# Synthesizing Queries via Interactive Sketching

Osbert Bastani University of Pennsylvania USA obastani@seas.upenn.edu Xin Zhang Peking University China xin@pku.edu.cn Armando Solar-Lezama MIT USA asolar@csail.mit.edu

## **Abstract**

We propose a novel approach to program synthesis, focusing on synthesizing database queries. At a high level, our proposed algorithm takes as input a sketch with soft constraints encoding user intent, and then iteratively interacts with the user to refine the sketch. At each step, our algorithm proposes a candidate refinement of the sketch, which the user can either accept or reject. By leveraging this rich form of user feedback, our algorithm is able to both resolve ambiguity in user intent and improve scalability. In particular, assuming the user provides accurate inputs and responses, then our algorithm is guaranteed to converge to the true program (i.e., one that the user approves) in polynomial time. We perform a qualitative evaluation of our algorithm, showing how it can be used to synthesize a variety of queries on a database of academic publications.

## 1 Introduction

Program synthesis has emerged as a promising way to help users write programs—e.g., it has been leveraged to generate highly optimized bit manipulation programs [31, 32], string processing programs [13, 25], and database queries [11, 19, 34, 36]. These techniques enable the user to focus on specifying their high-level intent. For example, one approach is for the user to provide a logical specification such as a sketch annotated with logical constraints [31, 32]. Then, the synthesizer automatically searches the space of programs to try and identify one that satisfies the given specification.

A key challenge is that the user may be unable to provide a precise, logical specification. One reason is that they may be an end-user who knows exactly what task they want to achieve, but is unfamiliar with logical specifications. In such cases, existing systems allow the user to provide *ambiguous specifications* such as input-output examples [11–13, 25, 34] or natural language descriptions [19, 36].

An alternative reason is that the user is familiar with logical specification, but uncertain about the task they want to achieve. For example, suppose a user wants to write a query to extract author names from a large database of academic publications. They may be uncertain about the task at hand—e.g., they may be unfamiliar with the database schema, so they are unable to write down a precise logical specification or even an ambiguous one. Alternatively, the task may be

inherently uncertain—e.g., the user might want to run the query on a database whose schema changes over time, or on multiple databases with different schemas. In these cases, the user may want to provide an *uncertain specification*—i.e., a precise, logical specification of not only their knowledge of their task, but also their uncertainty about their task.

We propose a specification language and corresponding program synthesis algorithm that enables users to express uncertain specifications. We assume the user is knowledgeable and can write both programs, but doing so requires significant effort to first resolve their uncertainty about their task. We focus on the setting of a user who is a data scientist trying to write a query to extract data from a large database. Their uncertainty may be because they are unfamiliar with the database schema, or because they want to run the same query on multiple databases with different schemas.

To use our system, the user provides a specification in the form of a sketch, which includes the structure of the query (including all select and project operations), but can have holes corresponding to tables (to be filled using a table constant or a sequence of inner-join operations) and columns (to be filled with a column name). In other words, the key portion of the query that the user can leave unspecified is the sequence of tables should be joined together to construct the desired flat table. The search space of such sequences can be very large—many datasets have dozens of tables and it is exponential in the number of tables to be joined.

To express their uncertainty, our language supports specifications in the form of soft constraints. In particular, the user can include soft constraints on expressions in the sketch; these constraints should encode the user's expectations about the value that should be obtained by evaluating the expression once all holes in the sketch have been filled. For example, to specify that a column should contain author names, the user might include a soft constraint saying that the values in that column are strings, some of which contain "Church".

Then, our synthesis algorithm fills the holes in the provided sketch, relying critically on user interaction to resolve uncertainty. We assume the user can recognize the correct query when they see it. This assumption may seem at odds with the fact that the user is uncertain about their intent, but manually resolving ambiguity can require significantly more effort than being actively guided through the process. In particular, our algorithm actively focuses the user's efforts on understanding portions of the database that are relevant to resolving the uncertainty they specified in their sketch.

aid	name		aid	pid
	Alan M. Tanina		0	0
	Alan M. Turing Alonzo Church		0	1
1			1	2
authors				
			writes	

pid	title	year			
0	Computability and $\lambda$ -definability	1937			
1	Intelligent machinery	1948			
2	A set of postulates for the foundation of logic	1932			
nublications					

**Figure 1.** Example of a database of computer science publications and their authors. This database includes three tables—one of authors ("authors"), one of published papers ("publications"), and one that links the previous two ("writes").

At a high level, our algorithm keeps track of a sketch (i.e., a program with holes), which is initialized to be the given sketch. Then, it iteratively interacts with the user to fill the holes in the sketch. At each iteration, our algorithm proposes a candidate refinement of the current sketch (i.e., proposes to fill a single hole with some expression) to the user, who either accepts the refinement if it matches the true program or rejects it otherwise. The added expression may contain new holes that must be filled in subsequent iterations. If the proposed refinement is accepted, then the current sketch is updated; otherwise, the current sketch remains unchanged. Either way, our algorithm continues to ask additional questions until the sketch is concrete (i.e., there are no remaining holes), at which point it is returned by our algorithm. We assume the user always correctly accepts or rejects the candidate refinement; in practice, they can backtrack if they realize their choices are incorrect.

Our algorithm enjoys two key advantages that derive from the rich feedback that it solicits from the user. First, it is guaranteed to find the true program (assuming the structure of the sketch is correct and the user always answers accurately). In particular, the true program must be attainable from some sequence of refinements of the given sketch, and the user affirms each refinement made by the algorithm, including the final program. This guarantee holds even if the constraints provided by the user are underspecified—i.e., it implicitly resolves ambiguity in the specified user intent. Of course, better specifications can lead to faster convergence.

Second, our algorithm is guaranteed to identify the true program in a number of iterations that is polynomial in the size of the true program. First, each iteration of our algorithm is efficient since our algorithm only has to search over the set of possible refinements, which is polynomial in the size of the sketch. Furthermore, our algorithm is guaranteed to identify the true program in a polynomial number of iterations. This

guarantee holds for any choice of questions; in practice, we try to minimize the overall number of questions asked.

The key challenge is designing a space of possible refinements such that the user can quickly and accurately respond to questions. Since we assume the user is a knowledgeable programmer, we assume they can recognize the correct program by inspecting both the program code and the data—e.g., in our setting, the user bases their decision on both the query and the information in the database. The key advantage is that they only need to understand the portions of the database that are relevant to questions asked by our algorithm. Assuming our algorithm asks good questions, then we can significantly reduce the user's workload.

To respond to a question, the user must be able to determine whether the proposed refinement is correct without seeing the whole program. Intuitively, the refinement should always be a concrete transformation of the input that the user can inspect and validate—e.g., a key refinement our algorithm can propose is to replace a hole with an inner join of a concrete table with another hole. Thus, the user only needs to check that the concrete table is part of the sequence of inner-join operations to construct the desired flat table.

Finally, we evaluate our approach on a dataset of queries to a database of academic publications [19]. We perform a small-scale user study to determine whether users can write reasonable sketches, as well as interact with our synthesis algorithm to refine a sketch. We find that 80% of sketches are successful, and 93% of questions were answered correctly, suggesting that users can use our system even with very little experience. Furthermore, we show that for all of our sketches, our algorithm quickly converges to the true query.

In summary, our contributions are:

- A novel formulation of program synthesis that interacts with the user to iteratively refine a user-provided specification in the form of a sketch (Section 3).
- An algorithm based on this approach in the context of synthesizing database queries (Section 4).
- An implementation of our algorithm in a tool called ISQL (Section 5), <sup>1</sup> and an evaluation that illustrates how our approach can be used to synthesize a variety of database queries (Section 6).

## 2 Overview

As a motivating example, consider a user (e.g., a data scientist) who wants to query a database of academic publications, including information such as conferences/journals, titles, abstracts, authors, citations, etc. Their goal may be to perform data exploration, to compute some basic statistics of the data, or to construct a flat table that will be used to train a machine learning model. Potential queries include which authors were active in a given year (e.g., published a paper in

<sup>&</sup>lt;sup>1</sup>ISQL stands for "interactive SQL"; we pronounce it "icicle".

```
SELECT ??c name:column
FROM (??t:table {(contains ??c_name:column ".*Church.*")
                 AND (1900 <= ??c_year:column <= 2020)})
WHERE ??c_year:column = 1948
                             sketch
               authors
??t:table
               INNER-JOIN ??t_new:table
               ON ??c new0:column = ??c new1:column
                  candidate production sequence
SELECT ??c_name:column
FROM (authors
      INNER-JOIN ??t_new:table
      ON ??c new0:column = ??c new1:column
      {(contains ??c_name:column ".*Church.*"
       AND (1900 <= ??c_year:column <= 2020)})
WHERE ??c_year:column = 1948
                           refinement
SELECT name
FROM (authors
      INNER-JOIN (writes
                  INNER-JOIN publications
                  ON writes.pid = publications.pid)
      ON authors.aid = writes.aid)
      {(contains name ".*Church.*") AND (1900 <= year <= 2020)})
WHERE publications.year = 1948
                           completion
```

**Figure 2.** Examples illustrating the concepts used in our algorithm. The goal is to select authors who published papers in 1948. The user provides the sketch (first row). Then, our algorithm proposes a candidate sequence of productions to use to refine the sketch (second row). If accepted by the user, the productions are applied to the sketch to obtain the refined sketch (third row). This interactive process continues iteratively until the sketch has no more holes (last row). Soft constraints are shown in red.

1948), which authors have the most citations, which authors have cited a given academic, etc.

While the user may be knowledgeable of SQL, they may not be familiar with the schema of the database. For example, many databases are poorly documented, including ones used for data science. As a consequence, it is challenging for the user to write the desired database query from scratch, since they have to spend time reading the database documentation and understanding its schema. However, if they are shown a candidate query, then they can easily check the corresponding tables and columns in the database to determine whether it is correct. This problem arises in data science tasks, where the user is a data scientist who is exploring a number of large datasets with the goal of deriving some kind of insight from the data. Many of these datasets are unfamiliar to the data scientist and poorly documented, making it time consuming to identify which tables and columns to use to fill in different parts of the desired query.

As a concrete example, consider the database shown in Figure 1, which includes three tables—one containing authors, one containing computer science publications, and a third that links the previous two. Suppose the user wants to select all authors who published a paper in 1948. To do so, the user can use the query on the last line of Figure 2. In this query, the writes table relates authors to publications; thus, the query joins publications with writes, and then joins the result with authors to obtain a flat table with both authors and their publications. Then, this query selects the authors that have a publication in the year 1948.

## 2.1 Initial User Specification

We assume that the user is familiar with the query language (e.g., SQL) and knows what the desired data should look like (e.g., publication years are mostly between 1900 and 2020) but is unfamiliar with the database schema and does not know where in the database the desired information is located. Furthermore, we assume the database is large and possibly poorly documented, which is typical of many datasets used in data science tasks. In particular, the user is capable of writing a skeleton of the query (e.g., the structure of the output they ultimately want to construct, and any aggregation operators they want to apply), but does not know which tables and columns to use inside this query.

In our example, we expect that the user knows the structure of the outermost select statement (which combines selecting rows with year 1948 and a project that retains only the author name column). The key challenge is that they do not know how to construct the flat table that includes both authors and publications. More precisely, the key challenge is that the user does not know the sequence of tables to inner-join to obtain this flat table. Our algorithm is designed to help the user discover this sequence.

**Sketches.** To this end, we assume the user is able to provide the sketch shown on the first row of Figure 2. This sketch outlines the structure of the select and project operators to be applied to the flat table. The flat table, along with the columns to be selected and projected, are left as holes in this sketch. In particular, there are three holes in this sketch:

- The hole ??c\_name: column is named c\_name, has type column, and corresponds to the unknown author name column.
- The hole ??t:table is named t, has type table, and corresponds to the unknown flat table.
- The hole ??c\_year: column is named c\_year, has type column and corresponds to the unknown publication year column.

In general, holes either have type column or table. The names associated with holes are used to link different holes in the sketch that are known to have the same value—e.g., in the example sketch, there are two holes named c\_year, which indicates that they must be filled using the same column.

aid	name	pid	title	year
0	Alan M. Turing	0		1937
0	Alan M. Turing	1		1948
1	Alonzo Church	2		1932

**Figure 3.** The table obtained by evaluating the expression (authors INNER-JOIN ...) in the complete program on the last line of Figure 2 on the database in Figure 1.

*Soft constraints.* In addition, the user can also provide specifications that encode how the sketch should be filled. Unlike traditional specifications, which provide hard constraints on the semantics of the sketch, these specifications are soft constraints that encode expectations that the user has about likely properties of the semantics. These soft constraints can be used to assign a score to a concrete program that indicates how well the program matches the user's expectation. <sup>2</sup>

In our example sketch, the user has provided three soft constraints on the expected semantics. In the context of database queries, these constraints encode the user's expectations about the properties of the data in a table constructed while evaluating the query. For example, consider the portion

```
??t_new:table
{(contains ??c_name:column ".*Church.*") ...}
```

of the sketch. This portion consists of a hole ??t\_new:table with name t\_new and type table. The soft constraint is the expression appearing in the curly braces  $\{\ldots\}$ . This constraint applies to the preceding expression—in this case, the hole named t\_new. Semantically, it maps the table t obtained by evaluating the preceding expression to a real-valued score  $s \in \mathbb{R}$ . In our example, the soft constraint contains (c,r) maps the table t to 1 if c contains a string matching the regular expression r, and 0 otherwise. Intuitively, this constraint encodes the user expectation that the table contains a column of strings, and that some of these strings (which should be author names) contain the substring "Church". This constraint is soft because they are interpreted in a way that does not prune a possible program just because it is violated.

We evaluate the soft constraint in the context of a *completion* of the sketch—i.e., a concrete program obtained by filling all the holes in the sketch with concrete expressions. Then, the expression filling the hole named  $t_n$ ew evaluates to some table t; then, we apply the soft constraint to t to obtain a score s. For example, in Figure 2, the concrete program on the last line is a completion of the sketch on the first line. In this completion, the expression above becomes

in which case (contains name ".\*Church.\*") evaluates to 1. In particular, the inner-join evaluates to the table shown in Figure 3. Then, the soft constraint says the user expects that some of the values in the name column of this table contain the substring "Church". Since one of the values in this column satisfies this property, we assign score 1.

Taken together, we can assign a score to a completion of a sketch by evaluating the completion, evaluating the soft constraints, and summing the resulting scores. Higher scoring completions correspond to concrete programs that are more likely according to the user-provided specification.

#### 2.2 Interactive Program Synthesis Algorithm

Our algorithm interacts with the user to determine how to fill the holes in the given sketch. It keeps track of a sketch P, which is initialized to the given sketch. At each iteration, it proposes a candidate *refinement* P' of P, which modifies P by filling a single hole in P with an expression. The user either accepts the refinement if it matches the "true" program that the user is aiming to write (in which case we update  $P \leftarrow P'$ ), or rejects it otherwise (in which case a different refinement is proposed). This process continues until P is concrete (i.e., it has no holes), at which point our algorithm returns P.

Selecting a candidate refinement. The key step performed by our algorithm on each iteration is to select a candidate refinement P' of P with which to question the user. In Figure 2, an example of a refinement is shown on the third line. This refinement is constructed from the sketch on the first line using the productions shown on the second line, which says that the hole named t should be replaced with the expression

```
authors INNER-JOIN ??t_new:table ON ...
```

Note that this expression contains new holes; if the user accepts this refinement, then our algorithm will need to fill in these holes on subsequent iterations.

Intuitively, our goal is to choose the question that elicits the most information about the true program—i.e., the question should cut down the search space by as much as possible. To formalize this intuition, we use the user-provided sketch to induce a probability distribution over completions of that sketch. Then, we can ask for the refinement that prunes the largest number of completions in expectation, where each completion is weighted by its probability.  $^3$  Computing this refinement is challenging due to the exponential size of the search space over programs. Instead, we use MCMC to randomly sample a finite number of completions  $\overline{P}$  of P, and then choose the refinement P' of P according to an estimate of the above objective using these samples.

<sup>&</sup>lt;sup>2</sup>We use soft constraints since the user might be uncertain about the actual contents of the database; we could easily include hard constraints if desired.

<sup>&</sup>lt;sup>3</sup>An optimal approach would be to optimize this objective over a sequence of questions; however, this approach quickly becomes intractable. In fact, the performance of the greedy strategy (in terms of expected number of questions used) is a log-factor of the performance of the optimal approach [7].

*User interaction.* Once our algorithm has selected a refinement P', it shows this refinement to the user. The user must either accept P' if the true program is a completion of P', or reject it otherwise. In our example in Figure 2, the true program on the last line is a refinement of the refinement on the third line; thus, the user accepts this refinement; then, P is updated to be P', and the interactive process continues.

A key constraint is that we need to choose the space of possible refinements to ensure that users can understand whether the refinement matches the true program. As discussed above, we assume that users are familiar with SQL, and the key challenge is making sure they can understand whether the tables and columns in the database are the right ones to use in various parts of the query.

In particular, as discussed above, the primary purpose of our algorithm is to determine the sequence of joins that are needed to construct the flat table that the user needs to perform subsequent tasks. In particular, there are two kinds of refinements considered by our algorithm: (i) specifying an inner-join expression of the form

$$?? \rightarrow t \text{ INNER-JOIN } ?? \text{ ON } ?? = ??$$

or a single table to use—i.e.,  $?? \rightarrow t$ , in which case a summary of t is shown to the user, or (ii) specifying which column to use in a select, project, or inner-join operation—i.e.,  $?? \Rightarrow c$ , in which case a summary of c is shown to the user.

In our example in Figure 2, the refinement on the third line is obtained from the sketch by filling the hole named t with the table authors (inner-joined with another, currently unspecified table). Thus, our algorithm would show a summary of the authors table in Figure 1 to the user (e.g., the first few rows of this table). This information suffices for the user to decide whether to accept the refinement, since they see that it includes author names that want included in the flat table constructed by the true program.

If the user accepts, we update P to equal P'. We continue the iterative process until P is concrete, at which point it returns P. Because we question the user at every step (including the final program P), we guarantee that the final program is the one desired by the user.

## 3 Sketch Language

We consider a domain-specific language (DSL)  $\mathcal{D}$  of database queries based on a fragment of SQL that only includes select, project, and inner-join operations. Its syntax is a context-free grammar  $\mathcal{G} = (V, \Sigma, \mathcal{R}, \mathcal{Q})$  with non-terminals V, terminals  $\Sigma$ , productions  $\mathcal{R}$ , and start symbol  $\mathcal{Q}$ . This grammar is shown in Figure 4. Projection of a table t onto a list of columns  $c_1, ..., c_n$  is denoted  $\Pi_{c_1, ..., c_n}(t)$ , selection of rows that satisfy a predicate  $\psi$  from table t is denoted  $\sigma_{\psi}(t)$ , and the inner-join of tables  $t_1$  and  $t_2$  on column  $c_1$  in  $t_1$  and  $c_2$  in  $t_2$  is denoted  $t_1 \bowtie_{c_1, c_2} t_2$ . The semantics  $\llbracket \cdot \rrbracket : \mathcal{L}(\mathcal{G}) \to \mathcal{T}$  maps programs  $P \in \mathcal{L}(\mathcal{G})$  to tables  $t \in \mathcal{T}$ . They ignore the soft constraints  $\Phi$ ; otherwise, they are standard, so we omit them.

$$\begin{array}{lll} Q & ::= \Pi_{C, \dots, C}(S) \; \{\Phi\} \\ S & ::= \sigma_{\Psi}(I) \; \{\Phi\} & T & ::= t_1 \mid \dots \mid t_n \\ I & ::= T \; \{\Phi\} \mid T \bowtie_{C,C} I \; \{\Phi\} & C & ::= c_1 \mid \dots \mid c_m \\ \Psi & ::= \operatorname{true} \mid C \; R \; V \mid \Psi \wedge \Psi \mid \Psi \vee \Psi & X & ::= x_1 \mid \dots \mid x_\ell \\ R & ::= \leq \mid < \mid = \mid > \mid \geq \\ \Phi & ::= \operatorname{true} \mid X \in C \mid \operatorname{contains}(c,r) \mid C \; U \; X \mid \Phi \wedge \Phi \\ U & ::= \lesssim \mid \; \simeq \mid \; \gtrsim \end{array}$$

**Figure 4.** Syntax of database queries, with start symbol Q. Here,  $t_1, ..., t_n$  are tables,  $c_1, ..., c_m$  are columns, and  $x_1, ..., x_\ell$  are constants (i.e., integers, floats, strings, and regular expressions). The language is based on [36], except there are no aggregation operations, and queries are normalized so projection and selection operations are executed last. We permit two kinds of holes: (i) a table in a sequence of innerjoin operations (corresponding to nonterminal I), and (ii) columns in any operation (corresponding to nonterminal C).

Note that we have constrained to expressions of the form

$$\Pi_{C,...,C}(\sigma_{\Psi}(T \bowtie_{C,C} \cdots \bowtie_{C,C} T)).$$

In general, by using the relational algebra, any composition of select, project, and inner-join operations can be equivalently expressed in this form. Since our focus is on the sequence of inner-join operations, we assume that the user will specify the structure of the project and select operations, and only leave tables in the inner-join operation on the inside as holes (columns can be left as holes anywhere in the query). Note that this grammar also includes soft constraints  $\phi$  on tables t (denoted t { $\phi$ }), which we discuss below.

**Notation.** We establish some standard notation. Consider a sequence  $\alpha = A_1...A_k \in (V \cup \Sigma)^*$ . Suppose that  $A_i \in V$ ; then, we can apply a production  $A_i \to A_{i1}...A_{ih}$  to obtain

$$\alpha' = A_1...A_{i-1}A_{i1}...A_{ih}A_{i+1}...A_k,$$

and we denote this relationship by  $\alpha \Rightarrow \alpha'$ . Furthermore, if there exists a sequence  $\alpha \Rightarrow \alpha' \Rightarrow ... \Rightarrow \alpha''$ , then we say  $\alpha''$  can be *derived* from  $\alpha$ , which we denote by  $\alpha \stackrel{*}{\Rightarrow} \alpha''$ .

We refer to a sequence  $\alpha$  such that  $A \stackrel{*}{\Rightarrow} \alpha$  for some nonterminal  $A \in V$  as an *expression*. We let  $\mathcal{L}(\mathcal{G},A)$  denote the *concrete* expressions  $\alpha \in \Sigma^*$  that can be derived from A—i.e.,  $A \stackrel{*}{\Rightarrow} \alpha$ . Note that  $\mathcal{L}(\mathcal{G},Q) = \mathcal{L}(\mathcal{G})$ —i.e., the space of programs defined by the grammar  $\mathcal{G}$  is the set of expressions that can be derived from the start symbol Q.

**Sketches.** Our algorithm keeps track of programs that have holes. In particular, a *sketch* [31, 32] is a a sequence  $P \in (V \cup \Sigma)^*$  such that P can be derived from Q—i.e.,  $Q \stackrel{*}{\Rightarrow} P$ . We refer to a nonterminal in P as a *hole*. We restrict to holes A

that are either A=I (i.e., tables in the sequence of inner-join operations) or A=C (i.e., columns). We associate a name s (i.e., a string) with each hole in P; a name identifies different holes that should be filled using identical expressions. A sketch is *complete* (also called a *concrete program*) if  $P \in \Sigma^*$ —i.e., it has no holes (which implies that  $P \in \mathcal{L}(\mathcal{G})$ ). In our language, the sketch on the first line of Figure 2 is written

$$P_{\text{author}} = \Pi_{C:\text{c\_name}} \left( \sigma_{C:\text{c\_year}=1948} \left( I : \text{t } \{\phi\} \right) \right)$$

$$\phi = (\text{contains}(C : \text{c\_name}, \text{"*Church.*"})$$

$$\land (C : 1900 \leq \text{c\_year} \leq 2020).$$

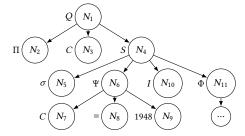
We have dropped the types from holes A, since they are determined by the value of the hole—i.e., table if A = I or column if A = C. Instead, we have used the notation A : s to denote the hole A with associated name s. We have also dropped soft constraints  $\{\phi\}$  when  $\phi = \text{true}$ .

**Abstract syntax trees.** Internally, our algorithm represents a sketch P using its abstract syntax tree (AST), which is a representation of the derivation of P in G. For convenience, we use P to denote both the sequence  $P \in (V \cup \Sigma)^*$  as well as its AST. We denote the nodes of P by nodes(P), the internal nodes by internal(P), and the leaves by leaves(P).

Each node N in P is associated with a symbol  $A_N \in V \cup \Sigma$ , which is the symbol associated with N in the derivation of P. An internal node is always labeled with a nonterminal. If P is complete, then each leaf node of P is labeled with a terminal; otherwise, a leaf node of P may be labeled with either a terminal or a nonterminal. Note that holes correspond to leaf nodes of P labeled with a nonterminal; we denote the set of holes by holes  $(P) \subseteq \text{leaves}(P)$ .

Finally, we use  $\alpha_N$  to denote the subtree of P at N. Note that  $\alpha_N$  can also be thought of as a subexpression  $\alpha_N \in (\Sigma \cup V)^*$  in P derived from  $A_N$ —i.e.,  $A_N \stackrel{*}{\Rightarrow} \alpha_N$ .

As an example,  $P_{\text{author}}$  corresponds to the AST



where we have omitted the subtree rooted at the child of  $N_{11}$ . Each node N in the AST is labeled with its corresponding symbol  $A_N$ . The holes shown are  $N_3$  (named c\_name),  $N_7$  (named c\_year) and  $N_{10}$  (named t) (there are additional holes not shown). An example of a subexpression  $\alpha_N$  is

$$\alpha_{N_5} = \sigma_{C:\texttt{c\_year} = 1948}(I:\texttt{t}~\{\phi\}).$$

**Refinements.** One sketch P' is a *refinement* of another one P if  $P \stackrel{*}{\Rightarrow} P'$ —i.e., P' can be obtained from P by filling in

the holes of P with expressions in  $\mathcal{G}$  (which may contain additional holes). Note that the nodes of P are a subset of the nodes of P'; we use  $\iota$ : nodes(P)  $\rightarrow$  nodes(P') to denote the natural injection from nodes of P to nodes of P'.

Furthermore, for a hole  $H \in \text{holes}(P)$ , we say H is *filled* with the expression  $\alpha_{\iota(H)}$ , where  $\alpha_{\iota(H)}$  is the subexpression of P' at  $\iota(H)$ . Note that conversely, given a set of expressions  $\{\alpha_H \mid H \in \text{holes}(P)\}$  such that  $A_H \stackrel{*}{\Rightarrow} \alpha_H$ , we can construct a refinement P' of P by replacing  $A_H$  with  $\alpha_H$  in P.

In our example in Figure 2, the sketch on the third line is obtained by applying the sequence of productions

$$\begin{split} I &\Rightarrow T \bowtie_{C:\mathtt{c\_new0},C:\mathtt{c\_new1}} I : \mathtt{t\_new} \\ &\Rightarrow \mathsf{authors} \bowtie_{C:\mathtt{c\_new0},C:\mathtt{c\_new1}} I : \mathtt{t\_new} \end{split}$$

to  $P_{\text{author}}$ . In particular, we then obtain the sketch

$$\begin{split} P_{\text{author}}' &= \Pi_{C:\text{c\_name}} \left( \sigma_{C:\text{c\_year}=1948} \left( t \; \{ \phi \} \right) \right) \\ & t = \text{authors} \bowtie_{C:\text{c\_new0},C:\text{c\_new1}} I: \texttt{t\_new} \end{split}$$

to  $P_{\rm author}$ —i.e.,  $P_{\rm author} \stackrel{*}{\Rightarrow} P'_{\rm author}$ . Note that  $P'_{\rm author}$  has three new holes c\_new0, c\_new1, and t\_new compared to  $P_{\rm author}$ . Furthermore, note that  $P_{\rm author}$  is obtained by filling the hole I: t in  $P_{\rm author}$  with the expression t.

**Completions.** Our algorithm assumes that the true program can be derived from the sketch provided by the user. In particular, a *completion*  $\overline{P}$  of P is a complete refinement of P—i.e., a refinement of P that has no holes; we denote the set of completions of P by  $\mathcal{P}_P$ . In this case, the subexpression  $\alpha_{\iota(H)}$  used to fill hole  $H \in \mathsf{holes}(P)$  is concrete (i.e., it has no holes). Thus, we have  $\alpha_{\iota(H)} \in \mathcal{L}(\mathcal{G}, A_H)$ .

In our example in Figure 2, the last line shows a completed sketch, which can equivalently be written

$$\begin{split} \overline{P}_{\text{author}} &= \Pi_{\text{name}} \left( \sigma_{\text{year}=1948} \left( \overline{t} \; \{ \overline{\phi} \} \right) \right) \\ & \overline{t} = \text{authors} \bowtie_{\text{aid,aid}} \text{writes} \bowtie_{\text{pid,pid}} \text{publications} \\ & \overline{\phi} = (\text{contains}(\text{name, ".*Church.*"}) \\ & \land (\text{year} \geq 1900) \land (\text{year} \leq 2020). \end{split}$$

This sketch is a completion of  $P_{\rm author}$  (and of  $P'_{\rm author}$ ). Note that  $\overline{P}_{\rm author}$  can be obtained from  $P_{\rm author}$  by filling hole t with  $\overline{t}$ , hole c\_name with name, and hole c\_year with year.

**Soft constraints.** Our language includes constraints of the form  $\alpha \sim \{\phi\}$ , where  $\alpha \in \mathcal{L}(\mathcal{G}, A)$  for some  $A \in \{Q, S, I\}$  is a table expression (i.e., a select, project, or inner-join operation), and  $\phi \in \mathcal{L}(\mathcal{G}, \Phi)$  is a *soft constraint* on tables  $t \in \mathcal{T}$ .

The DSL semantics  $[\![\cdot]\!]$  ignore these specifications. Instead, we define an additional semantics  $[\![\cdot]\!]_{\varphi}: \mathcal{L}(\mathcal{G}) \to \mathbb{R}$ ; these semantics are shown in Figure 5. In particular,  $[\![P]\!]_{\varphi}$  can be interpreted as a score encoding how well P satisfies the soft constraints; a high score means that P is satisfies the constraints very well. Furthermore, soft constraints encode

**Figure 5.** Semantics of our soft constraints. The first four lines are the semantics for table expressions t { $\phi$ }  $\in \mathcal{L}(\mathcal{G}, A)$  for  $A \in \{Q, S, I\}$ ; in this case,  $\llbracket \cdot \rrbracket_{\varphi} \in \mathbb{R}$ . The last four lines are the semantics for soft constraints  $\phi \in \mathcal{L}(\mathcal{G}, \Phi)$ ; in this case,  $\llbracket \cdot \rrbracket_{\varphi} \in \{\mathcal{T} \to \mathbb{R}\}$  is a mapping from tables to the reals. For a table  $t \in \mathcal{T}$  and a column c, t[c] is the list of values x in column c of t. For  $\llbracket u \rrbracket$ ,  $\llbracket \lesssim \rrbracket$  is  $\leq$ ,  $\llbracket \simeq \rrbracket$  is =, and  $\llbracket \gtrsim \rrbracket$  is  $\geq$ . We use  $\mathbb{I}$  to denote the indicator function, and match(x, x) to denote that string x matches regular expression x.

user expectations, so a high score means that P is a close match for what the user is expecting.

The soft constraints are on the values obtained when evaluating P. Intuitively, these semantics interpret  $\phi$  as a soft constraint on the value  $[\![\alpha]\!]$  obtained by evaluating the expression  $\alpha$  preceding  $\phi$  using the DSL semantics  $[\![\cdot]\!]$ . In particular,  $[\![\phi]\!]_{\varphi}: \mathcal{T} \to \mathbb{R}$  is a mapping from tables to real numbers. Then, the semantics for expressions of the form  $\alpha$   $\{\phi\}$  are obtained by applying  $\phi$  to  $\alpha$  to obtain  $[\![\phi]\!]_{\varphi} [\![\alpha]\!] \in \mathbb{R}$  (since  $[\![\alpha]\!] \in \mathcal{T}$ ). In addition, these semantics aggregate the values obtained from additional specifications in the expression  $\alpha$  by summing them together. In Figure 5, the first four lines show the semantics for expressions  $\alpha$   $\{\phi\}$ , and the last four lines show the semantics for soft constraints  $\phi$ .

The semantics for  $\phi=$  true always evaluates to 0; this choice is since true is the default specification, and evaluating to 0 ensures that these specifications do not affect the score  $[\![P]\!]_{\varphi}$ . For containment  $x\in c$ , the score is 1.0 if the value x is in column c of the given table t, and 0.0 otherwise. For an expression  $c\lesssim x$ , where x is a value and c is a column, the score is the fraction of values x' in column c of the table t that satisfies  $x'\leq x$ ; c and c are similar. In all these cases, c can be a integer, float, string, or regular expressions. For strings, inequalities are interpreted as lexicographical ordering. For regular expressions, we restrict to containment and approximate equality c; they are interpreted as typical regular expression matching (e.g., for containment, there must

Algorithm 1 Our interactive synthesis algorithm.

```
procedure InteractiveSynthesize(Sketch P)

\mathcal{N} \leftarrow \emptyset

while \neg \operatorname{IsComplete}(P) do

Q_P \leftarrow \operatorname{ConstructCandidateQueries}(P) \setminus \mathcal{N}

\mathcal{P} \leftarrow \operatorname{SampleCompletions}(\pi_{P,\mathcal{N}})

\hat{Q} \leftarrow \operatorname{arg\,max}_{Q \in Q_P} \widehat{\operatorname{score}}(Q; \mathcal{P})

b \leftarrow O(\hat{Q})

if b then P \leftarrow \hat{Q} else \mathcal{N} \leftarrow \mathcal{N} \cup \{\hat{Q}\}

end while

return P

end procedure
```

exist some  $x' \in t[c]$  such that  $x' \in \mathcal{L}(x)$ ). For conjunctions  $\phi \wedge \phi'$ , we add up the score based on  $\phi$  and the one based on  $\phi'$ . We make this choice since we typically use conjunctions to place multiple unrelated constraints on a table.

As an example, for the program  $\overline{P}_{\rm author}$ , there is a single soft constraint  $\overline{\phi}$ . This constraint applies to the table  $\overline{t}$ , which evaluates to the table shown in Figure 3. For this table, the soft constraint name  $\simeq 2$  evaluates to 0.33, year  $\gtrsim 1900$  to 1.0, and year  $\lesssim 2020$  to 1.0; thus,  $\|\overline{P}_{\rm author}\| = 2.33$ .

## 4 Interactive Synthesis Algorithm

Our algorithm takes as input a sketch P and returns a completion  $\overline{P}$  of P that satisfies the user intent. At a high level, it iteratively refines P, questioning the user at step to ensure that the refinement satisfies the user's intent. Our algorithm returns P once it is complete. Assuming the user answers correctly, it is guaranteed to return the true program after a number of iterations that is polynomial in the size of  $\overline{P}$ . Each iteration of our algorithm computes a question  $\hat{Q}$  (i.e., a candidate refinement of the current sketch P) in two steps:

- 1. Compute a set of candidate questions Q.
- 2. Compute the refinement  $\hat{Q} \in Q$  that maximizes the expected reduction in the size of the search space.

Then, our algorithm questions the user on the candidate refinement  $\hat{Q}$ . If the user accepts  $\hat{Q}$ , then our algorithm updates  $P \leftarrow \hat{Q}$ ; otherwise,  $\hat{Q}$  is added to a set of negative responses N that are avoided in subsequent iterations. By doing so, our algorithm maintains the invariant that the true program  $\overline{P}^*$  is a completion of the current sketch P. In particular, this invariant is preserved assuming the user provides a valid initial sketch and answers correctly at each iteration.

We describe our algorithm in more detail below. We initially ignore the impact of negative user responses, and then describe how to handle them in Section 4.3.

## 4.1 Constructing Candidate Questions

A *question* Q is a sketch that is a refinement of the current sketch P. We construct the *candidate questions*  $Q_P$  as follows:

- For each column hole C in P, we include the sketch obtained by filling C using a production  $C \Rightarrow c_i$  for some column  $c_i$ .
- For each table hole *I*, we include the sketches obtained by filling *I* using one of the sequences of productions

$$I \Rightarrow T \Rightarrow t_i \quad \text{or} \quad I \Rightarrow T \bowtie_{C,C} I \Rightarrow t_i \bowtie_{C,C} I$$

for some table  $t_i \in \mathcal{T}$ . In both cases, we implicitly use the default specification  $\Phi \Rightarrow$  true.

Note that all new holes correspond to nonterminal *I* or *C*.

Furthermore, we ensure that holes with the same name are filled using the same expressions. First, we include names on new holes in the productions used to fill holes—i.e.,

$$I \stackrel{*}{\Rightarrow} t_i \bowtie_{C:s_1,C:s_2} I:s_3.$$

Second, if we replace hole A: s, then we also replace all other holes named s using the same sequence of productions.

As an example,  $Q_{P_{\text{author}}}$  includes the sketch

$$\begin{split} P_{\text{author}}^{\prime\prime} &= \Pi_{\text{name}} \left( \sigma_{C:\text{c\_year}=1948} \left( I:\text{t} \left\{ \phi \right\} \right) \right) \\ \phi &= \left( \text{contains}(C:\text{c\_name}, \text{``.*Church.*''}) \right. \\ & \wedge \left( C:1900 \leqslant \text{c\_year} \leqslant 2020 \right). \end{split}$$

where both holes C: c\_name have been filled using name;  $Q_{P_{\text{author}}}$  also includes  $P'_{\text{author}}$ , among others.

#### 4.2 Computing Good Questions

Our algorithm selects a refinement  $\hat{Q} \in Q_P$  on which to question the user using a greedy active learning strategy [7] where we use sampling to estimate the score.

*User responses.* Our goal is to choose questions that cut off as much of the search space as possible. To do so, we need to define what part of the search space is cut off by a given user response. To this end, we represent the user as an oracle  $O:Q\to\mathbb{B}$ , where Q is the set of all possible questions and  $\mathbb{B}=\{\text{true}, \text{false}\}$ . We assume the user response O(Q) indicates whether the true program  $\overline{P}^*$  (i.e., the program the user desires) is a completion of the question Q—i.e.,  $O(Q)=\mathbb{I}[Q\stackrel{*}{\Rightarrow}\overline{P}^*]$ , where  $\mathbb{I}$  is the indicator function. For example,  $P'_{\text{author}}$  is a candidate question for  $P_{\text{author}}$ ; furthermore, assuming the true program is  $\overline{P}_{\text{author}}$ , then  $O(P'_{\text{author}})=\text{true}$  since  $\overline{P}_{\text{author}}$  is a completion of  $P'_{\text{author}}$ .

since  $\overline{P}_{\mathrm{author}}$  is a completion of  $P'_{\mathrm{author}}$ . Note that the search space is the set  $\mathcal{P}_P$  of completions of the current sketch P. If the user responds O(Q)= true, we remove completions  $\overline{P}\in\mathcal{P}_P$  such that  $\overline{P}\notin\mathcal{P}_Q$ —i.e., programs that are not completions of Q. Conversely, if O(Q)= false, we remove completions  $\overline{P}\in\mathcal{P}_P$  such that  $\overline{P}\in\mathcal{P}_Q$ .

**Scoring questions.** We score candidate questions  $Q \in Q_P$  based on the expected fraction of the search space that they cut off. To formalize this notion, we assign probabilities to completions  $\overline{P} \in \mathcal{P}_P$  of the current sketch P using the scores  $[\![\overline{P}]\!]_{\varphi}$  based on the soft constraints provided by the user. In

particular, we define a probability distribution  $\pi_P$  over the completions  $\overline{P}$  of P as follows:

$$\pi_P(\overline{P}) = \frac{1}{Z_P} e^{\llbracket \overline{P} \rrbracket_{\varphi}} \text{ where } Z_P = \sum_{\overline{p}'} e^{\llbracket \overline{P}' \in \mathcal{P}_P \rrbracket_{\varphi}}$$

where  $\mathcal{P}_P$  denotes all completions of P. In particular, completions of P with higher score have higher probability. Thus, programs that are more likely according to the user-provided soft constraints have higher probability.

Then, we score a candidate question based on the fraction of the search space  $\mathcal{P}_P$  of completions  $\overline{P}$  of P remaining based on the user's response, weighted by the probabilities  $\pi_P(\overline{P})$ . We weight according to these probabilities to focus on disambiguating among programs that are likely to be the true program  $\overline{P}^*$ . For a question  $Q \in Q_P$ , if the user accepts Q, then the fraction of the search space remaining is

$$\pi_+ = \mathbb{P}_{\overline{P} \sim \pi_P}[Q \overset{*}{\Longrightarrow} \overline{P}] = \sum_{\overline{P}} \mathbb{I}[Q \overset{*}{\Longrightarrow} \overline{P}] \cdot \pi_P(\overline{P}).$$

On the other hand, if the user rejects Q, then the fraction of the search space remaining is  $\pi_- = 1 - \pi_+$ .

However, we cannot know the fraction of the search space that is cut off by a question Q without knowing the user response O(Q). To address this issue, we interpret  $\pi_P(\overline{P})$  as the probability that the (unknown) true completion  $\overline{P}^*$  is  $\overline{P}$ . Then, we can compute the probability that the user responds true or false—in particular, the probability that the user responds true is exactly  $\pi_+$ , since  $\pi_+$  is the probability that  $\overline{P}^*$  is a refinement of Q. Analogously,  $\pi_-$  is the probability that the user responds false. Thus, the expected score is

$$score(Q; \pi_P) = \sum_{b \in \mathbb{B}} (search space pruned if O(Q) = b)$$

$$\times (probability that O(Q) = b)$$

$$= \pi_+ \cdot \pi_- + \pi_- \cdot \pi_+$$

$$= 2\pi_+ \cdot (1 - \pi_+).$$

The ideal question is one where  $\pi_+ = \pi_-$ —then, no matter how the user responds, it cuts the search space in half.

**Computing good completions.** Note that the score of a question is the expected fraction of the search space pruned if we ask the user that question. Thus, we would ideally choose the completion that maximizes the score

$$Q^* = \arg\max_{Q \in Q_P} \operatorname{score}(Q; \pi_P).$$

However, it is intractable to compute the score exactly due to the sum over completions  $\overline{P} \in \mathcal{P}_P$ . Instead, we use random samples  $\overline{P} \sim \pi_P$  to approximate the score—i.e., given a set  $\mathcal{P} \subseteq \mathcal{P}_P$  of i.i.d. samples from  $\pi_P$ , we use the approximation

$$score(Q; \pi_P) \approx \widehat{score}(Q; \mathcal{P}) = 2 \cdot \hat{\pi}_+ \cdot (1 - \hat{\pi}_+),$$

where 
$$\hat{\pi}_+ = \frac{1}{|\mathcal{P}|} \sum_{\overline{P} \in \mathcal{P}_P} \mathbb{I}[Q \xrightarrow{*} \overline{P}^*]$$
. Then, we choose question 
$$\hat{Q} = \operatorname*{arg\,max} \widehat{\operatorname{score}}(Q; \mathcal{P}).$$

**Sampling completions.** A remaining challenge is how to sample completions  $\overline{P} \sim \pi_{P_{\phi}}$ . The difficulty is that we can only compute the unnormalized probabilities  $\pi_{P_{\phi}}(\overline{P})$ . We use a standard approach to randomly sample completions—namely, the Metropolis-Hastings (MH) algorithm [6, 27]. We give a brief overview here. To use this algorithm, we must define (i) a way to sample an initial completion, and (ii) a way to sample a neighbor  $\overline{P}'$  of a completion  $\overline{P}$ .

To sample an initial completion, we independently sample an expression to fill each hole  $H \in \mathsf{holes}(P)$  of the current sketch P. In particular, for each hole H, we sample a random expression  $\alpha_H \sim \mathcal{L}(\mathcal{G}, A_H)$ , and then use  $\alpha_H$  to fill H. Once all the holes have been filled, we obtain a completion  $\overline{P}$ . Here, we think of  $\mathcal{L}(\mathcal{G}, A_H)$  as a probabilistic grammar in a standard way—i.e., by using the uniform distribution over productions for each nonterminal. Next, to sample a neighbor  $\overline{P}'$  of a completion  $\overline{P}$ , we uniformly randomly choose a single hole  $H \in \mathsf{holes}(P)$ , and replace the expression  $\alpha_H$  in  $\overline{P}$  with a newly sampled expression  $\alpha_H' \sim \mathcal{L}(\mathcal{G}, A_H)$ . This produces a modified completion  $\overline{P}'$ .

Then, MH starts by sampling an initial completion  $\overline{P}$ . Then, for a fixed number of steps, it samples a neighbor  $\overline{P}'$  of  $\overline{P}$ ; if the (unnormalized) probability  $\pi_P(\overline{P}')$  is larger than  $\pi_P(\overline{P})$  (i.e.,  $\overline{P}'$  better matches the soft constraints than  $\overline{P}$ ), then we update  $\overline{P} \leftarrow \overline{P}'$ . Otherwise, we still perform this update with some probability; this probability is computed to ensure that asymptotically,  $\overline{P}$  is a random sample from  $\pi_P$ .

## 4.3 Handling Negative Responses

So far, we have ignored the impact of user responses O(Q) = false on our search space. If a user responds O(Q), then the current sketch P does not change; instead, completions that do not match Q are removed from our search space. In particular, our algorithm keeps track of questions  $Q \in \mathcal{N}$  for which O(Q) = false. Then, our search space is actually

$$\mathcal{P}_{P,\mathcal{N}} = \{ \overline{P} \in \mathcal{P}_P \mid \forall Q \in \mathcal{N} : Q \stackrel{*}{\Rightarrow} \overline{P} \}.$$

In other words,  $\mathcal{P}_{P,N}$  omits completions that match any of the questions  $Q \in \mathcal{N}$ . Then, we modify  $\pi_P$  to take this restriction into account—i.e., we define the distribution  $\pi_{P,N}(\overline{P}) \propto \pi_P(\overline{P}) \cdot \mathbb{I}[\overline{P} \in \mathcal{P}_{P,N}]$  over completions  $\overline{P}$  of P. In other words,  $\pi_{P,N}$  is  $\pi_P$  conditioned on the event that  $\overline{P} \in \mathcal{P}_{P,N}$ .

Then, our algorithm remains the same, except we use  $\pi_{P,N}$  in place of  $\pi_P$  when sampling completions  $\mathcal{P}$ . We use rejection sampling to sample  $\overline{P} \sim \pi_{P,N}$ —i.e., we repeatedly obtain samples  $\overline{P} \sim \pi_P$  until we find one that satisfies the condition  $\mathbb{I}[\overline{P} \in \mathcal{P}_{P,N}]$ . To check this condition, we simply iterate over each  $O \in \mathcal{N}$  and check if  $\overline{P}$  is a completion of O.

#### 4.4 Overall Algorithm

Our algorithm is shown in Algorithm 1. It first constructs the set  $Q_P$  of candidate questions; in this step, it also removes questions  $Q \in \mathcal{N}$  for which the user has already responded negatively. Next, it samples an i.i.d. set of completions  $\mathcal{P}$  from  $\pi_{P,\mathcal{N}}$ . Then, it chooses the best completion  $\hat{Q}$  based on these samples. Finally, it questions the user on  $\hat{Q}$ . If  $O(\hat{Q}) = \text{true}$ , then it updates  $P \leftarrow \hat{P}$ ; otherwise, it adds  $\hat{Q}$  to the questions  $\mathcal{N}$  with negative responses. Finally, it iteratively continues this process until P is complete, at which point it returns P.

We prove that assuming the user provides a valid initial sketch and responds correctly, then our algorithm returns  $\overline{P}^*$ . We emphasize that our key contributions are our design decisions—i.e., the kind of input we require from the user. Our theoretical guarantees follow straightforwardly given these choices; wes give proofs in Appendix A. First, we prove that if our algorithm returns, then its return value is correct.

**Theorem 4.1.** Suppose (i) the initial sketch P provided by the user satisfies  $P \stackrel{*}{\Rightarrow} \overline{P}^*$ , and (ii) the user responses are  $O(Q) = \mathbb{I}[Q \stackrel{*}{\Rightarrow} \overline{P}^*]$ . Then, if our algorithm terminates, it returns  $\overline{P}^*$ .

Next, we prove that our algorithm is guaranteed to terminate. In fact, we prove that it is guaranteed to do so in a polynomial number of iterations.

**Theorem 4.2.** Our algorithm terminates after  $O((n+m) \cdot |\overline{P}^*|^2)$  iterations, where  $|\overline{P}^*|$  is the number of nodes in the AST of  $\overline{P}^*$ , n is the number of tables in the database, and m is the number of columns in the database.

This bound holds regardless of how our algorithm chooses questions. While polynomial, the number of iterations can still be large if the questions are chosen poorly—thus, in practice, choosing good questions is important.

## 5 Implementation

We have implemented our algorithm in a tool called ISQL. We briefly discuss our implementation, and give details in Appendix B. First, our implementation restricts the search space of complete programs—i.e., by restricting to inner-join operations on columns with matching keys and by imposing type constraints; by doing so, we reduce the number of iterations needed. We also modify the candidate questions for inner-join operations to fill columns involved with a join at the same time as a table—i.e.,  $t_i\bowtie_{C,C}I\stackrel{*}{\Rightarrow}t_i\bowtie_{c,c'}t_j$  and  $t_i\bowtie_{C,C}I\stackrel{*}{\Rightarrow}t_i\bowtie_{c,c'}t_j\bowtie_{C,C}I$ . Finally, we precompute the soft constraints so we do not have to evaluate database queries during the execution of our algorithm; this change slightly modifies the semantics  $\llbracket \cdot \rrbracket_{\omega}$ .

```
SELECT ??c_journal:column
   FROM (??t:table {(contains ??c_journal:column "TCS")
   AND (contains ??c_name:column ".*John.*")})
WHERE ??c_name:column = "H. V. Jagadish"
              What journals have H. V. Jagadish published in?
SELECT ??c_name:column
FROM (??t:table {(contains ??c_name:column ".*John.*")
                  AND (contains ??c_title:column ".*framework.*")})
WHERE ??c_title:column = "Other short communications"
    What is the authors of the paper titled "Other short communications"?
SELECT ??c title:column
FROM (??t:table {(contains ??c_keyword:column "Natural Language
                  AND (contains ??c_title:column ".*framework.*")})
WHERE ??c_keyword:column = "Natural Language"
     What are paper titles that contain the keywords "natural language"?
      SELECT ??c name:column
      FROM (??t:table {(contains ??c_name:column ".*John.*")
                         AND (1900 <= ??c_year:column <= 2020)
      WHERE ??c_year:column >= 2010
                    What authors published after 2010?
SELECT ??c name:column
 FROM (??t:table {(contains ??c_name:column ".*John.*")
                   AND (1900 <= ??c_year:column <= 2020)
AND (contains ??c_conference:column "PVLDB")})
 WHERE ??c_year:column = 2010 AND ??c_conference:column = "PVLDB"
                What authors published in PVLDB in 2010?
```

**Figure 6.** Examples of sketches in our first user study. The soft constraints in each sketch are shown in red. The names of the column holes indicate the name of the column used to fill it, but these are ignored by our algorithm and can be changed without affected the semantics of the sketch.

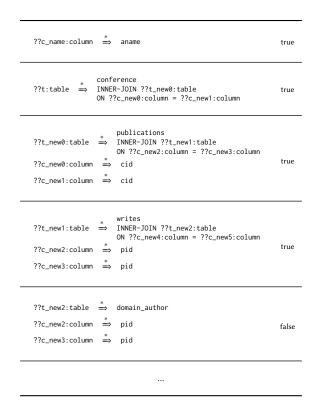
#### 6 Evaluation

We evaluate our approach on a database of academic publications [19]. This database has 16 tables; it comes with a dataset of 196 examples of SQL programs described in natural language. In our experience, writing SQL programs for this database is challenging since its schema is large and poorly documented. We aim to answer three key questions:

- Is it easy to write sketches in our language?
- Is it easy to respond to questions?
- How many iterations does our algorithm require?

We address the first two questions using a small-scale user studies where users write sketches or answer questions. We focus on evaluating feasibility of using our tool; additional experience would be needed for them to use it effectively.

**Writing sketches.** We performed a user study to evaluate whether users can write effective sketches; in particular, our goal is to evaluate whether they can write specifications for an SQL program without seeing the contents of the database.

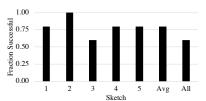


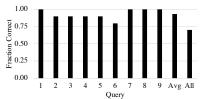
**Figure 7.** Sequence of questions for the first sketch in Figure 6. We show the question (left) and our response (right).

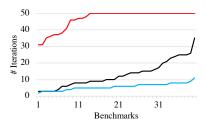
First, we selected natural language descriptions of five SQL programs, choosing each one to include a single new column compared to the previous ones for which the user needed to provide a soft constraint. For each of the five SQL programs, we wrote a corresponding "ground truth" sketch; for all of these, our tool terminated in at most ten iterations. Figure 6 shows each ground truth sketch along with a natural language description of the corresponding SQL program.

Each ground truth sketch includes at least one soft constraint for each column that appears in the sketch—i.e., each column in the project operation or in the logical formula of the select operation. Intuitively, this strategy ensures that each column hole is at least somewhat constrained, enabling our algorithm to quickly identify which columns it should include in the flat table along with the appropriate sequence of inner-join operations. Since we expect each of these columns to occur in the flat table constructed by the table hole; thus, we added these soft constraints on the table hole. For example, for the column of author names, we expected some author to be named John, so we added a soft constraint (contains ?? ".\*John.\*) to the table hole in that sketch.

Then, we recruited five users who were familiar with writing SQL programs. We showed them the natural language description for each of the five chosen SQL programs and asked them to write a sketch that captures its semantics. To train them to use our system, we asked them to read







**Figure 8.** Left: Fraction of user-written sketches that are successful at solving the task; "Avg" is the average across the five sketches, and "All" is the fraction of users that provided all successful sketches. Middle: Fraction of users that responded correctly to questions; "Avg" is the average across the nine questions, and "All" is the fraction of users that responded correctly to all questions. Right: Number of iterations used by our algorithm (black), a baseline that omits soft constraints (red), and an oracle that only asks correct questions (blue). The benchmarks are sorted by number of iterations (independently for each of the three curves). We time out benchmarks at 50 iterations or 1 hour.

through and understand a detailed example of a sketch for an SQL program. We used an SQL program for a database from another domain (in particular, an SQL program targeting a healthcare database) to ensure that there was little or no overlap in terms of the nature of the specific strategies used to write specifications, though it conveyed qualities that make a specification effective. In addition, we gave them the above guideline for writing effective sketches—i.e., that it should include at least one soft constraint for each column.

For this part of our evaluation, we did not provide users with interactive feedback for their sketch. One resulting issue was that since the users did not have access to a syntax checker, they made syntactic mistakes; we manually fixed these mistakes as long as doing so did not affect the intended semantics of the sketch. Then, we ran each of these sketches using ISQL, automatically responding to its questions based on whether they matched the ground truth sketch written by ourselves. We count a sketch as successful if it produces the correct SQL program in within two iterations of the number of iterations required when using the ground truth sketch.

In Figure 8 (left), we show the success rate of each sketch; the average success rate is 80%. Also, for three of the five users, all five sketches were successful. These results show that users can write successful sketches without significant training. Investigating the failure cases, we found they were all due to a failure to follow the guideline that a soft constraint should be included for each column used in the sketch. Since users could not interact with ISQL in this part of our study, we believe they were unable to internalize our guideline. In particular, one of our three users followed this strategy and wrote a successful sketch each of the five queries.

**Responding to questions.** Next, we evaluate whether users can respond to questions made by ISQL during sketch refinement. This study was significantly shorter, so we had a higher response rate of ten users. We show each user a description of an SQL program from the healthcare domain, along with an example of a question that should be accepted and one

that should be rejected. Then, using the last ground truth sketch in Figure 6, we asked them respond to the sequence of questions asked by ISQL. If all questions are answered correctly, there are nine total questions that the user must respond; the first five questions are shown in Figure 7. For seven of these, the correct response is "accept" (i.e., the proposed refinement was correct), and for two the response is "reject" (i.e., the proposed refinement was incorrect).

To help the user make their decision, we showed them the current sketch, the candidate refinement, and the first few rows of the relevant tables. For a column used to fill a column hole, we show them the table containing the candidate refinement. For an table hole, if the candidate refinement proposes to fill it with an inner-join operation, we show them both tables involved in that inner-join; otherwise, if it proposes filling it with a single table, we show them that table—e.g., for the second question, we show the table conference:

cid	cname	full_name	homepage
232	DAC	Design Automation Conference	
134	ATC	Autonomic and Trusted Computing	

Based on this information, it is apparent that this column is the correct one, so this question should be accepted.

For each of the nine questions, we evaluate how many users can answer them correctly, along with how many users answer *all* nine correctly; the latter measures the number of users that would obtain the correct sketch. Results are shown in Figure 8 (middle). Seven of the ten users correctly answered all questions; thus, our results suggest that users can easily respond to questions issued by ISQL.

**Iterations required.** We study the number of iterations required for our algorithm to converge to the true program  $\overline{P}^*$ . We run our algorithm using each of the sketches we wrote as input. We compare to a baseline where our algorithm is run with the soft constraints omitted. This baseline helps capture the size of the search space—in particular, it captures how our algorithm can use the soft constraints to cut down the

search space. We also compare to an oracle that only asks questions for which the user responds true—i.e., a measure of the minimum amount of work that the user must do.

We show results in Figure 8. Our algorithm terminates in fewer than 30 iterations in all but one case, and in fewer than 20 iterations in more than 75% of cases. The one case that took more than 30 iterations includes a sequence of five tables inner-joined together. In contrast, the baseline takes a large number of iterations—it times out on 27 of the 40 benchmarks (we time out after 50 iterations or 1 hour). Also, for the most part, our algorithm is only a factor of two worse than the oracle, and furthermore matches the oracle for easy benchmarks (i.e., those at the left-hand side of the plot). In contrast, the baseline is an order of magnitude worse even for easy benchmarks. Thus, our algorithm substantially cuts down the search space compared to the baseline.

Limitations. We briefly discuss a few of the limitations in our system. First, ISQL only implements a fragment of SQL—in particular, it omits aggregation operations. Extending our approach to work with these operations is straightforward; once implemented, the user can include these additional operations in the sketch (similarly to projection and selection operations), and our algorithm would fill in the holes with columns, tables, and inner-join operations.

Another limitation is the limited number of constraints that we provide. We have deliberately omitted soft constraints on words in the column names since these names may be misleading in practice. Nevertheless, it is easy to extend our system to include additional soft constraints.

Finally, our algorithm restricts to holes that are either sequences of inner-join operations (i.e., nonterminals I) or columns (i.e., nonterminals C). One particular place where a user might want to include a hole is for the constants that appear in the logical formula  $\psi$  in a select operation  $\sigma_{\psi}(t)$ . For example, in  $P_{\text{author}}$ , there is a select operation  $\sigma_{C:c\_year=1948}(...)$ . In this example, the user might not be sure how the current year is expressed, in which case they would be unable to write this sketch. However, this restriction is easy to address—the user can simply use a temporary logical formula  $\psi$  = true in the sketch to synthesize the program. Then, they can inspect the column in the synthesized program to determine the appropriate way to write  $\psi$ .

## 7 Related Work

**Program synthesis.** There has been recent interest in program synthesis. We can divide the literature along two dimensions: (i) the kind of user specifications that are used, and (ii) the search strategy. In terms of user specifications, there has been work on logical specifications [1, 24], sketches (i.e., a logical specification along with a sketch specifying the high-level structure of the code) [30–32], input-output examples [11–13, 23, 25, 34], and natural language [36].

Input-output examples and natural language specifications are especially prone to ambiguity—i.e., there are multiple programs with different semantics that satisfy the specification. In these cases, interaction has been used to resolve ambiguity [4, 5, 9, 14, 17, 19, 26, 33]. For the most part, these approaches have largely focused on obtaining additional input-output examples; as a consequence, they are typically heuristic and are unable to ensure correctness. One approach uses abstract input-output examples, which represent a potentially infinite set of concrete input-output examples; however, giving input-output examples is our setting is not practical, since we assume the user does not know the database schema. The most closely related work is [19], which interacts with the user to resolve ambiguity in user-provided natural language description. In contrast, our approach allows the user to provide more powerful specifications compared to natural language, and also has correctness guarantees.

In terms of search strategy, there has been work on deductive search [20], constraint-based search [30–32], search using version space algebras [13, 25], and (guided) enumerative search [1, 12]. We build on enumerative search in the context of syntax-guided synthesis [1], where the search is guided by a context-free grammar encoding the semantics of a domain-specific language. In contrast to existing approaches, however, our algorithm needs to sample programs rather than find a single one that satisfies the specification. On the one hand, our problem is more challenging since we need to generate many complete programs; on the other hand, it is easier since we do not need to find the true program on early iterations of our algorithm.

Next, there has been work on using AI to guide program synthesis [3, 10, 15, 18, 29], including the use of stochastic search [27, 28]. Our usage of stochastic search is different—we use it to generate samples from a distribution to use to select questions, rather than to optimize an objective. Finally, there has been work on synthesizing database queries from natural language [19, 36] and from examples [11, 21, 34, 35].

Refinement programming. There has been work on a refinement approach to writing programs [2, 8, 16, 22]. The goal is typically to develop programs that are correct by construction, where the correctness proofs are written alongside the program. In this approach, abstract components (which correspond to holes in this setting) become progressively more concrete. However, these approaches are largely manual and do not leverage interaction—i.e., the programmer has to manually both decide which abstract components to refine and implement the refinement. In contrast, our algorithm actively proposes candidate refinements, and the user only needs to accept or reject these candidates.

## 8 Conclusion

We have proposed a novel approach to interactively synthesizing database queries. Our approach is based on interactively refining a user-provided sketch with soft constraints on the expected semantics. By leveraging rich user feedback, our approach is able to provide strong correctness guarantees. We show how our algorithm can be used to synthesize queries for a database of academic publications. Future work includes implementing additional database operations and soft constraints, improving the sampling algorithm, learning a prior over programs to help guide the search, designing a user interface to facilitate interactions, and applying our approach to other domains.

#### References

- [1] Rajeev Alur, Rastislav Bodik, Garvit Juniwal, Milo MK Martin, Mukund Raghothaman, Sanjit A Seshia, Rishabh Singh, Armando Solar Lezama, Emina Torlak, and Abhishek Udupa. 2013. Syntax-guided synthesis. IEEE. 1–8.
- [2] Ralph Johan Back. 1978. On the correctness of refinement steps in program development. Department of Computer Science, University of Helsinki Helsinki, Finland.
- [3] Matej Balog, Alexander L Gaunt, Marc Brockschmidt, Sebastian Nowozin, and Daniel Tarlow. 2016. Deepcoder: Learning to write programs. In ICLR.
- [4] Osbert Bastani, Rahul Sharma, Alex Aiken, and Percy Liang. 2017. Synthesizing program input grammars. In PLDI, Vol. 52. ACM, 95– 110.
- [5] Osbert Bastani, Rahul Sharma, Alex Aiken, and Percy Liang. 2018. Active learning of points-to specifications. In Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation. ACM, 678–692.
- [6] Siddhartha Chib and Edward Greenberg. 1995. Understanding the metropolis-hastings algorithm. *The american statistician* 49, 4 (1995), 327–335.
- [7] Sanjoy Dasgupta. 2005. Analysis of a greedy active learning strategy. In Advances in neural information processing systems. 337–344.
- [8] Edsger W Dijkstra. 1968. A constructive approach to the problem of program correctness. BIT Numerical Mathematics 8, 3 (1968), 174–186.
- [9] Dana Drachsler-Cohen, Sharon Shoham, and Eran Yahav. 2017. Synthesis with abstract examples. In *International Conference on Computer Aided Verification*. 254–278.
- [10] Yu Feng, Ruben Martins, Osbert Bastani, and Isil Dillig. 2018. Program synthesis using conflict-driven learning. In Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation. ACM, 420–435.
- [11] Yu Feng, Ruben Martins, Jacob Van Geffen, Isil Dillig, and Swarat Chaudhuri. 2017. Component-based synthesis of table consolidation and transformation tasks from examples. In PLDI, Vol. 52. ACM, 422– 436.
- [12] John K. Feser, Swarat Chaudhuri, and Isil Dillig. 2015. Synthesizing Data Structure Transformations from Input-output Examples. ACM, 229–239.
- [13] Sumit Gulwani. 2011. Automating string processing in spreadsheets using input-output examples. ACM, 317–330.
- [14] Ruyi Ji, Jingjing Liang, Yingfei Xiong, Lu Zhang, and Zhenjiang Hu. 2020. Question selection for interactive program synthesis. In Proceedings of the 41st ACM SIGPLAN International Conference on Programming Language Design and Implementation. ACM, 1143–1158.
- [15] Ashwin Kalyan, Abhishek Mohta, Oleksandr Polozov, Dhruv Batra, Prateek Jain, and Sumit Gulwani. 2018. Neural-Guided Deductive

- Search for Real-Time Program Synthesis from Examples. In ICLR.
- [16] Gerwin Klein, Kevin Elphinstone, Gernot Heiser, June Andronick, David Cock, Philip Derrin, Dhammika Elkaduwe, Kai Engelhardt, Rafal Kolanski, Michael Norrish, et al. 2009. seL4: Formal verification of an OS kernel. In Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles. ACM, 207–220.
- [17] Vu Le, Daniel Perelman, Oleksandr Polozov, Mohammad Raza, Abhishek Udupa, and Sumit Gulwani. 2017. Interactive Program Synthesis. arXiv preprint arXiv:1703.03539 (2017).
- [18] Woosuk Lee, Kihong Heo, Rajeev Alur, and Mayur Naik. 2018. Accelerating search-based program synthesis using learned probabilistic models. In Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation. ACM, 436–449.
- [19] Fei Li and HV Jagadish. 2014. Constructing an interactive natural language interface for relational databases. Proceedings of the VLDB Endowment 8, 1 (2014), 73–84.
- [20] Zohar Manna and Richard Waldinger. 1986. A deductive approach to program synthesis. In Readings in artificial intelligence and software engineering. Elsevier, 3–34.
- [21] Ruben Martins, Jia Chen, Yanju Chen, Yu Feng, and Isil Dillig. [n. d.]. Trinity: An Extensible Synthesis Framework for Data Science. Proceedings of the VLDB Endowment 12, 12 ([n. d.]).
- [22] Carroll Morgan. 1994. Programming from specifications. Prentice Hall,.
- [23] Peter-Michael Osera and Steve Zdancewic. 2015. Type-and-exampledirected program synthesis. ACM, 619–630.
- [24] Nadia Polikarpova, Ivan Kuraj, and Armando Solar-Lezama. 2016. Program synthesis from polymorphic refinement types. (2016), 522–538.
- [25] Oleksandr Polozov and Sumit Gulwani. 2015. FlashMeta: a framework for inductive program synthesis. In ACM SIGPLAN Notices, Vol. 50. ACM, 107–126.
- [26] Yewen Pu, Zachery Miranda, Armando Solar-Lezama, and Leslie Pack Kaelbling. 2018. Selecting representative examples for program synthesis. In ICML.
- [27] Eric Schkufza, Rahul Sharma, and Alex Aiken. 2013. Stochastic superoptimization. In ASPLOS, Vol. 41. ACM, 305–316.
- [28] Eric Schkufza, Rahul Sharma, and Alex Aiken. 2014. Stochastic optimization of floating-point programs with tunable precision. In PLDI, Vol. 49. ACM, 53–64.
- [29] Xujie Si, Yuan Yang, Hanjun Dai, Mayur Naik, and Le Song. 2018. Learning a Meta-Solver for Syntax-Guided Program Synthesis. In ICLR.
- [30] Armando Solar-Lezama. 2009. The Sketching Approach to Program Synthesis.. In Proc. Asian Symposium on Programming Languages and Systems. Springer, 4–13.
- [31] Armando Solar-Lezama, Rodric M. Rabbah, Rastislav Bodík, and Kemal Ebcioglu. 2005. Programming by sketching for bit-streaming programs. ACM 281–294
- [32] Armando Solar-Lezama, Liviu Tancau, Rastislav Bodik, Sanjit Seshia, and Vijay Saraswat. 2006. Combinatorial Sketching for Finite Programs. ACM, 404–415.
- [33] Chenglong Wang, Alvin Cheung, and Rastislav Bodik. 2017. Interactive query synthesis from input-output examples. In Proceedings of the 2017 ACM International Conference on Management of Data. ACM, 1631– 1634.
- [34] Chenglong Wang, Alvin Cheung, and Rastislav Bodik. 2017. Synthesizing highly expressive SQL queries from input-output examples. ACM, 452–466
- [35] Navid Yaghmazadeh, Xinyu Wang, and Isil Dillig. 2018. Automated migration of hierarchical data to relational tables using programmingby-example. Proceedings of the VLDB Endowment 11, 5 (2018), 580–593.
- [36] Navid Yaghmazadeh, Yuepeng Wang, Isil Dillig, and Thomas Dillig. 2017. SQLizer: Query Synthesis from Natural Language. ACM, 63:1–63:26.

## A Proofs

**Proof of Theorem 4.1.** The key invariant maintained by our algorithm is that the true program is a completion of our current sketch P—i.e.,  $P \stackrel{*}{\Rightarrow} \overline{P}^*$ . We assume that the user provides a valid initial sketch, so this invariant holds at the beginning. Then, assuming the user answers queries correctly, if  $O(\hat{Q}) =$  true, we know that  $\overline{P}^*$  is a completion of  $\hat{Q}$ ; thus, the update  $P \leftarrow \hat{Q}$  maintains this invariant. As a consequence, we guarantee that our algorithm returns  $P = \overline{P}^*$ . In particular, at this point, P is both complete and satisfies  $P \stackrel{*}{\Rightarrow} \overline{P}^*$ , but the only completion of a complete sketch is itself; thus,  $P = \overline{P}^*$ .

**Proof of Theorem 4.2.** First, note that the number of successful iterations (i.e., iterations where  $O(\hat{Q}) = \text{true}$ ) is at most  $|\overline{P}^*|$ , since each successful iteration adds a node to the current sketch P, and we can add at most  $|\overline{P}^*|$  nodes total.

Next, the number of unsuccessful iterations before a successful iteration is at most the number of queries  $|Q_P|$ , where P is the current sketch. To bound  $|Q_P|$ , note that for each hole N in P, if  $A_N = I$ , then there are at most n ways to fill that hole (where n is the number of tables in the database), and if  $A_N = C$ , then there are at most m ways to fill that hole (where m is the number of columns in the database). Finally, there are at most  $|\overline{P}^*|$  holes in P. Thus, we have  $|Q_P| \leq (n+m) \cdot |\overline{P}^*|$ . As a consequence, the total number of iterations is  $O((n+m) \cdot |\overline{P}^*|^2)$ , as claimed.  $\square$ 

## **B** Implementation Details

We give details on our implementation.

**Restriction to joins on keys.** To constrain the search space, we restrict to inner-join operations on column pairs with matching keys (both primary key-foreign key joins and foreign key-foreign key joins). This restriction both improves performance and reduces the number of iterations needed.

*Modified candidate queries.* One challenge with our candidate queries  $Q_P$  is that if P contains an expression  $t_i \bowtie_{C,C} I$ , then one of the constructed queries is

$$t_i \bowtie_{C,C} I \stackrel{*}{\Rightarrow} t_i \bowtie_{C,C} t_j,$$

for some table  $t_j$ . However, it might be hard for a user decide if this refinement is correct without knowing in advance the values of the column holes C, C on which  $t_i$  and  $t_j$  are joined. Two other constructed queries are

$$t_i \bowtie_{C,C} I \Rightarrow t_i \bowtie_{c,C} I$$
  
 $t_i \bowtie_{C,C} I \Rightarrow t_i \bowtie_{C,c} I$ ,

where c is a column; as before, it might be hard for the user to know if either of these refinements are correct without knowing the other table I and the other column C. Thus, we require that these decisions be made together—i.e., for any

expression of the form  $t_i \bowtie_{C,C} I$ , we only consider refinements of the form

$$t_{i} \bowtie_{C,C} I \stackrel{*}{\Rightarrow} t_{i} \bowtie_{c,c'} t_{j}$$
$$t_{i} \bowtie_{C,C} I \stackrel{*}{\Rightarrow} t_{i} \bowtie_{c,c'} t_{j} \bowtie_{C,C} I$$

This modification changes Theorem 4.2—in particular, the dependence of the maximum possible number of iterations on the number of columns m and the number of tables n changes. Nevertheless, the number of iterations remains polynomial in these parameters.

*Type constraints.* So far, we have largely ignored the fact that values x in the sketch have types (i.e., integers, floats, strings, and regular expressions). Similarly, values in the database also have types (i.e., integers, floats, and strings). We can use these types to prune the space of programs. In particular, we impose these constraints on our grammar—i.e., in the expressions,  $X \in C$  and  $C \cup X$ , we require that the type of C and X be identical. A special case is when X is a regular expression, in which case we require that C have type string; in this case, we also require that C to be C. These constraints implicitly affect both our algorithm when constructing queries and when sampling completions.

Candidate queries from samples. We only consider candidate queries Q such that  $Q \stackrel{*}{\Rightarrow} \overline{P}$  for some sampled completion  $\overline{P} \in \mathcal{P}$  of the current sketch P. If Q does not satisfy this property, then it has estimated score  $\widehat{\text{score}}(Q;\mathcal{P})$ . In particular,  $\hat{\pi}_{-} = 0$  since  $\mathbb{I}[Q \stackrel{*}{\Rightarrow} \overline{P}] = 0$  for every  $\overline{P} \in \mathcal{P}$ . In other words, no sampled completion would lead to the user responding Q(Q) = true. Thus, we can safely ignore it.

**Precomputing soft constraints.** Note that the rules for scoring completions  $\overline{P}$  (i.e., computing  $[\![\overline{P}]\!]_{\varphi}$  in Figure 5 rely on the semantics  $[\![\cdot]\!]$  of subexpressions of  $\overline{P}$ . As part of our sampling procedure, we need to compute the score for a large number of completions  $\overline{P}$ . However, evaluating  $[\![\overline{P}]\!]$  can be computationally expensive.

Thus, we use an approximate approach to evaluating  $[\![\cdot]\!]_{\varphi}$ . At a high level, for a given database and user-provided sketch P, we precompute the values of the soft constraints  $\phi$  in P on every column c in the database. We assume tha columns are unique (i.e., no two tables have columns with the same name); we can achieve uniqueness by renaming columns. Then, instead of using  $[\![\cdot]\!]_{\approx}$  that only keeps track of the columns in tables. To apply a soft constraint  $\phi$ , we use the precomputed value if the corresponding column if it is in the table, and use  $-\infty$  otherwise.

In more detail, we precompute values for every *primitive* soft constraint  $\phi$  in the original sketch P—i.e.,  $\phi = x \in c$ , and  $\phi = c \ u \ x$  for some  $u \in \{ \leq, \approx, \geq \}$ . There are two cases. First,

if c is a column constant, then we precompute the value

$$\theta_{\phi} = \llbracket \phi \rrbracket_{\omega} t,$$

where t is the table that contains c. Second, if c in a hole, then for some column constant c', we let  $\phi_{c'}$  be the expression obtained by filling hole c in  $\phi$  with the production  $c \Rightarrow c'$ . Note that  $\phi_{c'}$  cannot have any more holes, since x and u are not allowed to be holes. Then, we precompute

$$\theta_{\phi_{c'}} = \llbracket \phi_{c'} \rrbracket_{\omega} t_{c'}$$

for every column c' in the database, where  $t_{c'}$  is the table containing c'. As an example, in  $P_{\rm author}$ , there are three primitive soft constraints  $C: c\_{\rm name} \simeq 2$ ,  $C: c\_{\rm year} \gtrsim 1900$ , and  $C: c\_{\rm year} \lesssim 2020$ . For the first constraint  $C: c\_{\rm name} \simeq 2$ , we precompute  $\theta_{\rm name} \simeq 2 = 0.5$ ,  $\theta_{\rm title} \simeq 2 = 0.33$ , etc. The others are similar. We also impose type constraints as described above for the modified candidate queries.

Then, we define the following approximate semantics for evaluating tables, which only keeps track of the columns in each table, and ignores the actual rows in the table:

$$\begin{split} \llbracket \Pi_{c_1,\dots,c_k}(s) \rrbracket_{\approx} &= \llbracket s \rrbracket_{\approx} \setminus \{ \llbracket c_1 \rrbracket, \dots, \llbracket c_r \rrbracket \} \\ \llbracket \sigma_{\psi}(i) \rrbracket_{\approx} &= \llbracket i \rrbracket_{\approx} \\ \llbracket t \rrbracket_{\approx} &= \alpha(t) \\ \llbracket t \bowtie_{c,c'} i \rrbracket_{\approx} &= (\llbracket t \rrbracket_{\approx} \setminus \{ \llbracket c \rrbracket \} ) \cup (\llbracket i \rrbracket_{\approx} \setminus \{ \llbracket c' \rrbracket \} ) \end{split}$$

where  $\alpha: t \mapsto (c_1, ..., c_k)$  maps a table to its columns.

Finally, we correspondingly modify the soft constraint semantics  $[\![\cdot]\!]_{\varphi}$  to obtain an approximate version  $[\![\cdot]\!]_{\varphi,\approx}$ . These semantics are identical to the semantics in Figure 5, except (i)  $[\![\cdot]\!]$  is replaced with  $[\![\cdot]\!]_{\approx}$ , and (ii) for the primitive soft constraints  $x \in c$  and  $x \ u \ c$ , we use the rules

in place of the ones in Figure 5.

This approximation actually changes the semantics of soft constraints. For example, for  $\overline{P}_{\rm author}$ , we have

$$\llbracket \overline{t} \rrbracket_{\approx} = (\text{aid, name, pid, title, year}),$$

where  $\overline{t}$  is the subexpression of  $\overline{P}_{\rm author}$ . Then, for the primitive soft constraint name  $\simeq 2$  in the soft constraints  $\overline{\phi}$  in  $\overline{P}_{\rm author}$ , we have

$$\label{eq:total_problem} [\![ \, \mathsf{name} \, \simeq \, 2 \, \, \overline{t} \, ]\!]_{\varphi,\approx} = \theta_{\mathsf{name} \simeq 2} = 0.5,$$

since name  $\in \llbracket \overline{t} \rrbracket_{\approx}$ . However, recall that  $\llbracket \text{name} \simeq 2 \ \overline{t} \rrbracket \rrbracket = 0.33$ , which shows that the semantics are different.

Intuitively, the difference is that  $[\![\cdot]\!]_{\varphi}$  is evaluated on the column observed during execution, whereas  $[\![\cdot]\!]_{\varphi,\approx}$  is evaluated on the original column. During execution, values in the column can be duplicated or deleted—e.g., due to inner-join operations with other columns or select operations. For example, in Figure 1, the value "Alan M. Turing" is duplicated

since he has two papers. Indeed, in some cases,  $[\![\cdot]\!]_{\varphi,\approx}$  may actually be more intuitive compared to  $[\![\cdot]\!]_{\varphi}$ .

**Modified scoring function.** When scoring programs  $\overline{P}$ , we also add a term based on the size of  $\overline{P}$ ; we measure the size of  $\overline{P}$  in terms of number of nodes in its AST, which we denote  $|\overline{P}|$ . In particular, we use

$$\widetilde{\text{score}}(Q; \pi_P) = \text{score}(Q; \pi_P) + \lambda \cdot |\overline{P}|,$$

where  $\lambda \in \mathbb{R}_{\geq 0}$  is a hyperparameter. We make a similar modification to  $\widehat{score}(Q; \pi_P)$ . In addition, we assign a score of  $-\infty$  to  $\overline{P}$  if there is some table operation in  $\overline{P}$  for which a column in that operation is not contained in the corresponding table. In particular, (i) for a project operation  $\Pi_{c_1,\ldots,c_k}(t)$ , we must have  $c_1,\ldots,c_k\in [\![t]\!]_{\approx}$ , (ii) for a select operation  $\sigma_{\psi}(t)$ , any column c appearing in  $\psi$  must satisfy  $c\in [\![t]\!]_{\approx}$ , and (iii) for an inner-join operation  $t\bowtie_{c,c'}t'$ , we must have  $c,c'\in [\![t]\!]_{\approx}$ .