: The similarity between sub-sentences in Spider-SS and Spider-CG generated by the same split algorithm under the deviation of one or two words.

Approach	Spider _D	CG-SUB _T	CG-SUB _D	CG-APP _T	CG-APP _D
RATSQL	72.0%	79.5%	72.0%	45.1%	47.2%
RATSQL _G	72.7%	80.9%	70.3%	45.2%	44.2%
RATSQL _{GN}	73.9%	90.2%	75.0%	67.8%	60.5%
RATSQL _{GNS}	74.5%	91.4%	76.7%	82.5%	68.3%

Table 5: Exact match accuracy on evaluation sets.

5.4 Model Results

Table 5 presents the exact match accuracy on the five different evaluation sets. Specifically, the RATSQL_{GNS} consistently outperforms other models. We use the sentence split algorithm to split every sentence before inputting the RATSQL_{GNS}. Although there are some un-similar splitting results, it did not prevent the RATSQL_{GNS} from getting good performance. RATSQL_{GNS} and RATSQL_{GN} have close results in the in-distribution dataset, i.e., Spider_D, CG-SUB_T, and CG-SUB_D, but the RATSQL_{GNS} significantly improves over RATSQL_{GN} on the OOD dataset, i.e., CG-APP_T and CG-APP_D. The results demonstrate that the sub-sentence-based method can improve the generalization performance.

The evaluation results of all models in Spider_D and CG-SUB_D are close, which further confirm that the sub-sentence substitution method generates in-distribution data in a cross-domain text-to-SQL setting. In the two OOD datasets, CG-APP_T and CG-APP_D, the performance of all models has dropped by about 10% to 30%.

Although the performance of all models on Spider_D is close, the performance of RATSQL_{GN} and RATSQL_{GNS} is significantly better in the rest four datasets. One of the reasons is that RATSQL and RATSQL_G are trained on SQL while RATSQL_{GN} and RATSQL_{GNS} use NatSQL for training. As discussed in Section 5.2, since the training data of Spider and Spider-SS are about 10% different, this leads to the performance degradation in the model trained on Spider when evaluated on the SQL generated by the NatSQL of Spider-SS, and vice versa. On the other hand, ex-

Approach	Spider _D	CG-SUB _T	CG-SUB _D	CG-APP _T	CG-APP _D
RATSQL _{GN}	75.8%	86.7%	78.0%	70.4%	68.9%
RATSQL _{GNS}	76.7%	88.3%	80.4%	78.8%	75.1%

Table 6: Execution match accuracy on evaluation sets.

periments in (Gan et al., 2021b) show that NatSQL improve the model performance in *extra hard* SQL. Therefore, RATSQL_{GN} and RATSQL_{GNS} suffer less performance degradation in CG-APP_T and CG-APP_D than RATSQL_G and RATSQL.

Since RATSQL and RATSQL_G do not support generating executable SQL, we compare the execution results between RATSQL_{GN} and RATSQL_{GNS} in Table 6. RATSQL_{GNS} again consistently outperform the RATSQL_{GN}.

6 Related Work

Text-to-SQL translation. To achieve text-to-SQL translation, researchers have built various benchmarks (Iyer et al., 2017; Ana-Maria Popescu et al., 2003; Tang and Mooney, 2000; Giordani and Moschitti, 2012; Li and Jagadish, 2014; Yaghmazadeh et al., 2017; Zhong et al., 2017; Yu et al., 2018b). In particular, most recent works focus on improving the performance on Spider benchmark (Yu et al., 2018b), where models are required to generate complex SQL in cross-domain setting. Among various model architectures (Guo et al., 2019; Zhang et al., 2019; Wang et al., 2020), most of them have implemented the pre-training method, including the latest state-of-the-art model (Scholak et al., 2021; Cao et al., 2021), where Yu et al. (2021) and Rubin and Berant (2021) augment the text-to-SQL data for pre-training.

Data augmentation for text-to-SQL models. Data augmentation has been commonly used for improving performance (Xiong and Sun, 2019; Li et al., 2019). In the context of text-to-SQL generation, Yu et al. (2018a) generate synthetic training samples from some pre-defined SQL and NL question templates. Parikh et al. (2020) introduces an table-to-text dataset with over 120,000 examples that proposes a controlled generation task: given a Wikipedia table and a set of highlighted table cells, produce a one-sentence description. Yu et al. (2021) sample from the given examples and then give a large number of tables to generate new synthetic examples. Shi et al. (2021) present a model pre-training framework that jointly learns representations of NL utterances and table schemas by

leveraging generation models to generate pre-train data.

Compositional generalization for semantic parsing. Compositional generalization for semantic parsing has captured wide attention recently (Finegan-Dollak et al., 2018; Oren et al., 2020; Furrer et al., 2020; Conklin et al., 2021). Most prior works on text-to-SQL tasks focus on the cross-domain generalization, which mainly assess how the models generalize the domain knowledge to new database schemas (Suhr et al., 2020; Gan et al., 2021a). On the other hand, Shaw et al. (2021) introduces TMCD splits for studying compositional generalization in semantic parsing, where they aim to maximize the divergence of SQL compounds between the training and test sets. However, this method does not support the benchmark construction for cross-domain out-of-distribution composition generalization evaluation.

Our model is inspired by prior works on neural parsers constructed to capture granular information from a whole. Yin et al. (2021) describe a span-level supervised attention loss that improves compositional generalization in semantic parsers. Herzig and Berant (2021) propose SpanBasedSP, a parser that predicts a span tree over an input utterance, and dramatically improves performance on splits that require compositional generalization. Chen et al. (2020) propose the Neural-Symbolic Stack machine (NeSS), which integrates a symbolic stack machine into a seq2seq generation framework, and learns a neural network as the controller to operate the machine.

7 Conclusion

We introduce Spider-SS and Spider-CG for measuring compositional generalization of text-to-SQL models. Specifically, Spider-SS is a human-curated sub-sentence-based text-to-SQL dataset built upon the Spider benchmark. Spider-CG is a synthetic text-to-SQL dataset constructed by substituting and appending sub-sentences of different samples, so that the training and test sets consist of different compositions of sub-sentences. We found that the performance of previous text-to-SQL models drop dramatically on the Spider-CG OOD examples, while modifying the RATSQL+GAP+NatSQL model to fit the segmented data of Spider-SS improves compositional generalization performance.

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649

References

Ana-Maria Popescu, Oren Etzioni, and Henry Kautz. 2003. [Towards a Theory of Natural Language Interfaces to Databases](#). In *Proceedings of the 8th International Conference on Intelligent User Interfaces*, pages 149–157.

Jacob Andreas. 2020. [Good-enough compositional data augmentation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7556–7566, Online. Association for Computational Linguistics.

Ruisheng Cao, Lu Chen, Zhi Chen, Yanbin Zhao, Su Zhu, and Kai Yu. 2021. [LGESQL: Line graph enhanced text-to-SQL model with mixed local and non-local relations](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2541–2555, Online. Association for Computational Linguistics.

Xinyun Chen, Chen Liang, Adams Wei Yu, Dawn Song, and Denny Zhou. 2020. [Compositional generalization via neural-symbolic stack machines](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1690–1701. Curran Associates, Inc.

Henry Conklin, Bailin Wang, Kenny Smith, and Ivan Titov. 2021. [Meta-learning to compositionally generalize](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3322–3335, Online. Association for Computational Linguistics.

Marie-Catherine de Marnee and Christopher D. Manning. 2016. [Stanford typed dependencies manual](#).

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Catherine Finegan-Dollak, Jonathan K. Kummerfeld, Li Zhang, Karthik Ramanathan, Sesh Sadasivam, Rui Zhang, and Dragomir Radev. 2018. [Improving text-to-SQL evaluation methodology](#). pages 351–360.

Daniel Furrer, Marc van Zee, Nathan Scales, and Nathanael Schärli. 2020. [Compositional generalization in semantic parsing: Pre-training vs. specialized architectures](#). *CoRR*, abs/2007.08970.

Yujian Gan, Xinyun Chen, and Matthew Purver. 2021a. [Exploring underexplored limitations of](#)

[cross-domain text-to-sql generalization](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 650
651
652

Yujian Gan, Xinyun Chen, Jinxia Xie, Matthew Purver, John R. Woodward, John Drake, and Qiaofu Zhang. 2021b. [Natural sql: Making sql easier to infer from natural language specifications](#). 653
654
655
656

Alessandra Giordani and Alessandro Moschitti. 2012. [Automatic Generation and Reranking of SQL-derived Answers to NL Questions](#). In *Proceedings of the Second International Conference on Trustworthy Eternal Systems via Evolving Software, Data and Knowledge*, pages 59–76. 657
658
659
660
661
662

Jiaqi Guo, Zecheng Zhan, Yan Gao, Yan Xiao, Jian-Guang Lou, Ting Liu, and Dongmei Zhang. 2019. [Towards Complex Text-to-SQL in Cross-Domain Database with Intermediate Representation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4524–4535, Florence, Italy. Association for Computational Linguistics. 663
664
665
666
667
668
669
670

Jonathan Herzig and Jonathan Berant. 2021. [Span-based semantic parsing for compositional generalization](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 908–921, Online. Association for Computational Linguistics. 671
672
673
674
675
676
677
678

Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, Jayant Krishnamurthy, and Luke Zettlemoyer. 2017. [Learning a Neural Semantic Parser from User Feedback](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 963–973. 679
680
681
682
683
684

Divyansh Kaushik, Eduard Hovy, and Zachary Lipton. 2020. [Learning the difference that makes a difference with counterfactually-augmented data](#). In *International Conference on Learning Representations*. 685
686
687
688

Brenden Lake and Marco Baroni. 2018. [Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks](#). In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2873–2882. PMLR. 689
690
691
692
693
694

Fei Li and H. V. Jagadish. 2014. [Constructing an interactive natural language interface for relational databases](#). *Proceedings of the VLDB Endowment*, 8(1):73–84. 695
696
697
698

Jingjing Li, Wenlu Wang, Wei Shinn Ku, Yingtao Tian, and Haixun Wang. 2019. [SpatialNLI: A spatial domain natural language interface to databases using spatial comprehension](#). In *GIS: Proceedings of the ACM International Symposium on Advances in Geographic Information Systems*, pages 339–348, New York, NY, USA. Association for Computing Machinery. 699
700
701
702
703
704
705
706

707	Inbar Oren, Jonathan Herzig, Nitish Gupta, Matt Gardner, and Jonathan Berant. 2020. Improving compositional generalization in semantic parsing . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2482–2495, Online. Association for Computational Linguistics.	
708		
709		
710		
711		
712		
713	Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. ToTTo: A controlled table-to-text generation dataset . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1173–1186, Online. Association for Computational Linguistics.	
714		
715		
716		
717		
718		
719		
720	Ohad Rubin and Jonathan Berant. 2021. SmBoP: Semi-autoregressive bottom-up semantic parsing . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 311–324, Online. Association for Computational Linguistics.	
721		
722		
723		
724		
725		
726		
727	Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. 2020. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization .	
728		
729		
730		
731	Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. Picard: Parsing incrementally for constrained auto-regressive decoding from language models .	
732		
733		
734		
735	Peter Shaw, Ming-Wei Chang, Panupong Pasupat, and Kristina Toutanova. 2021. Compositional generalization and natural language variation: Can a semantic parsing approach handle both? In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 922–938, Online. Association for Computational Linguistics.	
736		
737		
738		
739		
740		
741		
742		
743		
744	Peng Shi, Patrick Ng, Zhiguo Wang, Henghui Zhu, Alexander Hanbo Li, Jun Wang, Cicero Nogueira dos Santos, and Bing Xiang. 2021. Learning contextual representations for semantic parsing with generation-augmented pre-training . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 35(15):13806–13814.	
745		
746		
747		
748		
749		
750		
751	Alane Suhr, Ming-Wei Chang, Peter Shaw, and Kenton Lee. 2020. Exploring unexplored generalization challenges for cross-database semantic parsing . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8372–8388, Online. Association for Computational Linguistics.	
752		
753		
754		
755		
756		
757		
758	Lappoon R Tang and Raymond J Mooney. 2000. Automated Construction of Database Interfaces: Integrating Statistical and Relational Learning for Semantic Parsing . In <i>2000 Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora</i> , pages 133–141.	
759		
760		
761		
762		
763		
	Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2020. RAT-SQL: Relation-Aware Schema Encoding and Linking for Text-to-SQL Parsers . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7567–7578, Online. Association for Computational Linguistics.	764
		765
		766
		767
		768
		769
		770
	Hongvu Xiong and Ruixiao Sun. 2019. Transferable Natural Language Interface to Structured Queries Aided by Adversarial Generation . In <i>2019 IEEE 13th International Conference on Semantic Computing (ICSC)</i> , pages 255–262. IEEE.	771
		772
		773
		774
		775
	Navid Yaghmazadeh, Yuepeng Wang, Isil Dillig, and Thomas Dillig. 2017. SQLizer: Query Synthesis from Natural Language . In <i>International Conference on Object-Oriented Programming, Systems, Languages, and Applications, ACM</i> , pages 63:1–63:26.	776
		777
		778
		779
		780
		781
	Pengcheng Yin, Hao Fang, Graham Neubig, Adam Pauls, Emmanouil Antonios Platanios, Yu Su, Sam Thomson, and Jacob Andreas. 2021. Compositional generalization for neural semantic parsing via span-level supervised attention . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2810–2823, Online. Association for Computational Linguistics.	782
		783
		784
		785
		786
		787
		788
		789
		790
	Tao Yu, Chien-Sheng Wu, Xi Victoria Lin, Bailin Wang, Yi Chern Tan, Xinyi Yang, Dragomir Radev, Richard Socher, and Caiming Xiong. 2021. Grappa: Grammar-augmented pre-training for table semantic parsing .	791
		792
		793
		794
		795
	Tao Yu, Michihiro Yasunaga, Kai Yang, Rui Zhang, Dongxu Wang, Zifan Li, and Dragomir Radev. 2018a. SyntaxSQLNet: Syntax tree networks for complex and cross-domain text-to-SQL task . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1653–1663, Brussels, Belgium. Association for Computational Linguistics.	796
		797
		798
		799
		800
		801
		802
		803
	Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018b. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-SQL task . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3911–3921, Brussels, Belgium. Association for Computational Linguistics.	804
		805
		806
		807
		808
		809
		810
		811
		812
		813
	Rui Zhang, Tao Yu, Heyang Er, Sungrok Shim, Eric Xue, Xi Victoria Lin, Tianze Shi, Caiming Xiong, Richard Socher, and Dragomir Radev. 2019. Editing-based SQL query generation for cross-domain context-dependent questions . pages 5338–5349.	814
		815
		816
		817
		818
		819
	Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2SQL: Generating Structured Queries	820
		821

822 from Natural Language using Reinforcement Learn-
823 ing. *CoRR*, abs/1709.0.