

Structured Case-based Reasoning for Inference-time Adaptation of Text-to-SQL parsers

Abhijeet Awasthi, Soumen Chakrabarti, Sunita Sarawagi

Department of Computer Science and Engineering
Indian Institute of Technology Bombay, Mumbai, India
{awasthi,soumen,sunita}@cse.iitb.ac.in

Abstract

Inference-time adaptation methods for semantic parsing are useful for leveraging examples from newly-observed domains without repeated fine-tuning. Existing approaches typically bias the decoder by simply concatenating input-output example pairs (cases) from the new domain at the encoder’s input in a Seq-to-Seq model. Such methods cannot adequately leverage the structure of logical forms in the case examples. We propose StructCBR, a structured case-based reasoning approach, which leverages subtree-level similarity between logical forms of cases and candidate outputs, resulting in better decoder decisions. For the task of adapting Text-to-SQL models to unseen schemas, we show that exploiting case examples in a structured manner via StructCBR offers consistent performance improvements over prior inference-time adaptation methods across five different databases. To the best of our knowledge, we are the first to attempt inference-time adaptation of Text-to-SQL models, and harness trainable structured similarity between subqueries.

1 Introduction

Natural language interfaces to databases (Zelle and Mooney 1996; Tang and Mooney 2000; Popescu, Etzioni, and Kautz 2003) enable access to structured information for users who are not familiar with languages like SQL by parsing user provided text-queries into executable SQLs. Text-to-SQL semantic parsing is a challenging task that not only demands robust natural language understanding but simultaneously requires reasoning over the schema structure of the databases. Databases containing similar information (e.g. census in various countries) may be designed using diverse schema structures, thus making it hard for the model to generalize across schemas unseen during training. Hence, Text-to-SQL models often struggle to parse text queries for a new schema in a zero-shot manner (Suhr et al. 2020; Lee, Polozov, and Richardson 2021; Hazoom, Malik, and Bogin 2021). In practice, a small number of Text-to-SQL examples in the target schema are often essential for successful model adaptation. However, finetuning a Text-to-SQL model for each new database is not generally practical, for the following reasons: (i) Huge variation in database schema makes it tedious to collect sufficiently large finetuning datasets for each schema, while

finetuning on small datasets is unavoidably fraught with overfitting, catastrophic forgetting, and instability w.r.t. random seeds. (ii) Finetuning may take considerable time, preventing fast incorporation of new data into the model. (iii) Often, a single large-footprint model serves multiple clients with diverse databases at the same time. Fine-tuning a separate model for each database is considered too resource-intensive in such multi-tenant scenarios.

Therefore, we focus on fast online adaptation of Text-to-SQL models without parameter updates, until the next cycle of finetuning is deemed feasible. Recently, case-based reasoning (CBR), which utilizes a memory of past labeled examples as cases, has emerged as a promising paradigm of inference-time adaptation without finetuning (Das et al. 2020, 2021; Pasupat, Zhang, and Guu 2021; Gupta et al. 2021). CBR has been found effective for tasks like knowledge graph completion (KGC) (Das et al. 2020), question answering over knowledge bases (KBQA) (Das et al. 2021), task-oriented semantic parsing (Pasupat, Zhang, and Guu 2021; Gupta et al. 2021), translation (Khandelwal et al. 2021), and text-based games (Atzeni et al. 2022). However, many prior CBR approaches designed around Seq2Seq architectures simply concatenate input-output cases with the current input at the encoder (Das et al. 2021; Pasupat, Zhang, and Guu 2021; Gupta et al. 2021). These methods do not leverage the structure of logical forms (query plan trees) in case examples.

In response, we propose StructCBR, a *structured* CBR approach that directly exploits sub-tree level similarities between the candidate outputs and the case examples for adapting a Text-to-SQL model to a new schema. We start with SmBoP (Rubin and Berant 2021), a recent semi-auto-regressive architecture that decodes query trees bottom-up, respecting the structure of SQL grammar production rules, instead of left-to-right token-level decoding in Seq2Seq models (Guo et al. 2019; Wang et al. 2020; Scholak et al. 2021; Scholak, Schucher, and Bahdanau 2021). We implement a novel *structured case memory lookup* module to boost scores of promising candidate trees using sub-tree level similarity with case trees under similar input context. This similarity is trainable. We show that explicitly-learned structured memory lookup leads to more accurate transfer from cases, compared to prior inference-time adaptation methods such as ConcatCBR and

GTM (Khandelwal et al. 2020; Das et al. 2021; Khandelwal et al. 2021) that we implemented both on SmBoP, and other Seq2Seq Text-to-SQL architectures like T5-large.

We summarize our contributions as follows:

- 1) We propose StructCBR, which, to our knowledge, is the first inference-time adaptation method for Text-to-SQL parsing without parameter fine-tuning.
- 2) StructCBR incorporates a novel structured case memory and trainable query subtree similarity module that can boost scores of likely-correct outputs during inference. This is in contrast with earlier approaches like ConcatCBR and GTM.
- 3) We propose a trainable compositional sub-tree similarity function that is both more accurate and more efficient for scoring large search frontiers, compared to default whole-tree embeddings.
- 4) Through experiments with five database schemas (§ 5) of varying complexity, we observe that StructCBR is consistently better than prior inference-time adaptation methods on both SmBoP and sequence-based Text-to-SQL models.
- 5) We show that StructCBR provides almost instant adaptation to a target schema. In contrast, finetuning (§ 5) can be up to 500 times slower.

2 SmBoP preliminaries

We present a brief background on SmBoP here. Readers familiar with SmBoP can largely skip this section. SmBoP converts a natural language question $\bar{x} \in \mathcal{X}$ (called the ‘utterance’) targeting a database schema $\bar{s} \in \mathcal{S}$, to an SQL query $\hat{q} \in \mathcal{Q}$. We describe the major modules in SmBoP.

Utterance and schema encoding: Given token sequence $\bar{x} = [x_1, x_2, \dots, x_n]$ in the text query, and database schema $\bar{s} = [s_1, s_2, \dots, s_m]$ denoting table and column names, SmBoP jointly encodes them using a pre-trained Transformer like RoBERTa (Liu et al. 2020) followed by relation-aware-transformer (RAT) layers (Shaw, Uszkoreit, and Vaswani 2018; Wang et al. 2020; Scholak et al. 2021). We denote the output from the last encoder layer as $\bar{\mathbf{x}} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ and $\bar{\mathbf{s}} = [\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_m]$, representing the jointly encoded contextual embeddings of text tokens and schema elements respectively.

Decoding SQL output: Unlike standard sequence-based decoding (Wang et al. 2020; Scholak et al. 2021; Scholak, Schucher, and Bahdanau 2021), SmBoP decodes the SQL tree bottom-up and in layers. SmBoP views any SQL query as a height-balanced relational algebra tree converted using a special idempotent KEEP operator κ as shown in Figure 1. Given a beam size K , at decoding step t , the decoder beam B_t comprises K candidate sub-trees of height t from the bottom. At step $t + 1$, trees from B_t are grown either via unary operators (e.g. COUNT), or by combining two trees in B_t using a binary operator (e.g. >), as per the SQL grammar. The candidate trees at step $t + 1$ form a frontier set F_{t+1} and is of size $|F_{t+1}| = K^2|\mathcal{B}| + K|\mathcal{U}|$, where \mathcal{B} and \mathcal{U} represent the set of binary and unary operations respectively. SmBoP assigns each candidate tree $z \in F_{t+1}$ a score $s_\theta(z)$

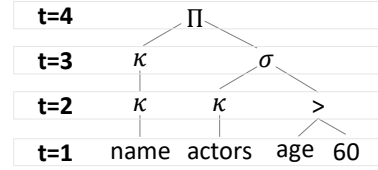


Figure 1: SmBoP (Rubin and Berant 2021) decodes SQL as a balanced relational algebra tree. At each level t , trees in the beam combine via unary or binary operators to form candidates of the next beam. StructCBR leverages CBR on generated sub-trees.

Text 1	<i>Give the code of the airport with the fewest number of flights</i>
SmBoP output	SELECT sourceairport FROM flights GROUP BY sourceairport ORDER BY SUM(flightno) ASC LIMIT 1
Correct SQL	SELECT airportcode FROM airports JOIN flights ON airportcode = sourceairport GROUP BY sourceairport ORDER BY COUNT(*) ASC LIMIT 1
Text 2	<i>What is the code of the airport that has the highest number of flights?</i>

Table 1: Illustration of the lack of generalization of Text-to-SQL to new schema.

(described below). The top- K highest scoring trees in F_{t+1} form the next beam B_{t+1} . This continues up to a maximum height T , when the highest scoring tree in B_T is output.

Scoring a tree: A tree $z = (z_b, z_\ell, z_r)$ consists of root operator z_b and subtrees z_ℓ, z_r . SmBoP encodes a variable-size tree z into two fixed dimensional vectors: (i) \mathbf{z} : an embedding of the tree computed recursively on the tree structure, where a transformer outputs $\mathbf{z} = \text{TX}_\theta([z_b, z_\ell, z_r])$; (ii) \mathbf{z}' : a contextual representation of z grounded in input text \bar{x} computed via a multiheaded cross attention module $\mathbf{z}' = \text{XAtt}_\theta(\mathbf{z}, \bar{\mathbf{x}})$. SmBoP computes the score of a tree $z \in F_{t+1}$, as follows:

$$s_\theta(z) = \mathbf{w}_{z_b}^T \text{FF}_\theta([\mathbf{z}_\ell; \mathbf{z}'_\ell; \mathbf{z}_r; \mathbf{z}'_r]) \quad (1)$$

where FF_θ is a feed forward network, and \mathbf{w}_{z_b} represents a learned embedding of operator z_b .

The model is trained using Text-SQL pairs from a set of training schema to maximize the likelihood of the correct subtrees at each beam. During inference, when presented with text utterances relating to a new database schema, the model often fails to discover the mapping of the text to schema names and relationships in the new schema. Table 1 presents an example where a SmBoP model trained on the Spider dataset (Yu et al. 2018) is deployed on a new schema about flights. On inspecting the predicted and correct SQL, we find that the model failed to reason that *number of flights* requires a `count(*)` instead of `sum(flightno)`. Now suppose an expert provides the correct SQL as additional information to be used during inference of subsequent queries. Consider a second query (shown as Text 2 in Table 1) that

also needs to reason about *number of flights*, and the default SmBoP makes similar errors (not shown). Only existing mechanism in SmBoP is to fine-tune parameters which could be time-consuming and unstable. In the next section we show how our method can instantaneously leverage test-time user labels to predict the correct SQL for Text 2. More such anecdotes appear in Table A4 of the Appendix.

3 Our proposed method: StructCBR

We aim to learn a Text-to-SQL model M , using a dataset $\mathcal{D}_{\text{train}}$ of Text-SQL pairs such that it is capable of **C1**: Independently translating the text queries \bar{x} to executable SQL programs \hat{q} , and **C2**: Utilizing a small set \mathcal{D}_{new} of Text-SQL pairs from a target schema s_{new} , to improve its own predictions during inference, without finetuning. In line with prior work (Das et al. 2020, 2021), we refer to the second capability **C2** as Case-based reasoning (CBR), and the dataset \mathcal{D}_{new} of Text-SQL pairs in the target schema as cases.

The StructCBR module leverages the similarity between gold subtrees that appear in similar contexts in the set of cases \mathcal{D}_{new} and the candidate subtrees in SmBoP’s frontier F_{t+1} , to boost the scores of likely-correct candidates at each decoding step $t + 1$. Consider a subtree z in the frontier F_{t+1} for an input text \bar{x} , a case-example with text question as \bar{x}^c , and the gold SQL tree as $\mathcal{Z}_{\text{gold}}^c$. Let z^c be a subtree of $\mathcal{Z}_{\text{gold}}^c$. The key idea of StructCBR is, if z and z^c are structurally similar, and appear in similar contexts w.r.t. \bar{x} and \bar{x}^c , then there is a strong evidence that the subtree z should also appear as a part of the gold tree $\mathcal{Z}_{\text{gold}}$ of \bar{x} . Figure 2 provides an illustration with $z = \text{age} > 60$ in the candidate frontier F_{t+1} , and a similarly structured case tree $z^c = \text{age} > 80$ appearing in a similar context \bar{x}^c (both contain the phrase *who are past*).

Even though the key idea of matching with case sub-trees is simple, several important design choices had to be made to ensure that CBR inter-operates efficiently with SmBoP’s own scoring, and consistently improves its performance in the presence of multiple cases of varying levels of relatedness. First, how should we compute the *contextual* similarity of a candidate tree z with a case tree, given that memory would also contain unrelated cases that would match wrong candidate trees? Second, how can we efficiently compute the similarity of all candidate trees with all entries in the case memory? Unlike Seq2Seq models that do not perform beam-search during training, SmBoP generates a large search frontier even during training. We elaborate on how our design tackles these challenges next.

Algorithm 1 presents the high-level pseudo code, with the text in blue font representing the StructCBR additions to the SmBoP model.

3.1 Choosing tree representations

We need to choose a representation of a tree z using which we can efficiently compute similarity with case trees. Just the structural similarity of z with a case z^c is not sufficient unless we also contextualize them on their respective inputs. Accordingly, we design an embedding function $G_\phi(z, \bar{x}) \mapsto \mathbb{R}^d$ that jointly encodes a candidate tree z corresponding to

Algorithm 1: SmBoP with StructCBR.

```

1 input:  $\bar{x}, \bar{s}, \mathcal{D}_{\text{new}}$ 
2  $M \leftarrow \text{CreateCaseMemory}(\mathcal{D}_{\text{new}})$  (§ 3.2)
3  $\bar{x}, \bar{s} \leftarrow \text{EncodeTextSchema}_\theta(\bar{x}, \bar{s})$ 
4  $B_0 \leftarrow \text{top-}K \text{ schema constants and DB values}$ 
5 for  $t \leftarrow 0 \dots T - 1$  do
6    $\mathbf{z} \leftarrow \text{CreateTreeReps}(z)$ 
7    $\mathbf{z}' \leftarrow \text{GroundTreeReps}(\mathbf{z}, \mathbf{x})$ 
8    $p_\theta \leftarrow \text{SmBoPScores}(\mathbf{z}, \mathbf{z}')$  (§ 2)
9    $G_\phi(z, \bar{x}) \leftarrow \text{JointReps}(\mathbf{z}, \mathbf{z}', \mathbf{x})$  (Eqn 2)
10   $\text{sim}_\phi \leftarrow \text{TreeSim}(G_\phi(z, \bar{x}), M)$  (Eqn 4)
11   $p_\phi \leftarrow \text{StructCBRScores}(\text{sim}_\phi, M)$  (Eqn 5)
12   $F_{t+1} \leftarrow \text{CombineScores}(p_\theta, p_\phi)$  (Eqn 6)
13   $B_{t+1} \leftarrow \text{top-}K(F_{t+1})$ 
14 return  $\text{argmax}_z(B_T)$ 

```

an input \bar{x} as a d dimensional vector. We train a separate transformer model TX_ϕ with parameters ϕ that takes as input four vectors: \mathbf{z} that encodes the structure of the tree z , \mathbf{z}' that is the contextual representation of z defined in § 2, an embedding \mathbf{w}_b of z ’s root node b , and $\text{pool}(\mathbf{x})$ a mean-pooled version of the input text representation \bar{x} :

$$G_\phi(z, \bar{x}) = \text{TX}_\phi([\mathbf{z}, \mathbf{z}', \mathbf{w}_b, \text{pool}(\mathbf{x})]). \quad (2)$$

This embedding captures both the structure and context and the parameters ϕ are trained to co-embed similar trees in matching contexts, while pulling apart pairs differing either structurally or contextually. For example, in Figure 2 if the query text was *Name all actors who are 60 or above*, then the similarity of candidate *age > 60* from the same case sub-tree should be reduced. Unlike recursive tree representations (Socher et al. 2013), here contextualization w.r.t. \bar{x} plays a critical role.

3.2 Case memory design

We construct a case memory M over the gold SQL trees $\{\mathcal{Z}_{\text{gold}}^c\}$ for *all* cases in \mathcal{D}_{new} . Corresponding to each node b of a gold tree $\mathcal{Z}_{\text{gold}}^c$ we form a subtree rooted at b and including the part of $\mathcal{Z}_{\text{gold}}^c$ below b . Thus, the size of the case memory is the total number of nodes over all gold trees in cases. The encoding $G_\phi(z^c, \bar{x}^c)$ of each subtree z^c for a case $(\bar{x}^c, \mathcal{Z}_{\text{gold}}^c)$ in \mathcal{D}_{new} is pre-computed using Equation 2 and stored in M .

3.3 Efficient tree similarity computation

We need to compute the similarity of each tree z in the frontier F_{t+1} with all case sub-trees $z^c \in M$. One way to compute similarity between trees z and z^c is based on ℓ_2 distance¹ between their G_ϕ representations as follows:

$$\text{sim}_\phi(z, z^c, \bar{x}, \bar{x}^c) = -\|G_\phi(z, \bar{x}) - G_\phi(z^c, \bar{x}^c)\|_2 \quad (3)$$

However, computing G_ϕ representations for each tree $z \in F_{t+1}$ entails large memory and compute costs since the frontier size $|F_{t+1}| = K^2|\mathcal{B}| + K|\mathcal{U}|$ is quadratic in beam-size

¹Like Khandelwal et al. (2020) we observed better results with ℓ_2 distance, in comparison to inner product.

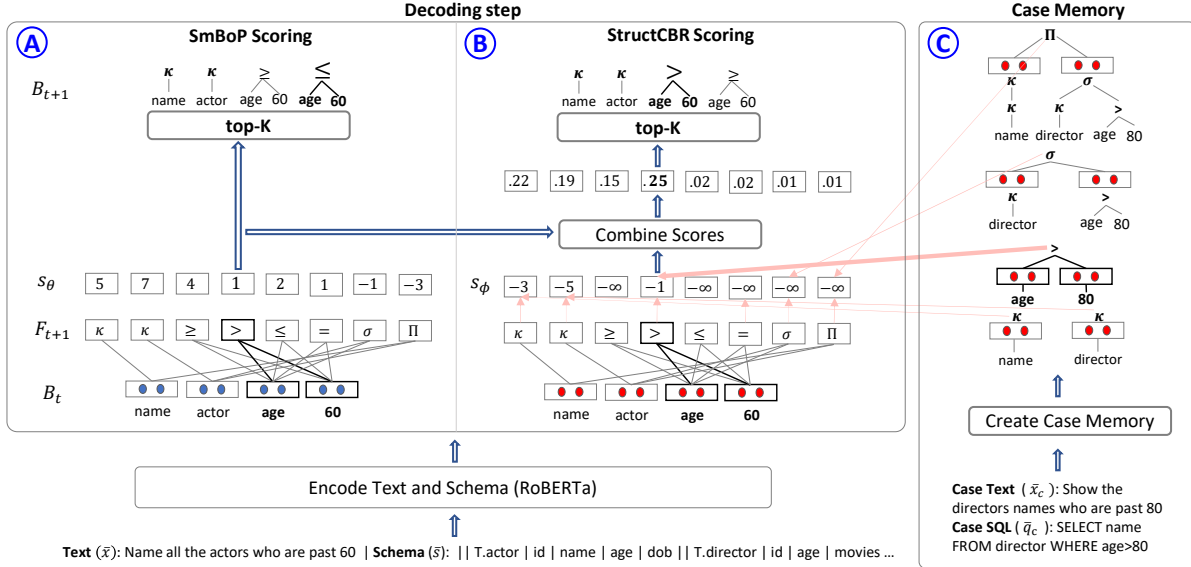


Figure 2: **Augmenting SmBoP with StructCBR (Structured Case-based Reasoning)**: In part (A), the top-K step in SmBoP scoring misses the correct sub-tree $\text{age} > 60$ due to a lower score (score=1) w.r.t. competing sub-trees in the frontier F_{t+1} like $\text{age} \geq 60$ (score=4) and $\text{age} \leq 60$ (score=2). In part (C), StructCBR creates a memory of all the sub-tree representations available in cases as described in § 3.2. In part (B), StructCBR scores the frontier candidates based on learned tree-similarities w.r.t. the sub-trees in cases as described in § 3.3 and § 3.4. For example, StructCBR boosts the score of $\text{age} > 60$ because of its high similarity with the case sub-tree $\text{age} > 80$ and similarity of context who are past. Thus, the top-K step applied on the combined SmBoP and StructCBR scores recovers the correct sub-trees that otherwise may get missed based on SmBoP’s scoring alone. For brevity, we consider only one case-example in this figure.

K . With the default K for SmBoP being 30, and size of the SmBoP grammar, this translates to around 23 thousand trees per frontier. Pruning the frontier F_{t+1} based on SmBoP scores alone resulted in poor performance. This led us to design an alternative compositional CBR scoring method that can more efficiently score all candidate trees in the frontier.

Our key idea is to compute the similarity between two trees compositionally as a function of similarity between their left and right children respectively. This requires only $\mathcal{O}(K)$ tree representations for the trees in beam B_t as against $K^2|B| + K|U|$ operations of the whole-tree approach. Recall that the trees $z \in F_{t+1}$ are formed by combining trees in beam B_t via SQL operations. A tree $z \in F_{t+1}$ can thus be represented as $z = (z_b, z_\ell, z_r)$ where $z_\ell, z_r \in B_t$ denote the left and right child respectively, and z_b is the root node combining z_ℓ and z_r . After the beam B_t is created, we compute the embedding G_ϕ for each tree in B_t using Equation (2). Now, the similarity between a candidate tree $z = (z_b, z_\ell, z_r)$ for an input text \bar{x} , and a case sub-tree $z^c = (z_b^c, z_\ell^c, z_r^c)$ on input text \bar{x}^c in memory M is computed as:

$$\widehat{\text{sim}}_\phi(z, z^c, \bar{x}, \bar{x}^c) = \text{sim}_\phi(z_\ell, z_\ell^c, \bar{x}, \bar{x}^c) + \text{sim}_\phi(z_r, z_r^c, \bar{x}, \bar{x}^c). \quad (4)$$

In Section 5, we show that using this similarity function provides better results by allowing the entire frontier to be scored more efficiently in comparison to computing similarities based on Equation 3 only for a subset of trees in the frontier pruned based on SmBoP scores.

3.4 Boosting SmBoP frontier with tree similarities

To compute an overall score of a candidate tree $z \in F_{t+1}$ based on its similarity with the case sub-trees in M , we aggregate over all the case sub-trees z^c with the same root node ($z_b = z_b^c$) using a logsumexp operator, which provides us a soft-maxed similarity of z w.r.t. case sub-trees.

$$s_\phi(z) = \log \sum_{c \in M \wedge z_b = z_b^c} \exp(\widehat{\text{sim}}_\phi(z, z^c, \bar{x}, \bar{x}^c)) \quad (5)$$

Now every candidate tree $z \in F_{t+1}$ has two scores: $s_\theta(z)$ assigned by default SmBoP and $s_\phi(z)$ computed by StructCBR. The scores $s_\theta(z)$ and $s_\phi(z)$ can lie in very different ranges. Summing them in a normalized probability space provided better results than summing the scores directly. Hence, we independently normalize $s_\theta(z)$ to $p_\theta(z)$ and $s_\phi(z)$ to $p_\phi(z)$ by a softmax operation applied over all trees in the frontier. The combined score of a frontier tree z is:

$$p(z) = (p_\theta(z) + p_\phi(z))/2. \quad (6)$$

3.5 Supervising StructCBR

During training, we assume availability of training data $\mathcal{D}_{\text{train}} = \{(\bar{x}_i, \bar{s}_i, \bar{q}_i)\}_{i=1}^N$ consisting of utterances \bar{x}_i on a schema \bar{s}_i , and the corresponding gold SQL queries \bar{q}_i . We first train the SmBoP model, parameterized as θ , using $\mathcal{D}_{\text{train}}$. The training objective of SmBoP for a single example maximizes the likelihood of sub-trees that are part of the tree $\mathcal{Z}_{\text{gold}}$

corresponding to gold SQL \bar{q} :

$$\mathcal{L}_\theta = - \sum_{t=0}^T \sum_{z_t \in \mathcal{Z}_{\text{gold}}} \log p_\theta(z_t). \quad (7)$$

Next, we introduce the StructCBR module parameterized as ϕ on top of the (now frozen) SmBoP model. We observed training the StructCBR parameters ϕ while freezing the learned SmBoP parameters θ to provide slightly better results in comparison to training both θ and ϕ jointly. The parameters ϕ are also learned using $\mathcal{D}_{\text{train}}$ by maximizing the likelihood of the gold subtrees as per the distributions p_ϕ and p through the following loss function:

$$\mathcal{L}_\phi = - \sum_{t=0}^T \sum_{z_t \in \mathcal{Z}_{\text{gold}}} \log p_\phi(z_t) + \log p(z_t) \quad (8)$$

The $-\log p_\phi(z_t)$ term maximizes the likelihood of gold trees w.r.t. the CBR distribution p_ϕ , independent of the SmBoP distribution p_θ . Similarly, the $-\log p(z_t)$ term maximizes the likelihood of the gold trees w.r.t. the combined distribution p (Eqn 6). During training, we design each training batch to contain C examples from same schema so that for a given train example, the remaining $C - 1$ examples serve as the cases from the same schema. We train with $C = 32$ and a batch-size of 64.

4 Related work

We review prior work on inference-time model adaptation for related tasks and also describe our adaptation of some of these works in the context of Text-to-SQL for comparisons with StructCBR.

Concatenating related examples with input: A common approach, that we call ConcatCBR, for utilizing cases during inference is to concatenate the input-output pair of each case along with the input text at the encoder of a Seq2Seq model. During training, the decoder is expected to learn to utilize the cases on the encoder side. Das et al. (2021) utilize ConcatCBR for question answering over knowledge bases, and Pasupat, Zhang, and Guu (2021); Gupta et al. (2021) utilize ConcatCBR for other semantic parsing tasks. ConcatCBR is similar to the retrieve and edit framework for structured outputs (Hashimoto et al. 2018) and machine translation (Hossain, Ghazvininejad, and Zettlemoyer 2020). For the Text-to-SQL task, we implement a ConcatCBR baseline that trains an SmBoP model to use retrieved Text-SQL examples concatenated with the input-text. During inference, the retrieval index is updated with the case-examples from the target schema.

Generalization through Memorization (GTM): Khandelwal et al. (2020, 2021) propose a memory look-up based method for adapting pre-trained language and machine translation models to a target domain. Given a target dataset, their method constructs a look-up index by using contextual embeddings from the pre-trained model as keys and the corresponding text tokens as values. During inference the model scores are interpolated with the similarity scores aggregated

over the nearest neighbours in the look-up index. For our Text-to-SQL set-up, we implement this baseline using a trained SmBoP model. We memorize the dataset \mathcal{D}_{new} in the target schema by creating a look-up index with embeddings of child subtrees from SmBoP as keys: $[\mathbf{z}_\ell; \mathbf{z}'_\ell; \mathbf{z}_r; \mathbf{z}'_r]$, and their parent nodes as values. During inference, the scores from the SmBoP model are interpolated with neighbour similarities in a way similar to Khandelwal et al. (2021). Unlike StructCBR and ConcatCBR, this baseline (GTM) does not explicitly train the SmBoP model for utilizing the cases during inference.

We discuss other related work in Appendix A.6.

5 Experiments

We evaluate StructCBR for adapting a Text-to-SQL model to five different target schemas without finetuning. The target schemas are chosen from varying domains. We compare StructCBR with prior inference-time adaptation methods discussed in § 4, and present an ablation study. We also show that StructCBR enables much faster adaptation of Text-to-SQL models in comparison to finetuning.

Datasets: We utilize Spider (Yu et al. 2018), which is a collection of Text-to-SQL examples covering 200 unique schemas. We use the train split of Spider as $\mathcal{D}_{\text{train}}$, for training all the models. $\mathcal{D}_{\text{train}}$ contains 7000 Text-SQL example pairs from 140 databases. For evaluation, we hold out the following five databases containing the most examples from Spider’s dev set ²: {world_1, car_1, cre_Doc_Template_Mgt, dog_kennels, flight_2}. The five evaluation databases do not occur in the train set, and belong to sufficiently different domains of varying difficulty. The remaining part of the dev set containing 576 examples is used for model selection while training on $\mathcal{D}_{\text{train}}$. We hold out 30 randomly selected examples from each of the five selected databases as \mathcal{D}_{new} (cases) for adaptation, and use the remaining examples as the test set, $\mathcal{D}_{\text{test}}$. The average size of $\mathcal{D}_{\text{test}}$ is 60, and varies from roughly 50 to 90 examples across the five schemas. To ensure robust evaluation, we report numbers averaged over three random $\mathcal{D}_{\text{new}}/\mathcal{D}_{\text{test}}$ splits. We also report the numbers micro-averaged over all the 300 test examples across the five schemas.

Evaluation metrics: Following prior work (Yu et al. 2018), we report Execution Accuracy (**EX**) and Exact-Set-Match Accuracy (**EM**) for all the methods. EX returns 1 if executing the gold query \bar{q} and the predicted query \hat{q} on the target database gives the same results. EM compares all the SQL clauses within \bar{q} and \hat{q} and returns 1 if all the clauses match, except possibly the DB-values (constants) in the SQL query. Most Text-to-SQL models utilize beam search, and return the top- K highest scoring candidates in the beam as the output. Hence, we also report the top- K versions of EM and EX metrics as BEM and BEX respectively, where K is the beam size. In our experiments, $K = 30$. BEM/BEX for a beam is 1, if at least one of the candidates in the beam has an EM/EX of 1.

²Spider’s test set is publicly inaccessible as of 08/15/2022.

Methods compared: We compare the accuracy of StructCBR after adaptation with the following methods: **(i)** SmBoP: The base model without any adaptation to benchmark the gains from different inference-time adaptation methods. **(ii)** ConcatCBR: The standard method of concatenating input-output case examples with the input-text. **(iii)** GTM: Mapping \mathcal{D}_{new} using SmBoP into a non-parametric memory for augmenting model’s predictions with inference-time memory look-ups similar to Khandelwal et al. (2020, 2021). We discussed ConcatCBR and GTM baselines in Section 4. All the baselines are implemented using SmBoP as the base model. In Appendix A.2 we also present ConcatCBR implemented on a T5-based Seq2Seq model.

Implementation details: We implemented StructCBR and baselines using AllenNLP (Gardner et al. 2017) and Transformers (Wolf et al. 2020) libraries. We utilize the authors’ implementation of SmBoP (Rubin and Berant 2021). Due to limited computing resources, we primarily experiment with the ROBERTA-BASE checkpoint for initializing the text encoder, followed by four RAT layers (Wang et al. 2020) to encode the schema structure. All other hyper-parameters are the set to their default values. The SmBoP model is trained on $\mathcal{D}_{\text{train}}$ for 60K steps with a batch size of 80, using the default learning rate (LR) of 1.86×10^{-4} . The GTM baseline utilizes the output of this model for memory look-ups. For ConcatCBR baseline we train the SmBoP model further for 60K steps with a LR of 5×10^{-5} , while concatenating the retrieved cases in the encoder’s input. StructCBR introduces 2.53% additional parameters (ϕ) over the SmBoP parameters (θ). We train the parameters ϕ on $\mathcal{D}_{\text{train}}$ using a batch size of 64 for 60K steps with the default LR of 1.86×10^{-4} . Additional training details are provided in Appendix A.4.

Overall Results: In Table 2, we compare StructCBR with different inference-time methods for adapting SmBoP based models on five different evaluation schemas. The performance of the unadapted SmBoP model varies from 46.1 EM to 84.3 EM across the five schemas indicating that the evaluation schemas are of varying difficulty. We find that StructCBR almost consistently offers substantial gains over the base SmBoP model w.r.t. all the metrics, and across all the five schemas. StructCBR gains upto 6.3 EM points and on average 4.6 EM points over the base SmBoP model, while achieving almost 4 times higher gains than best existing method. In contrast, the ConcatCBR method, which has been shown to work well for other semantic parsing tasks (Das et al. 2021; Pasupat, Zhang, and Guu 2021), provides positive EM gains for only three out of five schemas, and causes overall drop in micro-averaged EM over all the test instances. We also explored a ConcatCBR implementation based on T5-large model with a comparable base accuracy (Appendix A.2). Even compared to this method, we continue to observe higher and more consistent gains from StructCBR on SmBoP. GTM, another inference-time adaptation baseline, utilizes memory lookups similar to Khandelwal et al. (2021), and offers almost consistent but much smaller gains in comparison to StructCBR. The GTM baseline performs memory look-ups based on representations learned by the base SmBoP model whereas StructCBR is explicitly trained to perform sub-tree

	EM	EX	BEM	BEX
$s_{\text{new}} = \text{world_1} , \mathcal{D}_{\text{test}} = 89$				
SmBOP	46.1	36.3	55.4	49.0
ConcatCBR	-4.1	-4.9	-5.2	-6.3
GTM	+0.0	+0.0	+0.4	-0.4
StructCBR (Ours)	+2.6	+2.6	+3.8	+0.0
$s_{\text{new}} = \text{car_1} , \mathcal{D}_{\text{test}} = 60$				
SmBOP	43.3	45.6	50.6	53.9
ConcatCBR	+3.9	+3.9	+4.4	+3.4
GTM	+2.2	+2.2	+3.9	+2.8
StructCBR (Ours)	+6.1	+6.7	+8.3	+7.2
$s_{\text{new}} = \text{cre_Doc_Template_Mgt} , \mathcal{D}_{\text{test}} = 53$				
SmBOP	84.3	89.3	94.3	96.2
ConcatCBR	+5.7	+5.0	+1.9	+1.3
GTM	+3.2	+1.9	+1.9	+1.3
StructCBR (Ours)	+6.3	+3.8	+3.8	+2.5
$s_{\text{new}} = \text{dog_kennels} , \mathcal{D}_{\text{test}} = 49$				
SmBOP	66.6	59.2	80.3	74.8
ConcatCBR	+0.7	+0.0	-2.0	-1.3
GTM	+1.4	+2.7	+0.7	+3.4
StructCBR (Ours)	+3.4	+6.1	+3.4	+4.1
$s_{\text{new}} = \text{flight_2} , \mathcal{D}_{\text{test}} = 49$				
SmBOP	55.8	67.4	59.2	80.3
ConcatCBR	-8.8	-6.1	-2.0	-1.3
GTM	+0.0	-1.4	+2.1	+0.0
StructCBR (Ours)	+5.5	+4.1	+12.3	+4.8
Micro-Average , $ \mathcal{D}_{\text{test}} = 300$				
SmBOP	57.2	56.3	66.0	67.6
ConcatCBR	-0.8	-0.8	-1.0	-1.4
GTM	+1.2	+1.0	+1.7	+1.2
StructCBR (Ours)	+4.6	+4.4	+6.0	+3.4

Table 2: Comparison of StructCBR with prior inference time adaptation methods (ConcatCBR and GTM) on 5 different schemas. SmBoP row provides the performance of the unadapted model, while other rows report gains over SmBoP after adaptation. Micro-average refers to numbers averaged over all the test instances spanning across the five schemas. $|\mathcal{D}_{\text{test}}|$ refers to size of the test-set, and s_{new} provides the schema name.

level memory look-ups. In particular, StructCBR boosts the top- K EM score (BEM) by up to 12.3 points. With a large-sized SmBoP architecture we continue to observe consistent gains for most of the schemas. Section A.5 and Table A3 in Appendix provide results for adapting the large-sized models in the same setting as Table 2. Table A4 in Appendix provides some anecdotal examples that were originally mispredicted by the base SmBoP model, but were fixed by StructCBR.

Impact of number of cases: We run another set of experiments by reducing the number of cases from 30 to 20. Figure 3 shows that gains from both GTM and StructCBR are smaller with fewer cases, as expected. Interestingly, Struct-

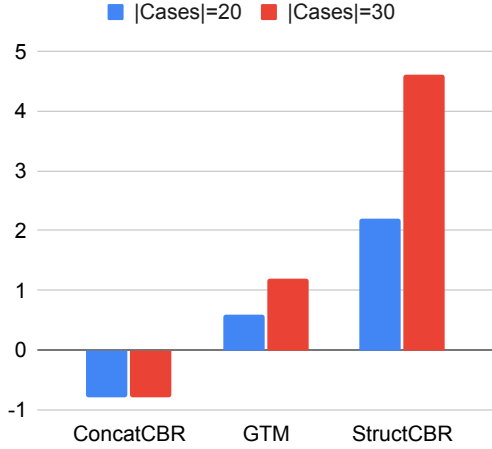


Figure 3: Impact of case size on gains in EM (y-axis) of adapted models (x-axis). With 20 cases, StructCBR still outperforms GTM with 30 case examples.

CBR with 20 cases still outperforms GTM with 30 case examples. For ConcatCBR we do not observe significant changes because it additionally augments the case memory with examples retrieved from Spider’s train set. Not augmenting cases via retrieval resulted in even worse performance of ConcatCBR.

Justification for tree similarity: In Section 3.3, we argued that directly computing tree similarities (sim_ϕ) using Equation (3) was inefficient, and required pruning to be practical. Instead in StructCBR we compute tree similarity ($\widehat{\text{sim}}_\phi$ as per Equation 4) more efficiently as a composition of similarity of its children, and does not require pruning. Table 3 justifies our design by comparing results of scoring a *pruned frontier* containing top-5K trees using sim_ϕ , with scoring the *entire frontier* using $\widehat{\text{sim}}_\phi$. Efficiently scoring the entire frontier provides us better results on 4 out 5 schemas and a micro-averaged gain of 2.2 points in EM.

	world_1	car_1	cre_Doc	dog_kenn	flights_2	Micro avg
Whole-tree	43.4	48.9	88.7	76.2	55.1	59.6
Ours	48.7	49.4	90.6	70.1	61.2	61.8

Table 3: Scoring *pruned frontier* as per similarity of whole trees as in Equation (3) vs. our scoring of the *entire frontier* as per similarity composed from subtrees, as in Equation (4). Better performance on four out of five schemas justifies our design of tree similarity.

Comparison with fine-tuning: Adapting transformer models via finetuning often provides significant accuracy improvements, thanks to recent advances in optimization of transformer models (Huang et al. 2020; Xu et al. 2021). However, finetuning is not viable when a model needs to be adapted instantaneously from just a few data points, e.g.

Adaptation Method	Time(s)↓	EM%↑
SmBOP (Unadapted)	0.0	48.3
StructCBR	0.1	50.6
Finetuning (1 epochs)	5.0	48.3
Finetuning (2 epochs)	10.0	47.5
Finetuning (5 epochs)	25.0	48.3
Finetuning (10 epochs)	50.0	50.2
Finetuning (20 epochs)	100.0	50.9
Finetuning (100 epochs)	500.0	52.1

Table 4: Comparing adaptation time and EM accuracy of StructCBR and finetuning for different number of epochs on 30 cases of world_1 schema. We report wall clock times in seconds. All the numbers were averaged over 3 runs. Finetuning took atleast 500x more time (10 epochs) to achieve EM gains that are almost instantly (0.1s) achieved by StructCBR

quickly incorporating expert feedback in the form of a few examples. In Table 4, we show that StructCBR serves the purpose of instantaneously improving model performance (+2.6 pts EM), while being roughly 50× faster than finetuning for a single epoch, and 5000× faster than finetuning for 100 epochs of 30 case examples from world_1 schema. Finetuning required atleast 10 epochs to outperform StructCBR’s accuracy. Each epoch involved four parameter update steps of batch-size 8. We note that applying StructCBR over SmBoP incurs a small increase ($\sim 1.2\times$) in inference time per example. Overall, we observe that on three out of five schemas StructCBR instantly offers more than 50% of gains achieved by finetuning for 100 epochs, and 43% of gains on average across all schemas (Table A2 in appendix). This establishes StructCBR as a method for instantly utilizing available case examples for fast model adaptation until the next cycle of finetuning becomes affordable.

6 Conclusion and Future Work

We presented StructCBR, a method for instant adaptation of Text-to-SQL models without finetuning. We show that utilizing case examples in a more structured way via sub-tree level look-ups offers better performance in comparison to the standard method of concatenating case examples with input text into a Seq2Seq encoder. We find that explicitly learning to perform memory look-ups provides larger gains in comparison to look-ups using a pre-trained model. Finally, we show that StructCBR enables much faster model adaptation in comparison to finetuning, potentially allowing instantaneous adaptation to expert feedback provided in form of a few examples. We propose StructCBR as a faster alternative to model adaptation, until the next finetuning cycle is deemed feasible. Despite its speed, there remains an accuracy gap between StructCBR and sufficient finetuning, which might be narrowed by more sophisticated similarity networks. Our exploration of StructCBR focused only on the Text-to-SQL task. In future we wish to explore the effectiveness of StructCBR for other semantic parsing tasks.

References

- Atzeni, M.; Dhuliawala, S. Z.; Murugesan, K.; and SACHAN, M. 2022. Case-based Reasoning for Better Generalization in Text-Adventure Games. In *International Conference on Learning Representations*.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.
- Chen, M.; Tworek, J.; Jun, H.; Yuan, Q.; de Oliveira Pinto, H. P.; Kaplan, J.; Edwards, H.; Burda, Y.; Joseph, N.; Brockman, G.; Ray, A.; Puri, R.; Krueger, G.; Petrov, M.; Khlaaf, H.; Sastry, G.; Mishkin, P.; Chan, B.; Gray, S.; Ryder, N.; Pavlov, M.; Power, A.; Kaiser, L.; Bavarian, M.; Winter, C.; Tillet, P.; Such, F. P.; Cummings, D.; Plappert, M.; Chantzis, F.; Barnes, E.; Herbert-Voss, A.; Guss, W. H.; Nichol, A.; Paino, A.; Tezak, N.; Tang, J.; Babuschkin, I.; Balaji, S.; Jain, S.; Saunders, W.; Hesse, C.; Carr, A. N.; Leike, J.; Achiam, J.; Misra, V.; Morikawa, E.; Radford, A.; Knight, M.; Brundage, M.; Murati, M.; Mayer, K.; Welinder, P.; McGrew, B.; Amodei, D.; McCandlish, S.; Sutskever, I.; and Zaremba, W. 2021. Evaluating Large Language Models Trained on Code.
- Das, R.; Godbole, A.; Monath, N.; Zaheer, M.; and McCallum, A. 2020. Probabilistic Case-based Reasoning for Open-World Knowledge Graph Completion. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, 4752–4765.
- Das, R.; Zaheer, M.; Thai, D.; Godbole, A.; Perez, E.; Lee, J.-Y.; Tan, L.; Polymenakos, L.; and McCallum, A. 2021. Case-based Reasoning for Natural Language Queries over Knowledge Bases. *arXiv preprint arXiv:2104.08762*.
- Gardner, M.; Grus, J.; Neumann, M.; Tafjord, O.; Dasigi, P.; Liu, N. F.; Peters, M.; Schmitz, M.; and Zettlemoyer, L. S. 2017. AllenNLP: A Deep Semantic Natural Language Processing Platform.
- Guo, J.; Zhan, Z.; Gao, Y.; Xiao, Y.; Lou, J.-G.; Liu, T.; and Zhang, D. 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*, 4524–4535. Florence, Italy: Association for Computational Linguistics.
- Gupta, V.; Shrivastava, A.; Sagar, A.; Aghajanyan, A.; and Savenkov, D. 2021. RETRONLU: Retrieval Augmented Task-Oriented Semantic Parsing. *arXiv preprint arXiv:2109.10410*.
- Hashimoto, T. B.; Guu, K.; Oren, Y.; and Liang, P. S. 2018. A retrieve-and-edit framework for predicting structured outputs. *Advances in Neural Information Processing Systems*, 31.
- Hazoom, M.; Malik, V.; and Bogin, B. 2021. Text-to-SQL in the Wild: A Naturally-Occurring Dataset Based on Stack Exchange Data. In *Proceedings of the 1st Workshop on Natural Language Processing for Programming (NLP4Prog 2021)*, 77–87. Online: Association for Computational Linguistics.
- Hossain, N.; Ghazvininejad, M.; and Zettlemoyer, L. 2020. Simple and effective retrieve-edit-rerank text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2532–2538.
- Huang, X. S.; Perez, F.; Ba, J.; and Volkovs, M. 2020. Improving Transformer Optimization Through Better Initialization. In III, H. D.; and Singh, A., eds., *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, 4475–4483. PMLR.
- Khandelwal, U.; Fan, A.; Jurafsky, D.; Zettlemoyer, L.; and Lewis, M. 2021. Nearest Neighbor Machine Translation. In *International Conference on Learning Representations*.
- Khandelwal, U.; Levy, O.; Jurafsky, D.; Zettlemoyer, L.; and Lewis, M. 2020. Generalization through Memorization: Nearest Neighbor Language Models. In *International Conference on Learning Representations*.
- Le Scao, T.; and Rush, A. M. 2021. How many data points is a prompt worth? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2627–2636.
- Lee, C.-H.; Polozov, O.; and Richardson, M. 2021. KaggleDBQA: Realistic Evaluation of Text-to-SQL Parsers. 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)*, 2261–2273. Online: Association for Computational Linguistics.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2020. Ro{BERT}: A Robustly Optimized {BERT} Pretraining Approach.
- Min, S.; Lyu, X.; Holtzman, A.; Artetxe, M.; Lewis, M.; Hajishirzi, H.; and Zettlemoyer, L. 2022. Rethinking the Role of Demonstrations: What Makes In-Context Learning Work? *arXiv preprint arXiv:2202.12837*.
- Pasupat, P.; Zhang, Y.; and Guu, K. 2021. Controllable Semantic Parsing via Retrieval Augmentation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 7683–7698.
- Pawlik, M.; and Augsten, N. 2015. Efficient computation of the tree edit distance. *ACM Transactions on Database Systems (TODS)*, 40(1): 1–40.
- Pawlik, M.; and Augsten, N. 2016. Tree edit distance: Robust and memory-efficient. *Information Systems*, 56: 157–173.
- Perez, E.; Kiela, D.; and Cho, K. 2021. True few-shot learning with language models. *Advances in Neural Information Processing Systems*, 34.
- Poesia, G.; Polozov, O.; Le, V.; Tiwari, A.; Soares, G.; Meek, C.; and Gulwani, S. 2022. Synchromesh: Reliable code generation from pre-trained language models. *arXiv preprint arXiv:2201.11227*.
- Popescu, A.-M.; Etzioni, O.; and Kautz, H. 2003. Towards a Theory of Natural Language Interfaces to Databases. In *Proceedings of the 8th International Conference on Intelligent User Interfaces*, 149–157.

- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research*, 21(140): 1–67.
- Rubin, O.; and Berant, J. 2021. SmBoP: Semi-autoregressive Bottom-up Semantic Parsing. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 311–324.
- Scholak, T.; Li, R.; Bahdanau, D.; de Vries, H.; and Pal, C. 2021. DuoRAT: Towards Simpler Text-to-SQL Models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 1313–1321.
- Scholak, T.; Schucher, N.; and Bahdanau, D. 2021. PICARD - Parsing Incrementally for Constrained Auto-Regressive Decoding from Language Models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Shaw, P.; Uszkoreit, J.; and Vaswani, A. 2018. Self-Attention with Relative Position Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, 464–468. New Orleans, Louisiana: Association for Computational Linguistics.
- Shin, T.; Razeghi, Y.; Logan IV, R. L.; Wallace, E.; and Singh, S. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 4222–4235.
- Socher, R.; Chen, D.; Manning, C. D.; and Ng, A. 2013. Reasoning with neural tensor networks for knowledge base completion. *Advances in neural information processing systems*, 26.
- Suhr, A.; Chang, M.-W.; Shaw, P.; and Lee, K. 2020. Exploring unexplored generalization challenges for cross-database semantic parsing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 8372–8388.
- Tang, L. R.; and Mooney, R. J. 2000. Automated Construction of Database Interfaces: Integrating Statistical and Relational Learning for Semantic Parsing. In *2000 Joint SIG-DAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, 133–141. Hong Kong, China: Association for Computational Linguistics.
- Wang, B.; Shin, R.; Liu, X.; Polozov, O.; and Richardson, M. 2020. RAT-SQL: Relation-Aware Schema Encoding and Linking for Text-to-SQL Parsers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 7567–7578.
- Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; Davidson, J.; Shleifer, S.; von Platen, P.; Ma, C.; Jernite, Y.; Plu, J.; Xu, C.; Scao, T. L.; Gugger, S.; Drame, M.; Lhoest, Q.; and Rush, A. M. 2020. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 38–45. Online: Association for Computational Linguistics.
- Xie, T.; Wu, C. H.; Shi, P.; Zhong, R.; Scholak, T.; Yasunaga, M.; Wu, C.-S.; Zhong, M.; Yin, P.; Wang, S. I.; et al. 2022. UnifiedSKG: Unifying and Multi-Tasking Structured Knowledge Grounding with Text-to-Text Language Models. *arXiv preprint arXiv:2201.05966*.
- Xu, P.; Kumar, D.; Yang, W.; Zi, W.; Tang, K.; Huang, C.; Cheung, J. C. K.; Prince, S. J.; and Cao, Y. 2021. Optimizing Deeper Transformers on Small Datasets. 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)*, 2089–2102. Online: Association for Computational Linguistics.
- Yu, T.; Wu, C.-S.; Lin, X. V.; Tan, Y. C.; Yang, X.; Radev, D.; Xiong, C.; et al. 2020. GraPPa: Grammar-Augmented Pre-Training for Table Semantic Parsing. In *International Conference on Learning Representations*.
- Yu, T.; Zhang, R.; Yang, K.; Yasunaga, M.; Wang, D.; Li, Z.; Ma, J.; Li, I.; Yao, Q.; Roman, S.; et al. 2018. Spider: A Large-Scale Human-Labeled Dataset for Complex and Cross-Domain Semantic Parsing and Text-to-SQL Task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 3911–3921.
- Zelle, J. M.; and Mooney, R. J. 1996. Learning to Parse Database Queries using Inductive Logic Programming. In *AAAI/IAAI*, 1050–1055. Portland, OR: AAAI Press/MIT Press.
- Zhao, Z.; Wallace, E.; Feng, S.; Klein, D.; and Singh, S. 2021. Calibrate Before Use: Improving Few-shot Performance of Language Models. In Meila, M.; and Zhang, T., eds., *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, 12697–12706. PMLR.

A Appendix

A.1 Details of the retriever used with ConcatCBR

The ConcatCBR baseline discussed in Section 4 first retrieves Text-SQL pairs for a given utterance \bar{x} , which are then concatenated with the utterance \bar{x} as an input to model’s encoder. Similar to Das et al. (2021); Poesia et al. (2022), we train a RoBERTa-BASE based sentence embedding model $\mathbf{E} : \mathcal{X} \mapsto \mathbb{R}^d$, that retrieves Text-SQL pairs $\{\bar{x}_r, \bar{q}_r\}$ from cases or training data having higher cosine similarity with the input utterance \bar{x} . The retriever is trained independent of Text-to-SQL models. For any two Text-SQL pairs (\bar{x}_i, \bar{q}_i) and (\bar{x}_j, \bar{q}_j) in the train set, we utilize normalized tree-edit-distance $\text{TED}(\bar{q}_i, \bar{q}_j) \in [0, 1]$ between queries (\bar{q}_i, \bar{q}_j) to supervise the cosine-similarity scores between sentence embeddings $\text{CSim}(\mathbf{E}(\bar{x}_i), \mathbf{E}(\bar{x}_j))$, such that pairs with low tree-edit-distance have higher cosine-similarity. We utilize APTED library (Pawlik and Augsten 2015, 2016) to compute tree-edit-distance between the relational algebra trees of the SQL queries. We modify the cost function of tree-edit-distance to ignore the leaf values and constants in the tree so that structurally similar trees are assigned lower distance. We normalize tree-edit-distance by size of the larger of two trees, so it lies in range $[0, 1]$.

More concretely, for a given Text-SQL pair in a training batch (\bar{x}_i, \bar{q}_i) we sample 15 more Text-SQL pairs $\{(\bar{x}_j, \bar{q}_j)\}_{j=1}^{15}$. Then we compute tree-edit-distances $\{\text{TED}(\bar{q}_i, \bar{q}_j)\}_{j=1}^{15}$, and cosine similarities $\{\text{CSim}(\mathbf{E}(\bar{x}_i), \mathbf{E}(\bar{x}_j))\}_{j=1}^{15}$ between sentences embeddings. The cosine-similarity scores are now supervised using tree-edit-distances as per loss in Equation 9.

$$w_{i,j} = \frac{\exp(1 - 2 \text{TED}(\bar{q}_i, \bar{q}_j))}{\sum_j \exp(1 - 2 \text{TED}(\bar{q}_i, \bar{q}_j))} \quad (9)$$

$$\mathcal{L}_{\mathbf{E}} = - \sum_{i,j} w_{i,j} \log \frac{\exp(\text{CSim}(\mathbf{E}(\bar{x}_i), \mathbf{E}(\bar{x}_j)))}{\sum_j \exp(\text{CSim}(\mathbf{E}(\bar{x}_i), \mathbf{E}(\bar{x}_j)))}$$

During inference on a new schema, we update the retrieval index with available cases from the new schema. Now, given a text query \bar{x} , we retrieve the 5 most similar examples as per cosine-similarity.

A.2 ConcatCBR baseline using T5

For a fair comparison, all the baselines in Section 5 are built on the top of the SmBoP architecture that utilizes non-auto-regressive bottom-up decoding. However, ConcatCBR architectures in prior works (Das et al. 2021; Pasupat, Zhang, and Guu 2021) utilize the standard auto-regressive Seq2Seq architectures. Hence, to further ensure a fair comparison, we explored ConcatCBR baselines built on the top of T5-large (Raffel et al. 2020) models. We utilize the implementation of UnifiedSKG library (Xie et al. 2022) for T5-large based Text-to-SQL models³ and modify it to concatenate input-output examples at T5’s encoder to get the ConcatCBR baseline. The T5-large models were initialized with an LM-Adapted⁴ version. The default T5-large based Text-to-SQL

³https://github.com/HKUNLP/UnifiedSKG/blob/main/configure/Salesforce/T5_large_finetune_spider_with_cell_value.cfg

⁴https://github.com/google-research/text-to-text-transfer-transformer/blob/main/released_checkpoints.md#lm-adapted-t511m100k

	world _1	car _1	cre _Doc	dog _kenn	flights _2	avg
EM						
T5	46.4	40.5	83	61.2	66	57.3
+ConcatCBR	+2.3	+1.1	-2.5	+10.9	-3.4	+1.7
EX						
T5	43.1	41.7	92.5	66.7	77.6	61
+ConcatCBR	+4.1	+2.8	-8.2	+7.5	-5.4	+0.7

Table A1: EM and EX performance of a T5-large based Text-to-SQL model and gains from its ConcatCBR extension. Similar to Table 2, we observe that ConcatCBR does not provide consistent gains over the unadapted model. For two schemas the gains are negative, and gains micro-averaged (avg) across five schemas are small.

model finetuned on Spider dataset has a comparable EM accuracy of 57.3 w.r.t. SmBoP’s 57.2 micro-averaged across five schemas.

Table A1 presents the EM and EX accuracies of the default T5-large based Text-to-SQL model, and gains obtained by its ConcatCBR extension. Similar to Table 2, gains from ConcatCBR are not consistent. For two out of five databases, we observe a drop in EM and EX. The micro-averaged gains across all five schemas are also small. Thus, our observations remain largely consistent with ConcatCBR baseline based on SmBoP.

A.3 StructCBR gains relative to finetuning

Table A2 shows gains of StructCBR over SmBoP relative to gains of finetuning over SmBoP. Relative Gain = $\frac{\text{+StructCBR} + \text{Finetuning}}{\text{+Finetuning}} \times 100$. StructCBR instantly obtains 20% to 83% of EM gains achieved by finetuning for 100 epochs across different schemas. For 3 out of 5 schemas relative gains are more than 50%. Micro-averaged across all examples, StructCBR achieves 42.8% of gains relative to finetuning. It is important to note here that we do not position StructCBR as an alternative to finetuning, but instead as a method of utilizing available cases for instantaneously improving model’s performance until the next cycle of finetuning becomes affordable.

	world _1	car _1	cre _Doc	dog _kenn	flights _2	avg
SmBOP	46	43.3	84.3	66.6	55.8	57.2
+StructCBR	+2.6	+6.1	+6.3	+3.4	+5.5	+4.6
+Finetuning	+6.0	+11.7	+7.6	+4.8	+27.2	+10.7
Relative Gain (%)	43.9	52.3	83.7	71.5	20.1	42.8

Table A2: StructCBR instantly obtains 20% to 83% EM gains achieved by finetuning for 100 epochs across different schemas. On 3 out 5 schemas gains are higher than 50% of gains from finetuning. Averaged (avg) across all examples, gains are 42.8% of finetuning.

A.4 Additional hyperparameter and training details

The baseline SmBoP model is based on RoBERTa-base architecture and contains four RAT layers. The total number of

trainable parameters in this more is roughly 133M. Adding StructCBR module on top of SmBoP introduces only 2.53% additional parameters. In particular, the transformer model TX_ϕ used for computing G_ϕ representations in Equation 2 is composed of two transformer blocks of hidden_size=256, attention_heads=8, feedforward_dim=1024. All the experiments were performed on a NVIDIA RTX 3060 GPU, and a NVIDIA RTX A6000. Training times for various Text-to-SQL models usually varied from 12 hours to 16 hours.

A.5 Results with a larger model size

Constrained by limited computing resources for training large models, we conducted all the experiments and ablations in Section 5 using SmBoP models initialized with a pretrained RoBERTa-BASE checkpoint (Liu et al. 2020) followed by 4 RAT layers (Wang et al. 2020). In this section, we validate our key results on a larger model size, by repeating the experiment reported in Table 2 of Section 5. More specifically, we replace RoBERTa-BASE in SmBoP with a pretrained RoBERTa-LARGE checkpoint based on grammar augmented pre-training (Yu et al. 2020) as used by Rubin and Berant (2021). The number of RAT layers is also increased from 4 to 8. We refer to the larger SmBoP architecture as SmBoP-LARGE (367M parameters) and the base-sized SmBoP model used in Section 5 as SmBoP-BASE (133M parameters). Similar to Table 2, Table A3 compares different inference-time methods for adapting the SmBoP-LARGE model to new schemas. The performance of SmBoP-BASE is reported just as a reference. For each adaptation method (ConcatCBR, GTM, StructCBR), we report it’s gains over the SmBoP-LARGE model. We consistently observe positive gains from StructCBR across all the metrics on four out of five schemas. In contrast, the gains from prior inference time adaptation methods are not consistently positive across different metrics and schemas. It is also interesting to note that the gains in top-K accuracy metrics (BEM and BEX) from StructCBR are significantly higher than prior methods and consistently positive across all the schemas. Better performance in top-K metrics indicates that boosting candidate subtree scores via StructCBR improves the recall of correct subtrees during beam decoding.

A.6 Additional Related Work

Incontext Learning with LLMs: Learning new tasks via carefully designed exemplar prompts (Shin et al. 2020; Le Scao and Rush 2021; Perez, Kiela, and Cho 2021) for large language models (LLMs) like GPT-3 (Brown et al. 2020) and Codex (Chen et al. 2021) has shown promise for many tasks (Zhao et al. 2021; Min et al. 2022) without requiring to finetune model parameters. For Text-to-SQL, Poesia et al. (2022) retrieve and encode input-specific prompts in a way similar to the ConcatCBR method. They differ from ConcatCBR as they do not explicitly train the LLM to perform CBR with the extracted prompts. While their method achieves impressive performance by using LLMs, training significantly (100x) smaller Text-to-SQL models on task-specific datasets like Spider outperforms their LLMs by considerable margins.

	EM	EX	BEM	BEX
$s_{\text{new}} = \text{world_1}$, $ \mathcal{D}_{\text{test}} = 89$				
SmBoP-BASE	46.1	36.3	55.4	49.0
SmBoP-LARGE	50.6	42.3	58.1	49.8
ConcatCBR	-2.3	-0.4	-1.9	+0.7
GTM	-0.8	-0.8	+0.7	+1.1
StructCBR (Ours)	-1.5	-0.7	+3.4	+2.2
$s_{\text{new}} = \text{car_1}$, $ \mathcal{D}_{\text{test}} = 60$				
SmBoP-BASE	43.3	45.6	50.6	53.9
SmBoP-LARGE	53.3	57.2	66.7	67.8
ConcatCBR	+3.9	-6.1	-2.2	-2.8
GTM	+6.6	+3.3	+3.3	-0.5
StructCBR (Ours)	+3.3	+3.3	+1.1	+1.7
$s_{\text{new}} = \text{cre_Doc_Template_Mgt}$, $ \mathcal{D}_{\text{test}} = 53$				
SmBoP-BASE	84.3	89.3	94.3	96.2
SmBoP-LARGE	90.6	95.6	96.2	98.1
ConcatCBR	+3.7	+0.6	+3.2	+1.3
GTM	-1.9	-1.2	+1.3	+0.6
StructCBR (Ours)	+3.1	+1.9	+3.2	+1.9
$s_{\text{new}} = \text{dog_kennels}$, $ \mathcal{D}_{\text{test}} = 49$				
SmBoP-BASE	66.6	59.2	80.3	74.8
SmBoP-LARGE	70.1	62.6	81.0	74.1
ConcatCBR	+1.3	+2.7	+2.7	+6.2
GTM	-0.7	-0.7	-1.4	+4.8
StructCBR (Ours)	+2.1	+3.4	+5.4	+8.2
$s_{\text{new}} = \text{flight_2}$, $ \mathcal{D}_{\text{test}} = 49$				
SmBoP-BASE	55.8	67.4	59.2	80.3
SmBoP-LARGE	59.9	65.3	66.0	75.5
ConcatCBR	+6.1	+10.9	+5.4	+11.6
GTM	+1.3	+0.7	+4.7	+7.5
StructCBR (Ours)	+4.1	+3.4	+14.9	+12.2
Micro-Average , $ \mathcal{D}_{\text{test}} = 300$				
SmBoP-BASE	57.2	56.3	66.0	67.6
SmBoP-LARGE	62.6	61.8	71.3	70.1
ConcatCBR	+2.3	+1.0	+1.1	+2.8
GTM	+1.2	+0.2	+1.9	+2.3
StructCBR (Ours)	+2.1	+1.9	+5.3	+4.7

Table A3: Results using SmBoP-LARGE: We compare StructCBR with prior inference time adaptation methods (ConcatCBR and GTM) for adapting SmBoP-LARGE on 5 different schemas. SmBoP-BASE and SmBoP-LARGE rows provide the performance of unadapted model SmBoP models, while other rows report gains over SmBoP-LARGE after adaptation. SmBoP-BASE numbers are just of reference. Micro-average refers to numbers averaged over all the test instances spanning across the five schemas. $|\mathcal{D}_{\text{test}}|$ refers to size of the test-set, and s_{new} provides the schema name.

A.7 Anecdotes

We include some anecdotal examples where StructCBR fixes the mistakes made by SmBoP by utilizing case examples.

Text	What is the code of airport that has the highest number of flights?
Incorrect SQL	<code>SELECT flights.sourceairport FROM flights GROUP BY flights.sourceairport ORDER BY SUM(flights.flightno) DESC LIMIT 1</code>
Corrected SQL	<code>SELECT airports.airportcode FROM airports JOIN flights ON airports.airportcode = flights.sourceairport GROUP BY flights.sourceairport ORDER BY COUNT(*) DESC LIMIT 1</code>
Case Text	Give the code of the airport with the least flights.
Case SQL	<code>SELECT airports.airportcode FROM airports JOIN flights ON airports.airportcode = flights.sourceairport GROUP BY flights.sourceairport ORDER BY COUNT(*) ASC LIMIT 1</code>
Mistake(s) Fixed	Correct column selection, Add missing join, Replace SUM by COUNT
Text	Which airlines have a flight with source airport AHD?
Incorrect SQL	<code>SELECT flights.airline FROM airlines JOIN flights ON airlines.uid = flights.airline JOIN airports ON flights.sourceairport = airports.airportcode WHERE airports.airportname = 'AHD'</code>
Corrected SQL	<code>SELECT airlines.airline FROM airlines JOIN flights ON airlines.uid = flights.airline WHERE flights.sourceairport = 'AHD'</code>
Case Text	What are airlines that have flights arriving at airport 'AHD'?
Case SQL	<code>SELECT airlines.airline FROM airlines JOIN flights ON airlines.uid = flights.airline WHERE flights.destairport = 'AHD'</code>
Mistake(s) Fixed	Drop additional condition on Join
Text	What are the population and life expectancies in Brazil?
Incorrect SQL	<code>SELECT country.population , country.lifeexpectancy FROM city JOIN country ON city.countrycode = country.code WHERE country.name = 'Brazil'</code>
Corrected SQL	<code>SELECT country.lifeexpectancy , country.population FROM country WHERE country.name = 'Brazil'</code>
Case Text	What are the region and population of Angola?
Case SQL	<code>SELECT Population , Region FROM country WHERE Name = 'Angola'</code>
Mistake(s) Fixed	Drop the join condition
Text	Return the codes of countries that do not speak English and do not have Republics for governments.
Incorrect SQL	<code>SELECT country.code FROM country WHERE countrylanguage.language = 'English' INTERSECT SELECT countrylanguage.countrycode FROM country WHERE country.governmentform != 'Republic'</code>
Corrected SQL	<code>SELECT country.code FROM country WHERE country.governmentform != 'Republic' EXCEPT SELECT countrylanguage.countrycode FROM countrylanguage WHERE countrylanguage.language = 'English'</code>
Case Text	What are the codes of the countries that do not speak English and whose government forms are not Republic?
Case SQL	<code>SELECT country.code FROM country WHERE country.governmentform != 'Republic' EXCEPT SELECT countrylanguage.countrycode FROM countrylanguage WHERE countrylanguage.language = 'English'</code>
Mistake(s) Fixed	Set operation and Equality

Table A4: Anecdotal examples where mistakes in output SQLs generated by SmBoP are fixed with help of StructCBR through related examples in case memory. We show only one of the related examples from cases for brevity