

# Synthesizing Structurally Rich SQL Queries

## Abstract

SQL is central to relational database systems and has a large number of users. Though being expressive and powerful, SQL is hard to write since its declarative language design forces users to write structurally complex queries for many common tasks, such as moving average or argmax. As we have observed in many online Q&A forums that users often ask help from experts by presenting input-output examples, it would be desirable to build a programming-by-example (PBE) system for SQL to increase the usability of the language. However, due to the structurally rich nature of SQL, none of the existing synthesis algorithm can be satisfyingly adopted to solve SQL query synthesis task.

To address this synthesis challenge, we designed 1) a symmetry-reduced subset of SQL as the target synthesis language, and 2) a new inductive synthesis algorithm that uses abstract representations to increase the efficiency of search space traversal. Our key insight in designing the synthesis algorithm is that we can use abstract representations to reduce the original search problem and bridge the information from the input examples and output example so that we can perform searching and pruning in a more efficient way. With this new synthesis algorithm, we have built a PBE system *Scythe* and evaluated it on a set of 60 benchmarks from the online Q&A forum. The evaluation result indicates that our algorithm is efficient enough to enable us to solve most problems within seconds.

## 1. Introduction

Relational database serves an important role in modern data management, and SQL is one of the most commonly used query language. While SQL is designed declaratively and easy-to-use for many basic tasks (e.g. selecting some columns or rows from a table), it also forces users to write structurally complex queries with nested subqueries for some other tasks that cannot be easily expressed declaratively in SQL (e.g. the ArgMax problem, querying the row with greatest or least values for each group). As a result, correctly formulating such queries remains highly challenging for non-expert users: in Stack Overflow, there is a frequently asked question tag “greatest-n-per-group”<sup>a</sup> with more than 1,600 posts only for the SQL ArgMax problem, and there are even more questions posted under tags “sql”, “mysql” or “postgresql” asking for help on writing such structural complex queries.

In many online Q&A forums like Stack Overflow and Tek-Tips, users ask questions by providing input-output examples to experts: they briefly describe their question, and provide several input tables and an output table to demonstrate the query behavior. Especially for tasks involving analytical queries (e.g., argmax, moving average), users choose to demonstrate their problems using examples since pure English descriptions are often too abstract or too ambiguous for expert to correctly understand. As a result, existing natural language query systems [20] cannot solve these problems and it would be highly desirable to develop programming-by-example

(PBE) systems to enable automatic synthesis of queries from these examples to help solve the challenge in SQL formulation.

Although many successful PBE systems have been built for different domains including spreadsheet data transformation [12; 30; 14; 16] and text extraction [18], none of existing PBE systems for SQL [37; 35; 14] can satisfyingly solve the problem, since the subset of SQL used in these systems is too restrict to solve many real world problems.

The challenge of building programming-by-example systems for SQL is that synthesizing structurally rich SQL queries (e.g., queries with aggregations, subqueries) using traditional synthesis algorithms is prohibitively expensive. The three most common approaches to synthesis are enumerative, constraint solver-based, and version space algebra (VSA)-based. Enumerative search approaches [26; 25; 9] requires evaluating all intermediate enumerated results to enable memoization and equivalent-class partitioning, but SQL queries are relatively expensive to evaluate (comparing arithmetic expression) and more memory demanding to memoize (since intermediate results are tables). VSA-based approaches [27] requires that operators in the target language can be considered independently to each other to work efficiently, as that is a necessary condition for the algorithm to efficiently learn each operator independently from the output example. Finally, constraint-solver based approaches employ SMT solvers to synthesize programs, yet they do not work well in SQL since no efficient solver has been developed for the theory of relations.

Our key insight to the challenges brought up by the structural complexity of SQL language is that we built intermediate abstract representations to step-wisely reduce the problem into a simpler problem that we could easily utilize information from both input and output examples to efficiently prune and traverse each search space efficiently to find candidate queries.

First, we designed the *SynthSQL* grammar, a SQL grammar with reduced syntactical symmetries, to reduce the problem of searching SQL queries to searching *SynthSQL* queries that are consistent with input-output examples. Second, we designed a *forward abstraction algorithm* and a data structure *summary table* to divide the query space into numerous subspaces and store them compactly to represent all possible queries that can be built from user inputs. Particularly, this intermediate representation helps reduce the original problem of searching consistent queries into a simpler problem of searching for bitvectors and integer lists satisfying some property. Third, for each subspace, we designed a data structure *trace* and a *backward inference algorithm* that uses information from the output example to further simplify the search problem in a way that the output example can be used in pruning the search space. Last, we designed a *merge algorithm* that uses information stored in both forward-backward abstract representations to efficiently derive queries that are consistent with input-output examples from each subspace.

Based on the core synthesis algorithm described above, we built a programming-by-example system *Scythe* to address SQL query synthesis problem. *Scythe* also contains a heuristic ranking algorithm to rank the candidates, and an interface that allows user

<sup>a</sup> <http://stackoverflow.com/questions/tagged/greatest-n-per-group>

to modify input-output examples to refine the results when none of the top ranked queries are correct. Scythe performs efficiently in practice and is able to solve a rich set of problems in the real world.

**Contributions** This paper makes the following contributions:

- We devised a new SQL grammar SynthSQL that reduces symmetries in SQL to increase search efficiency.
- We designed new synthesis algorithms that use abstract representations to reduce the original synthesis problem to simpler search problem to overcome the challenges in synthesizing structurally complex SQL queries.
- We implemented our algorithm, and evaluated it on a set of 60 real world benchmarks collected from online Q&A forums. The evaluation result shows that our synthesis algorithm is very practical and can solve 55 real world problems within seconds.

## 2. Motivating Examples

To better understand the SQL synthesis task, we present two representative motivating example picked up from Stack Overflow to demonstrate the key challenges in synthesizing structurally rich SQL queries and our insights in addressing these challenges.

**EXAMPLE 1** Our first example is an ArgMax question. In this example, the user intends to find rows from input table whose Value field equals to the maximum Value below 8 for each User group. The user’s post<sup>b</sup> contains the following problem descriptions, and an expert provided a solution sometime after the problem was posted, and the expert used a subquery in the Join-clause to write the correct answer.

- *Description*: “I have a table (left) and I want to select the maximum value lower than the threshold 8 for each user, I want result to be as follows (right).”
- *Input example (left) and output example (right)*:

User	Phone	Value
Peter	0	1
Peter	456	2
Peter	456	3
Paul	456	7
Paul	789	10

col1	col2	col3
Peter	456	3
Paul	456	7

- *Solution provided by an expert*:

```
Select t1.User, t1.Phone, t1.Value
From (Select User, Max(Value) As MaxValue
      From input
      Where Value < 8
      Group By User) As t2 Join input As t1
Where t1.User = t2.User And t1.Value = t2.MaxValue
```

As shown by the expert, to solve the problem we first need to use an aggregation query to calculate the maximum Value under 8 for each User, and join the aggregation result with the input table on User and Value fields to find the phone number associated with these maximum Value.

**EXAMPLE 2** This is an example that involves using a subquery in the Exists-clause to write the correct answer. In this example, the user intends to find phone calls in the Call table such that its Id is not in the Phone.Book table. The user’s post<sup>c</sup> contains the following problem descriptions, and an expert provided the following solution.

<sup>b</sup> <http://stackoverflow.com/questions/33063073/postgresql-max-value-for-every-user>

<sup>c</sup> <http://stackoverflow.com/questions/367863/sql-find-records-from-one-table-which-dont-exist-in-another>

- *Description*: “How do I find out which calls were made by people whose Phone.Number is not in the Phone.Book table?”
- *Input Example*: The input example contains two tables: Phone.Book (left) and Call (right).

Id	Date	Phone
1	0945	1111
2	0950	2222
3	1045	3333

Id	Name	Phone
1	John	1111
2	Jane	2222

- *Output Example*:

Id	Date	Phone
3	1045	3333

- *Solution provided by an expert*:

```
Select *
From Call As t1
Where Not Exists
(Select *
 From Phone_Book As t2
 Where t2.Phone = t1.Phone)
```

The solution is to use Select query on the Call table, and the filter uses a Not Exists-clause to check whether there exists a rows in the table Phone.Book such that the Phone value is equal to the current Phone number.

From these examples (and more in our benchmark, as will be shown in Section 10), we made the following observations:

- Users ask questions by providing input-output example pairs (where the input example can contain multiple tables, and the output example is a single table) as well as a set of constants and aggregators in text (e.g., integer “8”, aggregator “Max” in the first example), and the desirable query typically contains no arbitrary constants or aggregators outside of user provided ones.
- Although the output example can be simple (e.g., output tables in both examples are simply subtables of the input tables), the correct queries can be very different and contain complex query structures (e.g., subquery, aggregation, Exists-clause).

These observations motivate that our synthesis system should 1) use input-output example as well as constants, aggregators provided by the user to guide the search direction and 2) search in a large search space considering rich features in SQL even when the output example is simple.

However, the search space can be very large when rich features (subquery, aggregation etc) are considered. For example, given the query skeleton (a query whose filter predicates in On, Where and Having are represented as holes) in Figure 1, even when we restrict filters that can be filled in the holes to conjunctive filter predicates with maximum length 2, ignoring Exists-clauses, the number of syntactically different queries is already greater than 2,682,720. For queries that contain multiple Join or Exists clauses, the number of syntactically different queries represented by such skeletons can be as high as  $10^{30}$ . Furthermore, due to the simplicity of the output examples, it is typically very hard to learn directly from the output example about the structure of the direct queries, which makes the search space very challenging for pruning purposes.

### 2.1 Our Approach

Our key insight in solving the aforementioned challenges is to build intermediate abstractions to reduce the original problem of searching SQL queries to some simpler search problems that can be efficiently guided by the output examples. Concretely, we solve the original problem with the following reduction steps.

**Problem Reduction 1.** Our first reduction of the search problem is that we find a core grammar of SQL language, SynthSQL, that contains less symmetry as our target language to reduce the size of

```

Select ...
From (Select User, Max(Value) As MaxValue
      From input
      Where Value < 8
      Group By User
      Having ...) As t2
Join (Select * From input Where ...) As t1
On ...

```

**Figure 1.** A query skeleton for the running example, where the ellipsis refers to holes in the skeleton (both in projection predicates and filter predicates).

our search space. SynthSQL grammar contains the following two features: 1) queries in SynthSQL are also SQL [23] queries and they can be directly executed in SQL interpreters, 2) compared to SQL, SynthSQL grammar contains additional syntactical restrictions to reduce symmetry in the search space (i.e., having fewer semantic equivalent queries).

For example, SynthSQL does not support casting a table to a value, and the following query (which is also a correct solution for EXAMPLE 1) will not be searched during the synthesis process.

```

Select t1.User, t1.Phone, t1.Value
From input As t1
Where t1.Value = (
  Select Max(Value)
  From input
  Where Value < 8 And User = t1.User)

```

However, as we will show in Section 4, benefiting from the semantics-preserving SQL rewriting rules [4], many of unsupported operators (e.g., In, All, casting table to values) can be desugared to operators in SynthSQL. In other words, the above query can be rewrote into a Join query, as shown in the solution of EXAMPLE 1. As a result, SynthSQL remains highly expressive and can capture many real world cases.

With the definition of the LightSQL, we can present our synthesis problem as below:

- (*Problem 1*). Given the grammar of SynthSQL, search for all queries that are consistent with input-output examples.

**Problem Reduction 2.** Our second reduction of the problem is to use an intermediate abstraction, query skeleton, to reduce the problem of searching syntactical components (filter predicates and projection statements) of SynthSQL queries into the problem of searching the values of these syntactical components.

The key idea of this reduction is based on the observation that many syntactically different filters have the same value in a query skeleton. For example, for the query skeleton in Figure 1, the predicate “ $t1.Value=t2.MaxValue$ ” has the same effect as “ $t1.Value \geq t2.MaxValue$ ” when they are used in the On-clause. If we can avoid this type of repetition during the search process, search efficiency will be significantly boosted. In fact, especially for ‘Exists’-predicates, the reduction rate between the number of syntactically different primitive filter predicates<sup>d</sup> and the number of different values in a query skeleton can be as high as  $1000\times$ .

Our system performs this reduction in the following two steps. First, our system enumerates all different query skeletons and divides the whole search space of SynthSQL queries into numerous subspaces, each specified by a query skeleton. Second, we encode each primitive filter predicate into a bitvector and encode each projection statement as a list of integers based on their values in the

<sup>d</sup> We use “primitive filter predicates” to refer to non-conjunctive filters, such as “ $x < y$ ”, “Exists (...)” and their negations. Since disjunction tend to be uncommonly used, SynthSQL does not provide its support and all filter predicates can be viewed as conjunctions of primitive filter predicates.

query skeleton, and thus reduce the search problem to a search problem in the domain of bitvectors and integer lists.

**Example** Figure 2 shows three examples of how this reduction is performed. For the first query skeleton,  $\bar{Q}_1 = \text{‘Select ... From input Where ...’}$  (Figure 4), we use the table input in the skeleton (we call such table the *core table* of a skeleton) to encode all primitive filter predicates and projection statement. For example, the three filter predicates “Phone < Value”, “Phone <= Value”, and “Phone < 8” all evaluate to the bitvector [10000] since only they evaluate to True on the first row of  $T_1$  and False on all other rows. The filter predicate “Value < 8” evaluates to [11110] since it evaluates to True on the first four rows but false on the last one.

In addition, each projection clause has a natural encoding to integer lists based on the indexes of the columns in the projection statement: the statement “Select User, Phone” can be mapped to an integer list [1, 2], since User is the first column of input and Phone is the second column of the table input, and another projection statement “Select User, Value, Phone” is mapped to an integer list [1, 3, 2].

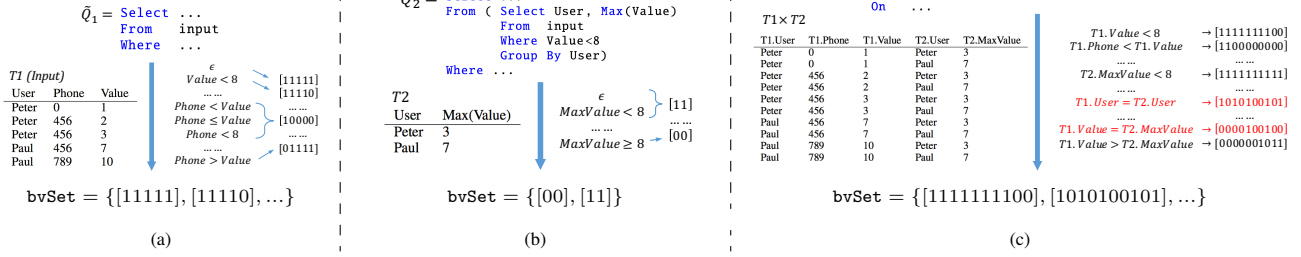
Figure 2(b,c) are two other examples demonstrating the encodings for skeletons of an aggregation query and a join query skeleton. For aggregation queries, only primitive filter predicates in the Having clause can be encoded (as the aggregation result is dependent to filter predicates inside group by and thus they are part of the skeleton) and they are encoded based on the table evaluated from the query inside. For join queries, primitive filter predicates from both subqueries and On-clause are encoded based on the Cartesian product result of the two tables in the sub queries.

Based on these encodings, we have the following correspondences between the original problem and the new problem: 1) a primitive filter predicate is mapped as a bitvector  $\beta$  (and filtering a table  $T$  with  $\beta$  is to extract rows in  $T$  whose index bit in  $\beta$  is set to 1), 2) conjunctions of filter predicates are mapped to conjunctions of bitvectors (i.e., combining a list of bitvectors using the bitwise ‘and’ operator, for example the conjunction of {[110], [011]} is the bitvector [010]), and 3) a projection statement corresponds to an integer list  $\gamma = [i_1, \dots, i_n]$  (and projecting a table with  $\gamma$  is to select columns in the table, whose indexes are in  $\gamma$ ). The original problem is then reduced as follows (Problem 1’ below is an alternative presentation of Problem 1):

- (*Problem 1’*). For each query skeleton, search for all pairs consisting of 1) a list primitive filter predicates and 2) a projection statement, such that filling the holes with these filter predicates and projection statements can generate a query that are consistent with input-output examples.
- (*Problem 2*). For each query skeleton, search for all pairs of form  $(\bar{\beta}, \gamma)$  consisting of 1) a list of bit vectors (encoded from primitive filter predicates) and 2) an integer list (encoded from a projection statement) satisfying the following checking condition: filtering the core table of the skeleton with the conjunction of all bitvectors and projecting the result with the  $\gamma$  can produce  $T_{out}$  (the output example).

With this reduction, as long as we can find all such pairs consisting of a list of bitvectors and an integer list, we can decode them to obtain filter predicates projection statements and then further generate all candidate queries within the search space.

**Example** The pair  $(\{[0000100100], [101010010101]\}, [1, 2, 3])$  for the problem in Figure 2(c) (highlighted as red) is a pair that can generate the output example in EXAMPLE 1 (by filtering the core table  $T_1 \times T_2$  with the combination result of the two bitvectors and projecting with the the filter). And the pair can be decoded to the following two predicates in SQL.



**Figure 2.** Encoding primitive filter predicates in query skeletons into bit vectors, based on their values on the core table of the skeleton.

- Filter: “ $T_1 \text{ Value} = T_2 \text{ MaxValue}$  And  $T_1 \text{ User} = T_1 \text{ User}$ ”
- Projection: “Select  $T_1 \text{ User}, T_1 \text{ Phone}, T_1 \text{ Value}$ ”

By filling the holes with these predicates, we can generate the desired query as presented in the solution of Example 1.

**Problem Reduction 3.** The third part of our problem reduction is to use information from the output example to simplify the checking condition in the reduced problem 2 to further increase the search efficiency.

The main drawback of problem 2 is that we need to enumerate and evaluate all  $(\bar{\beta}, \gamma)$  pairs on the core table of a query skeleton to select candidates, and this process can be computationally expensive due to the large number of  $(\bar{\beta}, \gamma)$  pairs and the extensive evaluation time for checking table equivalence.

Our solution to this problem is to evaluate from the output example and the core of the skeleton a set of bitvector-integer list pairs of form  $P = \{(\bar{\beta}, \gamma), \dots\}$ , containing all possible ways to obtain the output example  $T_{out}$  from the core table  $T$  by filtering and projection, (i.e., a bitvector  $\beta'$  and an integer list  $\gamma'$  can filter and project  $T$  to  $T_{out}$  if and only if the pair  $\langle \beta', \gamma' \rangle \in P$ ).

**Example** For example, between the table  $T_1$  and  $T_{out}$  in Figure 2(c), the set  $P_1 = \{(\langle [00110], [1, 2, 3] \rangle)\}$  contains all possible ways to obtain  $T_{out}$  from  $T_1$ : the only way is to filter it with the bv-filter  $[00110]$  followed by a projection on the first three columns. Similarly, such set for  $T_2$  and  $T_{out}$  is  $P_2 = \emptyset$ , since there is no way to evaluate  $T_2$  to  $T_{out}$  by only using filtering and projection.

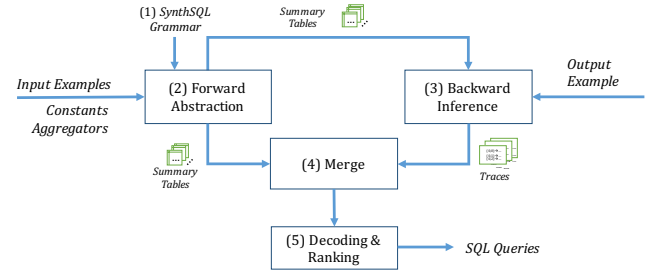
As a result, we can use the set  $P$  inferred between the core table of a skeleton and the  $T_{out}$  to further reduce the synthesis problem.

- (Problem 3). For each query skeleton, searching for all pairs of form  $(\bar{\beta}, \gamma)$  consisting of 1) a list of bit vectors (encoded from primitive filter predicates) and 2) an integer list (encoded from a projection statement) that satisfies the following checking condition: suppose  $\beta'$  is the result generated by the conjunction of  $\bar{\beta}$ , then the pair  $\langle \beta', \gamma \rangle \in P$ , inferred between the core table of the skeleton and  $T_{out}$ .

With this reduction, we can efficiently use such set  $P$  to guide searching  $(\bar{\beta}, \gamma)$  pairs in a query skeleton, for example, if  $P = \emptyset$ , we can directly prune the whole search space of the query with such skeleton, and furthermore, with a set  $P \neq \emptyset$ , we only need to search for combinations of bitvectors whose conjunction result is contained in a pair in  $P$ , and we can get  $\gamma$  for free.

## 2.2 System Pipeline

Based on these insights, we have designed algorithms and data structures to support the reduction and the final stage search process.



**Figure 3.** System overview.

As shown in Figure 3, the system synthesizes queries within SynSQL grammar, and it takes four inputs: 1) a set of input examples (denoted as  $I$ ), 2) a set of constants (optional) 3) a set of aggregators (optional), and 4) an output example (denoted as  $T_{out}$ ). The output of the systems is a list of ranked queries that are consistent with the given input-output examples.

The system first takes the inputs from the user and uses the forward enumeration algorithm to lazily enumerate all query skeletons in the space of SynSQL queries, and compactly stores them using the summary table data structure in order to perform the second reduction of our problem. Then for each query skeleton, the backward inference algorithm infers the set  $P$  between the core table of the skeleton and the output example (and compactly stores them as *traces* for the third reduction. Third, our algorithm uses a *merge* algorithm to efficiently search aforementioned  $(\bar{\beta}, \gamma)$  pairs for each skeleton using the information from the last two modules. Last, all these  $(\bar{\beta}, \gamma)$  pairs are decoded into SQL queries and returned to the user after ranking them based on a heuristic.

## 3. Preliminaries

Before introducing our approach, we present terminologies used in this paper and formally define the representation of table.

- Our system inputs includes 1) a set of input tables, 2) a set of constants, 3) a set of aggregators, and 4) an output table. We use the term *Input-output examples* to refer to the tables that are provided to the system by a user: the input example is a set of tables, denoted as  $I$ , and the output example is a single table, denoted as  $T_{out}$  (our system takes only one input-output example pair at a time).
- A *consistent query* refers to a query  $Q$  that is consistent with the provided input-output examples, i.e.,  $Q(I) = T_{out}$ .

$T ::= (\text{schema}, \text{content})$	(Table)
$\text{schema} ::= [c_1 : \tau_1, \dots, c_m : \tau_m]$	(Schema)
$\text{content} ::= [r_1, \dots, r_n]$	(Content)
$r ::= [v_1, \dots, v_m]$	(Row)
$\tau ::= \text{int} \mid \text{double} \mid \text{string} \mid \text{time} \mid \text{date}$	(Value Types)
$\beta ::= [b_1 \dots b_n]$	(Bv-Filters)
$\gamma ::= [i_1, \dots, i_k]$	(Projectors)

$T_1 \times T_2 := (\text{schema}(T_1) \uplus \text{schema}(T_2),$   
 $\quad [r_1 \uplus r'_1, \dots, r_1 \uplus r'_{n_2}, \dots, r_{n_1} \uplus r'_1, \dots, r_{n_1} \uplus r'_{n_2}])$   
 $\quad \text{where } \text{Content}(T_1) = [r_1, \dots, r_{n_1}]$   
 $\quad \text{and } \text{Content}(T_2) = [r'_1, \dots, r'_{n_2}]$   
 $\text{filter}(T, \beta) :=$   
 $\quad (\text{schema}(T), [r_i \mid \forall r_i \in \text{content}(T) \text{ and } \beta[i] = 1])$   
 $\text{proj}(T, \gamma) :=$   
 $\quad (\text{subList}(\text{schema}(T), \gamma), [\text{subList}(r, \gamma) \mid \forall r \in \text{content}(T)])$   
 $\text{deDup}(T) := (\text{schema}(T), \text{toList}(\{r_i \mid \forall r_i \in \text{content}(T)\}))$   
 $\text{subList}(L, [i_1, \dots, i_k]) := [v_i \mid v_i \in L \text{ and } i \in [i_1, \dots, i_k]]$

**Figure 4.** Table representation and operations on tables, where ‘ $\uplus$ ’ is the list-append operator,  $\oplus$  is the bitwise-‘and’ operator,  $\text{subList}(L, \gamma)$  extracts elements in  $L$  according to indexes stored in  $\gamma$ , and  $\text{toList}(S)$  converts the set  $S$  to a list.

- Bit-wise and operator  $\oplus$  is defined as below<sup>e</sup>, and is promoted to bitvectors as  $[b_1 \dots b_n] \oplus [b'_1 \dots b'_n] = [c_1 \dots c_n]$  where  $c_i = b_i \oplus b'_i$ . Furthermore, we use  $\bigoplus_{\beta \in \tilde{\beta}} \beta$  to refer to  $\bigoplus_{\beta \in \tilde{\beta}} \beta$ .

$$b_1 \oplus b_2 = \begin{cases} 1 & b_1 = 1 \text{ or } b_2 = 1 \\ 0 & \text{otherwise} \end{cases}$$

### 3.1 Relational Table

The definition of relational tables is presented in Figure 4. Relational table is defined as a pair of *schema* and *content*. The schema of a table is a list of pairs of column name and column type, the  $i$ -th element of the schema  $c_i : \tau_i$  specifies that column  $i$  has name  $c_i$ , and all values in column  $i$  of the table are of type  $\tau$ . The content of a table is a list of rows, and each row is a list of values, where the size of each row is required to be the same as the size of the schema. We use  $\text{schema}(T)$  and  $\text{content}(T)$  to refer to the schema and content of a table  $T$ , respectively.

For a table with  $n$  rows and  $m$  columns, the rows are indexed from 1 to  $n$  and columns from 1 to  $m$ , and we use coordinate  $(i, j)$  to refer to the cell in the  $i$ -th row and  $j$ -th column of the table.

Besides, we also define bv-filters ( $\beta$ ) and projectors ( $\gamma$ ) on a table: the former is a bitvector while the latter is a list of integers. The function  $\text{filter}(T, \beta)$  applies a bv-filter  $\beta = [b_1 \dots b_n]$  to a table  $T$ , in a way that row  $i$  is kept if and only if  $b_i$  equals to 1, and the function  $\text{proj}(T, \gamma)$  refers to projecting table  $T$  with a projector  $\gamma$ , in a way that column  $i$  is kept if and only if  $i \in \gamma$ . Other operators on tables include Cartesian product ‘ $\times$ ’, and  $\text{deDup}$ , which removes duplicate rows in a table.

Since we use bag semantics instead of set semantics (following SQL conventions), duplicate rows are allowed in a table and two tables are equivalent if and only if their contents are contained in each other, and the same rows in two tables have the same multiplicity (although row order in the table does not necessary need to be the same for equivalence).

<sup>e</sup> We use the notation ‘ $\oplus$ ’ to represent bit-and for bv-filters in order to distinguish it from logical ‘and’ operators used in other places of the paper.

## 4. SynthSQL Grammar

### 4.1 Syntax

The syntax of SynthSQL is presented in Figure 5. As presented in Section 2, SynthSQL is a subset of the SQL grammar with reduced symmetry while remains expressive enough to capture many real world problems. Many SQL queries using operators outside of SynthSQL grammar can be semantically-equivalently reduced to queries in SynthSQL using SQL rewriting rules.

The top level construct of a SynthSQL query  $Q$  is either an intermediate query  $q$  or a projection query built from  $q$ . An intermediate query  $q$  can be one of 1) named table reference, 2) filter query, 3) join query, 4) aggregation query and 5) rename query (where the table name and all columns are required to be specified). For SynthSQL, we use the term *primitive filter predicates* refers to filters of the form “ $v \text{ binop } v$ ” or “(not) exists  $q_e$ ”, and use the term *conjunction of filter predicates* to refer to filters of the form “ $f \text{ and } f$ ”.

Besides the above definitions, SynthSQL has the following additional restrictions on filter formulation.

- For primitive join predicates of the form  $v_1 \text{ binop } v_2$  in a query construct “ $q_1 \text{ join } q_2 \text{ on } f$ ”,  $v_1$  must refer to a column in  $q_1$  while  $v_2$  must refer to a column in  $q_2$ .
- For primitive filter predicates of form  $v_1 \text{ binop } v_2$  in a Having-clause of an aggregation query,  $v_1$  can only refer to an aggregation column  $\alpha(c'_i)$ .
- At least one of filter predicates of form  $v_1 \text{ binop } v_2$  in a nested-subquery in Exists-clause have the property that  $v_1$  refers a correlated column in outer level query.

$Q ::=$	$q$	
	$\mid$ <b>select</b> (distinct) $c_1, \dots, c_n$	(Projection)
	$\mid$ <b>from</b> $q$	
$q ::=$	$T$	(Table Reference)
	$\mid$ <b>select</b> * <b>from</b> $T$ <b>where</b> $f$	(Filtering)
	$\mid$ $q_1 \text{ join } q_2 \text{ on } f$	(Join)
	$\mid$ <b>select</b> $c_1, \dots, c_n, \alpha(c'_1), \dots, \alpha(c'_m)$	(Aggregation)
	$\mid$ <b>from</b> $q$	
	$\mid$ <b>group by</b> $c_1, \dots, c_n$	
	$\mid$ <b>having</b> $f$	
	$\mid$ $q \text{ as name}[c_1, \dots, c_n]$	(Rename)
$f ::=$	$f \text{ and } f$	
	$\mid$ $v \text{ binop } v$	
	$\mid$ (not) exists $q_e$	
$q_e ::=$	<b>select</b> * <b>from</b> $T$ <b>where</b> $f$	
$v ::=$	$c$	
	$\mid$ constant	
$\alpha ::=$	max, min, sum, avg, count, concat	
$\text{binop} ::=$	=, >, <, >=, <=, <>	

**Figure 5.** SynthSQL syntax, where meta-variable  $c$  ranges over column names,  $T$  ranges over table names from user inputs,  $\text{name}$  ranges over newly assigned table names.

Although SynthSQL grammar contains only a restricted set operators in SQL, it remains very expressive, benefiting from semantics-preserving query rewriting rules [4]. For example, although SynthSQL does not directly support ‘In’, ‘All’ operators in filter, or casting table to values, SQL queries written with these syntax constructs can be transformed to queries in SynthSQL, and thus are still supported. For example, subquery-unnesting rewriting rule [10] is able to convert table to value casting queries to queries that have subqueries in the Join-clause.

These restrictions on the syntax of SynthSQL reduce *symmetry* in the search space, i.e., there are fewer queries that are semantically equivalent but with different syntactical forms in the search

space comparing to SQL, which reduce the complexity of the process of search space traversal.

**Limitations** There are also queries in SQL that cannot be semantics-equivalently transformed to SynthSQL: e.g., Union, Outer Join, disjunction filter predicates and arithmetic expressions. However, as we can see in evaluation, these operator occurs less common in on-line question posts and our language support is expressive enough to capture most real-world use cases.

## 4.2 Semantics

Since SynthSQL queries are also native SQL queries, they can be interpreted with SQL interpreters. Particularly, we use bag-semantics of SQL in this paper, which is formally defined by Negri et al. [23], and commonly used in many modern database systems including MySQL, SQL Server, etc.

## 5. Synthesis Algorithm

We first present our synthesis algorithm. The algorithm takes four inputs: 1) a set of input example tables  $I$ , 2) an output example table  $T_{out}$ , 3) a set of constants  $consts$ , and 4) a set of aggregators. The algorithm synthesizes a list of ranked queries in SynthSQL grammar that are consistent with the input-output example. The synthesis algorithm contains the following steps.

1. (*Forward Abstraction*) The algorithm iterates over the complexity of the skeletons, starting from 1, and then enumerates all possible query skeletons within the complexity level that can be built from the input tables, constants and aggregators provided by user (line 4). (Since there are infinitely many queries consistent with input-output examples, we use the variable *complexity* to decide when to communicate the results to the user.)
2. (*Backward Inference*) For each summary table  $\rho$  generated from forward enumeration, infer the trace  $\sigma$  between  $\text{coreTable}(\rho)$  (i.e., the core table of  $\rho$ ) and the output example  $T_{out}$  (line 7).
3. (*Merge*) Use the trace  $\sigma$  to guide the search in the subspace specified by the summary table  $\rho$  to obtain a set of filter-projector pairs in  $\rho$ , which are encodings of all consistent queries in the subspace specified by  $\rho$  (line 8).
4. (*Decode*) Decode each  $(\bar{\beta}, \gamma)$  pair on the summary table  $\rho$  to obtain a list of candidate queries that are consistent with the given input-output example (lines 9,10).
5. (*Ranking and User Interaction*) Rank all candidate queries synthesized in this complexity level (if not empty) and present them to the user (lines 11,13). If the user accepts the answer, the algorithm terminates, otherwise continue searching for more complex queries in the search space (lines 13). (A user can also provide new examples to the system to restart the synthesis process on the new inputs, when none of these synthesized queries is correct.)

In the following sections, we formally define the intermediate data structures and elaborate on the synthesis modules used in the algorithm.

## 6. Forward Abstraction

The first module of our algorithm is forward abstraction, where we enumerate all query skeletons in the search space and construct summary tables from them to divide the space of SynthSQL into subspaces and reduce the problem of searching queries within the subspace into a simpler problem of searching for bv-filters and projectors on a table.

```

1 Synthesis( $I, T_{out}, consts, aggregators$ )
2    $complexity \leftarrow 1$ 
3   while True:
4      $S \leftarrow \text{forwardEnum}(I, consts, aggregators, complexity)$ 
5      $candidates \leftarrow \emptyset$ 
6     foreach  $\rho \in S$ :
7        $\sigma \leftarrow \text{backwardInference}(\text{coreTable}(\rho), T_{out})$ 
8        $R \leftarrow \text{merge}(\rho, \sigma)$ 
9       foreach  $(\bar{\beta}, \gamma) \in R$ :
10         $candidates \leftarrow candidates \cup \text{decode}(\rho, \bar{\beta}, \gamma)$ 
11       $candidates \leftarrow \text{rank}(candidates)$ 
12      if  $candidates \neq \emptyset$ :
13         $feedback \leftarrow \text{communicate}(candidates)$ 
14        if  $feedback = \text{'Accepted'}$ :
15          break;
16       $complexity \leftarrow complexity + 1$ 

```

Figure 6. Synthesis Algorithm.

$$\begin{aligned}
\rho &::= \text{ST}_p(T, B, Q) && (\text{Primitive Summary Table}) \\
&\quad | \text{ST}_c(\rho_1, \rho_2, B) && (\text{Compound Summary Table}) \\
B &::= \{\beta_1, \dots, \beta_k\} && (\text{Bv-filter set})
\end{aligned}$$

$$\begin{aligned}
\text{coreTable}(\text{ST}_p(T, B, Q)) &:= T \\
\text{coreTable}(\text{ST}_c(\rho_1, \rho_2, B)) &:= \text{coreTable}(\rho_1) \times \text{coreTable}(\rho_2) \\
\text{projSet}(\rho) &:= \{[i_1, \dots, i_k] \mid 1 \leq k, \\
&\quad 1 \leq i_k \leq \text{colNum}(\text{coreTable}(T))\} \\
\text{bvSet}(\text{ST}_p(T, B, Q)) &:= B \\
\text{bvSet}(\text{ST}_c(\rho_1, \rho_2, B)) &:= \\
&\quad B \cup \{\text{pullUpL}(\beta, T_1, T_2) \mid \beta \in \text{bvSet}(\rho_1)\} \\
&\quad \cup \{\text{pullUpR}(\beta, T_1, T_2) \mid \beta \in \text{bvSet}(\rho_2)\} \\
&\quad \text{where } T_i = \text{coreTable}(\rho_i)_{i=1,2} \\
\text{pullUpL}([b_1 \dots b_n], T_1, T_2) &:= [b_{11} \dots b_{1m} \dots b_{n1} \dots b_{nm}], \\
&\quad \text{where } b_{ij} = b_{i, (i \in [1, n], j \in [1, m])} \text{ and } m = \text{rowNum}(T_2) \\
\text{pullUpR}([b_1 \dots b_m], T_1, T_2) &:= [b_{11} \dots b_{1m} \dots b_{n1} \dots b_{nm}], \\
&\quad \text{where } b_{ij} = b_{j, (i \in [1, n], j \in [1, m])} \text{ and } n = \text{rowNum}(T_1)
\end{aligned}$$

Figure 7. Definition of summary tables and operations on them, where *maxLen* stands for the length limit for conjunctive filters (predefined in the system).

### 6.1 Data Structure: Summary Table

Summary tables are defined in Figure 7 and they are the data structures used for representing these subspaces.

A primitive summary table  $\rho = \text{ST}_p(T, B, Q)$  represents the subspace specified by a query skeleton of form “Select...From Q Where...”, where  $Q$  is either an aggregation query or a reference to an input table, and it consists of three components: 1)  $T$  is the table evaluated from the query  $Q$ , 2)  $B$  is a bv-filter set, containing values of all primitive filter predicates that can be filled into the hole in the query skeleton and 3) the query  $Q$  inside the skeleton.

A compound summary table  $\rho = \text{ST}_c(\rho_1, \rho_2, B)$  captures all queries with the skeleton “Select...From  $\tilde{Q}_1$  Join  $\tilde{Q}_2$  On...”, where  $\tilde{Q}_1$  is the skeleton represented by  $\rho_1$ ,  $\tilde{Q}_2$  is the skeleton represented by  $\rho_2$  (where the projection clauses of both sub-summary tables are “Select \*”, according to the SynthSQL grammar). The bv-filter set  $B$  stores values of all primitive join predicates that can be filled in the On clause, evaluated on  $\text{coreTable}(\rho)$  (which is defined as Cartesian product of core tables of sub summary tables).

### 6.2 Enumeration Algorithm

The forward enumeration starts from input tables, constant set and aggregator set provided by the user, and a complexity control



```

1 forwardEnum( $I, consts, aggregators, k$ )
2   if  $k = 1$ :
3     return enumFilteringQuery( $I, consts, maxLen$ )
4    $S \leftarrow$  forwardEnum( $I, consts, aggregators, k - 1$ )
5    $S_A \leftarrow$  enumAggrQuery( $S, aggregators$ )
6    $S_J \leftarrow \emptyset$ 
7   forall  $(i, j)$  s.t.  $i + j = k$ :
8      $S_1 \leftarrow$  forwardEnum( $I, consts, aggregators, i$ ),
9      $S_2 \leftarrow$  forwardEnum( $I, consts, aggregators, j$ ),
10     $S_J \leftarrow S_J \cup$  enumJoinQuery( $S_1, S_2$ )
11  return  $S_A \cup S_J$ 

```

**Figure 8.** Enumeration algorithm for complexity level  $k$ .

integer, and it returns all query summary tables of complexity level  $k$  tables that can be constructed from these inputs.

The input of the forward enumeration algorithm (Algorithm 8) includes 1) a set of input example  $I$ , 2) a set of constants, 3) a set of aggregators and 4) an integer  $k$ , specifying the complexity level of the summary tables to be enumerated. And the enumeration algorithm returns all summary tables of complexity level  $k$ : if  $k$  equals 0, the result would be all summary tables built from input examples (line 3), otherwise, the algorithm builds summary tables from aggregation queries or join queries by recursively invoking the enumeration function (lines 4-11).

The three enumeration modules used in the algorithm are presented below (Figure 9).

- **enumFilteringQuery** enumerates and encodes primitive filters for all tables provided by the user, including both binary comparison predicates (i.e.,  $v_1 \text{ op } v_2$ ) and exists-filters. Binary comparison predicates are generated by **enumPrimFilters**, where  $v_1$  is a column name, and  $v_2$  can be either a column name or a constant with same type as  $v_1$  from *consts*. For Exists-predicates, to avoid re-enumerating these clauses, instead of directly enumerating them, we first construct a sketch of the clause using the function **enumExistsSkeleton**, and instantiate them to filter predicates by substituting columns in the table with holes in a skeleton query to form a correlated column reference.
- **enumAggrQuery** enumerates queries with aggregations with given summary tables and aggregators. It first instantiates all conjunctive bv-filters represented by each summary table (lines 3,4), then evaluate each bv-filter on the core table to obtain a table  $T$ . The algorithm then enumerates and evaluates all aggregation queries that can be build from  $T$ , and build a primitive summary table to store the enumeration result.
- **enumJoinQuery** takes two sets of summary tables as input and build compound summary tables from them. For every  $\rho_1, \rho_2$  such that each  $\rho$  one summary table set, the algorithm enumerates all primitive join predicates between these two sub-summary tables and encodes them on  $\text{coreTable}(\rho_1) \times \text{coreTable}(\rho_2)$  as bit vectors to build a summary table.

At the end of the forward enumeration process, a set of summary tables is generated, and with these abstract representations, we can perform the following reduction on the original search problem of consistent query within the search space of SynthSQL.

### 6.3 Problem Reduction

Before introducing how the original problem can be reduced with summary tables, we first introduce one key operator on a summary table  $\rho$ : **bvSet**( $\rho$ ). The function **bvSet**( $\rho$ ) computes the values of all primitive filter predicates that can be filled in the query skeleton: for primitive summary tables (evaluated on  $\text{coreTable}(\text{row})$ ). For primitive summary tables, the result is the set  $B$  stored in the summary table, since  $B$  already contains the values of all

filter predicates that can be filled in the “where...” clause of the skeleton. While for compound summary tables, since the set  $B$  only stores the values of join predicates encoded on  $\text{coreTable}(\rho)$  but not filter predicates inside sub summary tables, we need to generate the encoding of these filter predicates from the two sub summary tables.

We solve this problem by using the two “pullUp” operators: since these filter predicates are already encoded in sub summary tables (evaluated on  $\text{coreTable}(\rho_1)$  and  $\text{coreTable}(\rho_2)$ ), we can reuse these encodings to generate their corresponding encodings on  $\text{coreTable}(\rho)$ . The correctness of these “pullUp” operators are based on the following two equations, which indicates that as long as  $\beta$  is an predicate encoded on one of the sub summary table, the pulling up results (e.g.,  $\text{pullUpL}(\beta, T_1, T_2)$ ) has the same effect as  $\beta$  in the query skeleton, and thus it forms a correct encoding on  $\text{coreTable}(\rho)$ .

$$\begin{aligned} \text{filter}(T_1 \times T_2, \text{pullUpL}(\beta, T_1, T_2)) &= \text{filter}(T_1, \beta) \times T_2 \\ \text{filter}(T_1 \times T_2, \text{pullUpR}(\beta, T_1, T_2)) &= T_1 \times \text{filter}(T_2, \beta) \end{aligned}$$

In addition, the function **projSet**( $\rho$ ) enumerates all combinations of integers between 1 and  $\text{colNum}(\text{coreTable}(\rho))$ , representing all projectors on the core table of  $\rho$ .

**Problem 2.** Given a summary table  $\rho$ , the problem of traversing the subspace to find all candidate queries that are consistent with input-output examples is reduced to the problem of finding all  $(\bar{\beta}, \gamma)$  pairs such that:

$$\begin{aligned} \bar{\beta} &\subseteq \text{bvSet}(\rho), \\ \gamma &\in \text{projSet}(\rho), \\ \text{proj}(\text{filter}(\text{coreTable}(\rho), \bigoplus \bar{\beta}), \gamma) &= T_{out}. \end{aligned}$$

And if we can find all these  $(\bar{\beta}, \gamma)$  pairs, we can 1) decode each primitive filter and the projection predicate to a syntactical predicate, and 2) fill the holes in the skeleton with these decoded predicates to generate queries that are consistent with input-output examples. The correctness of this reduction is presented by the following lemma.

**Lemma 1.** (Complete of Summary Tables) Suppose  $Q$  is a query expressible in SynthSQL within inputs  $(I, consts, aggregators)$ , then there exists a summary table  $\rho$  from the enumeration result **forwardEnum**( $I, consts, aggregators, k$ ), for some integer  $k$ , s.t. exists  $\bar{\beta} \subseteq \text{bvSet}(\rho)$ ,  $\gamma \in \text{projSet}(\rho)$  and the equation  $\text{proj}(\text{filter}(T, \bigoplus \bar{\beta}), \gamma) = \text{eval}(Q)$  holds (where  $T = \text{baseTable}(\rho)$ ).

*Proof.* We first prove that for each constructor of  $q$  in SynthSQL grammar, there exists a summary table  $\rho$  and a set of bv-filter  $\bar{\beta} \in \text{bvSet}(\rho)$  such that  $\text{filter}(T, \bigoplus \bar{\beta}) = \text{eval}(Q)$ . This proof can be done by induction on the constructor of  $q$  in SynthSQL grammar:

- If  $q := T_0$  (the first constructor of  $q$ ) where  $T_0 \in I$  (the input table set), it is obvious that  $\text{ST}_p(T_0, B, T_0)$  is in the result set of **enumFilteringQuery**( $I, const, holes$ ) and the filters  $\bar{\beta} = \{[11...1]\}$  satisfy the property (which is the default empty filter added into the set  $B$ ).
- If  $q := \text{select } * \text{ from } T_0 \text{ where } f$ . Then consider the summary table  $\rho = \text{ST}_p(T_0, B, T_0)$ . If the filter  $f$  is a primitive filter  $f_0$ , then there exists  $\beta_0 \in B$  such that  $\text{eval}(T_0, f_0) = \beta_0$  as the encoding filters will encode all primitive filters into bv-filters and we have  $\rho, \bar{\beta} = \{\beta_0\}$  satisfying the result. If  $f$  is a compound primitive filter, i.e.  $f = f_1 \text{ and } \dots \text{ and } f_n$ , we can find the bv-filter set  $\bar{\beta} = \{\beta_1, \dots, \beta_n\}$  satisfying the property, where  $\beta_i$  is the encoding result of  $f_i$  on  $T_0$ .
- If  $q$  is an aggregation query (with nested subquery  $q_0$ ), by induction, there exists a summary table  $\rho_0$  and a set of bv-filters  $\bar{\beta}_0 \in \text{bvSet}(\rho_0)$  s.t.  $\text{filter}(T_0, \bigoplus \bar{\beta}_0) = \text{eval}(q_0)$ .

Since all group by fields are enumerated in `enumAggrQuery` on `eval(q0)`, there exists a summary table  $ST_p(T_1, B, Q)$  such that  $T_1$  is the aggregation result of  $T_0$  on the group by fields. Similar to the second case, each filter in **having** clause has a corresponding encoding and the property is satisfied.

- If  $q$  is a join query of the form  $q = q_1 \text{ join } q_2 \text{ where } f$ , by induction, there exists  $\rho_1, \rho_2, \beta_1, \beta_2$ , such that  $\text{eval}(q_i) = \text{filter}(T_i, \oplus \beta_i)$ . Since a compound summary table  $\rho = (\rho_1, \rho_2, B)$  is generated in `enumCompoundTable` algorithm and  $\text{eval}("q_1 \text{ join } q_2") = \text{filter}(T_1 \times T_2, \oplus \beta'_1 \oplus \oplus \beta'_2)$ , where  $\beta'_1 = \{\text{pullUpL}(\beta_1, T_1, T_2) \mid \beta_1 \in \beta_1\}$  ( $\text{pullUpR}$  for  $\beta'_2$ ). Furthermore, since the filter  $f$  has a set of corresponding encodings  $\beta$  in  $B$ , we have the property that  $\text{eval}(q) = \text{filter}(\text{eval}("q_1 \text{ join } q_2"), \oplus \beta)$  and then we can further expand the expression into  $\text{eval}(q) = \text{filter}(\text{filter}(T_1 \times T_2, \oplus \beta'_1 \oplus \oplus \beta'_2), \oplus \beta) = \text{filter}(T_1 \times T_2, \oplus \beta'_1 \oplus \oplus \beta'_2 \oplus \oplus \beta)$ . As a result, the property also holds for this case.

With all three cases proved, the lemma is proved.  $\square$

**Challenges in Problem 2.** Although Problem 2 is a simpler search problem comparing to the original problem, directly enumerating all these pairs remains challenging, since 1) the number of all  $(\beta, \gamma)$  pairs we need to visit in the search space is large, and 2) the computation we need for checking the whether a pair  $(\beta, \gamma)$  is a candidate relatively expensive (as we need to evaluate all pairs on the table  $T_1$  and compare the result against  $T_{out}$ ). These challenges exist mainly because  $T_{out}$  is not effectively used in the search process. In the next section, we introduce how  $T_{out}$  can be effectively used to further reduce the complexity of the search problem.

## 7. Backward Inference

To overcome the challenges in Problem 2, our key insight is to build a simpler checking condition in a backward manner, i.e., using the output example, substitute the original one to enable more efficient pruning and checking. Concretely, given a summary table  $\rho$ , the backward inference algorithm contains two main components:

- For each summary table  $\rho$ , we build the trace  $\sigma$  (as will be shown later) between `coreTable`( $\rho$ ) and the output example  $T_{out}$  to capture how `coreTable`( $\rho$ ) can be filtered and projected to generate  $T_{out}$ .
- From each trace  $\sigma$ , we infer a set  $P = \{(\beta, \gamma), \dots\}$  such that a pair  $(\beta, \gamma)$  that satisfies the original checking condition if and only if  $(\oplus \beta, \gamma) \in P$ . And we use this condition to further reduce the original search problem.

### 7.1 Data Structure: Trace

Trace is the data structure built between two tables  $T_1$  and  $T_2$ , and it compactly represents all possible ways to obtain  $T_2$  from  $T_1$  with `filter` and `proj` operators defined in Figure 4.

The definition of trace in presented is Figure 10: a trace  $\sigma$  built from two tables  $T_1$  and  $T_2$  is a set of mappings between cell coordinates in  $T_2$  and cell coordinates in  $T_1$ . Each mapping  $\omega \mapsto \Omega$  maps a cell  $\omega$  in  $T_2$  to a set of cells  $\{\omega'_1, \dots, \omega'_m\}$  in  $T_1$ , where  $\omega = (i, j)$  represent the coordinate of a cell in  $T_2$  ( $i$  for its row index and  $j$  for its column index) and  $\omega'_1, \dots, \omega'_m$  are cells in  $T_1$  such that values in these cells are equal to the value in  $\omega$ .

### 7.2 Inference Algorithm

The algorithm to build a mapping instance  $M$  from two tables  $T_1$  and  $T_2$  is presented in Figure 11. The algorithm iterates over all cells in  $T_2$ , and for each cell  $\omega$ , it searches all cells  $\omega'_1, \dots, \omega'_k$  in  $T_1$

such that their values are equal to the value in  $\omega$  and add a mapping  $\omega \mapsto \{\omega'_1, \dots, \omega'_k\}$  into the mapping instance.

### 7.3 Problem Reduction

The key operation on a trace  $\sigma$  is `fullInstantiate`( $\sigma$ ): the operator extract a set of filter-projector pairs<sup>f</sup>  $P$  from  $\sigma$ , where each pair  $\langle \beta, \gamma \rangle = \langle [b_1 \dots b_n], [i_1, \dots, i_m] \rangle \in P$  contains a bv-filter and a projector on  $T_1$ , and it satisfies the equation  $\text{proj}(\text{filter}(T_1, \beta), \gamma) = T_2$ . Furthermore, the set  $P = \text{fullInstantiate}(\sigma)$  is *complete* if it contains all  $\langle \beta, \gamma \rangle$  pairs satisfying the equation above, as presented below in Lemma 2.

**Lemma 2.** (Complete Representation of Mapping Instances) *Given two tables  $T_1$  and  $T_2$ , suppose  $M$  is the mapping instance built from  $T_1$  and  $T_2$  using the backward inference algorithm defined in Figure 11,  $\beta$  and  $\gamma$  are two bitvectors such that  $|\beta| = \text{rowNum}(T_1)$ ,  $|\gamma| = \text{colNum}(T_1)$ , then  $\text{proj}(\text{filter}(T_1, \beta), \gamma) = T_2$  if and only if  $\langle \beta, \gamma \rangle \in \text{fullInstantiate}(M)$ .*

*Proof.* The proof of the lemma contains the following two parts:

- “ $\Rightarrow$ ”: Suppose  $\beta = [i_1, \dots, i_n]$  and  $\gamma = [j_1, \dots, j_m]$  is an bv-filter and projector pair such that  $\text{proj}(\text{filter}(T_1, \beta), \gamma) = T_2$ . By definition of `proj` and `filter`, we know that for each  $k \in [1, \dots, n]$ ,  $l \in [1, \dots, m]$ ,  $\text{val}((i_k, j_l), T_1) = \text{val}((k, l), T_2)$ . Thus, for all pairs  $(i_k, j_l)$ ,  $(i_k, j_l) \in \text{Img}(\sigma_{T_1, T_2}((i, j)))$  and  $\beta \gamma$  can be derived from the backward inference algorithm.
- “ $\Leftarrow$ ”: Suppose  $\beta = [i_1, \dots, i_n]$  and  $\gamma = [j_1, \dots, j_m]$  are generated from backward inference algorithm. By definition, for each  $k \in [1, \dots, n]$ ,  $l \in [1, \dots, m]$ ,  $\text{val}((i_k, j_l), T_1) = \text{val}((k, l), T_2)$ , and thus filtering  $T_1$  using  $\beta$  and projecting the result with  $\gamma$  results in  $T_2$ .

With both directions of “if” and “only if” proved, the lemma is proved.  $\square$

**Problem 3.** Given a summary table  $\rho$ , and the trace  $\sigma$  inferred between `coreTable`( $\rho$ ) and  $T_{out}$ , Problem 2 can be further reduced to the problem of finding all  $(\beta, \gamma)$  pairs such that:

$$\begin{aligned} \beta &\subseteq \text{bvSet}(\rho), \\ \gamma &\in \text{projSet}(\rho), \\ (\oplus \beta, \gamma) &\in \text{fullInstantiate}(\sigma). \end{aligned}$$

**Challenges in problem 3.** Although the second reduction can efficiently help guide the search process of  $(\beta, \gamma)$  in a space summary  $\rho$ , directly instantiating all  $\langle \beta, \gamma \rangle$  from  $\rho$  from  $\sigma$  when the core table of the  $\rho$  is the Cartesian product of two (or more) tables (e.g.,  $\rho_3$  in Figure 2(c)) can be prohibitively expensive, as the set  $P = \text{fullInstantiate}(\sigma)$  can contain an extremely large number of elements in some cases.

As an illustration, consider the following table fraction formed by the Cartesian product of an input example with a table  $T'$  with  $r$  rows. (the input and output table are from Example 1 of Section 2)

$r$ rows	...	...	...	...
	Peter	456	3	...
	...	...	...	...
$r$ rows	Peter	456	3	...
	Paul	456	7	...
	...	...	...	...
	Paul	456	7	...
	...	...	...	...

The problems stems from the Cartesian product operation. Originally, the target filter-projector set  $P_1$  built between  $T_1$  and the output example ( $T_{out}$ ) contains exactly 1 element (as

<sup>f</sup>We use  $\langle \dots, \dots \rangle$  to represent these pairs to enable distinguishing them against the  $(\beta, \gamma)$  pairs mentioned in the last section.



```

1 enumFilteringQuery( $I, \text{consts}, \text{holeNum}$ )
  // Input: input table set ( $I$ ), constant set ( $\text{consts}$ ), and an integer ( $\text{holeNum}$ )
  // indicating the number of correlated values in exists-clauses
2  $S \leftarrow \emptyset$ ;
3  $K \leftarrow \text{enumExistsSkeleton}(I, \text{consts}, \text{holeNum})$ ;
4 foreach  $T \in I$ :
5    $F \leftarrow \text{enumFilters}(\text{schema}(T), \text{schema}(T) \cup \text{consts})$ ;
6    $F \leftarrow F \cup \text{enumExistsFilters}(\text{schema}(T), K, \text{holeNum})$ ;
7    $B \leftarrow \text{encodeFilters}(F, T)$ ;
8    $S \leftarrow S \cup \{\text{ST}_p(T, B, T)\}$ ;
9 return  $S$ ;

1 enumAggrQuery( $S, \text{aggregators}, \text{consts}$ )
  // Input: summary tables, aggregators, and constant set.
2  $R \leftarrow \emptyset$ ;
3 foreach  $\rho \in S$ :
4   forall  $\bar{\beta} \subseteq \text{bvSet}(\rho)$ :
5      $T \leftarrow \text{TFilter}(\text{coreTable}(\rho), \bigoplus \bar{\beta})$ ;
6     foreach  $\{c_1, \dots, c_k\} \subseteq \text{columns}(T)$ :
7       foreach  $c_t \in \text{columns}(T), \alpha \in \text{aggregators}$ :
8          $Q' \leftarrow \text{"Select } c_1, \dots, c_k, \alpha(c_t) ;$ 
9            $\text{From decode}(\rho, \bar{\beta}, [1, \dots, \text{colNum}(T)])$ ;
10           $\text{Group By } c_1, \dots, c_k"$ ;
11          $T' \leftarrow \text{eval}(Q')$ ;
12          $F \leftarrow \text{enumFilters}(\{c_t\}, \text{schema}(T) \cup \text{consts})$ ;
13          $B \leftarrow \text{encodeFilters}(F, T')$ ;
14          $R \leftarrow R \cup \{\text{ST}_p(T, B, Q')\}$ ;
15 return  $R$ ;

1 enumJoinQuery( $S_1, S_2, \text{consts}$ )
  // Input: two summary table sets  $S_1, S_2$ , constant set ( $\text{consts}$ )
2  $R \leftarrow \emptyset$ ;
3 foreach  $\rho_1 \in S_1, \rho_2 \in S_2$ :
4    $T_1 \leftarrow \text{BaseTable}(\rho_1), T_2 \leftarrow \text{BaseTable}(\rho_2)$ ;
5    $F \leftarrow \text{enumFilters}(\text{schema}(T_1), \text{schema}(T_2))$ ;
6    $B \leftarrow \text{encodeFilters}(F, T_1 \times T_2)$ ;
7    $R \leftarrow R \cup \{\text{ST}_c(\rho_1, \rho_2, B)\}$ ;
8 return  $R$ ;

1 enumExistsSkeleton( $I, \text{consts}, \text{holeNum}$ )
  // Input: input tables, constant set, number of holes in exists-clause.
  // if  $\text{holeNum} = 0$ : return  $\emptyset$ ;
2  $H \leftarrow \emptyset$ ;
3 foreach  $\tau \in \text{Types}$ :
4    $H \leftarrow H \cup \{h_1 : \tau, \dots, h_{\text{holeNum}} : \tau\}$ ;
5 return  $\text{enumFilteringQuery}(I, H \cup \text{consts}, 0)$ ;

1 enumFilters( $L, R$ )
  // Input: two sets of values  $L$  and  $R$  and a boolean
2  $F \leftarrow \emptyset$ ;
3 foreach  $v_1 : \tau_1 \in L, v_2 : \tau_2 \in R$ :
4   if  $\tau_1 \neq \tau_2$ : continue;
5   foreach  $op \in \text{lookupOperators}(\tau_1)$ :
6      $F \leftarrow F \cup \{v_1 \text{ op } v_2\}$ ;
7 return  $F$ ;

1 enumExistsFilters( $V, K$ )
  // Input: table  $T$ , a set of constants  $\text{consts}$ , and a boolean
2  $F \leftarrow \emptyset$ ;
3 foreach  $\kappa \in K$ :
4    $[h_1 : \tau_1, \dots, h_n : \tau_n] \leftarrow \text{holes}(\kappa)$ ;
5   forall  $[v_1 : \tau'_1, \dots, v_n : \tau'_n] \subseteq V$ :
6     if  $\exists \tau_i \neq \tau'_i (i \in [1, n])$ :
7       continue;
8      $Q \leftarrow [h_1 \mapsto v_1, \dots, h_n \mapsto v_n] \kappa$ ;
9      $F \leftarrow F \cup \{\text{"Exists } Q'\}$ ;
10 return  $F$ ;

1 encodeFilters( $F, T$ )
  // Input: a set of filters on  $T$ , and a table  $T$ 
2  $B \leftarrow \emptyset, n \leftarrow \text{rowNum}(T)$ ;
3 foreach  $f \in F$ :
4    $\beta \leftarrow [b_1 \dots b_n]$  where  $b_i = \text{eval}(f, r_i)_{i \in [1, n]}$ ;
5    $B \leftarrow B \cup \{\beta\}$ ;
6 return  $B$ ;

```

**Figure 9.** Modules in forward enumeration algorithm, where we use function  $\text{eval}(Q)$  to refer to evaluating a query  $Q$  using SynthsQL interpreter,  $\text{eval}(f, r)$  to refer to evaluating a filter on a table row with SQL interpreter, and  $\text{decode}(\rho, \beta, \gamma)$  is the function to decode bv-filter  $\beta$  and projector  $\gamma$  on a summary table  $\rho$  to a query, as will be shown in Figure 14.

```

 $\sigma_{T_1, T_2} ::= \{\omega \mapsto \Omega \mid \omega \in \text{coord}(T_2),$ 
 $\Omega = \{\omega' \in \text{coord}(T_1) \mid \text{val}(T_1, \omega') = \text{val}(T_2, \omega)\}\}$ 

fullInstantiate( $\sigma$ ) ::=
 $\{(\beta, \gamma) \mid \text{bitvecToList}(\beta) = [i_1, \dots, i_{n_2}], \gamma = [j_1, \dots, j_{m_2}],$ 
 $\forall x \in [i_1, \dots, i_{n_2}], y \in [j_1, \dots, j_{m_2}]. (i_x, j_y) \in \sigma(x, y)\}$ 

partialInstantiate( $\sigma, \text{split}$ ) ::=
 $\{(\beta, \gamma) \mid \text{bitvecToList}(\beta) = [i_1, \dots, i_{n_2}], \gamma = [j_1, \dots, j_{m_2}],$ 
 $\forall x \in [i_1, \dots, i_{n_2}], y \in [j_1, \dots, j_{m_2}]. (i_x, j_y) \in \sigma(x, y),$ 
 $\exists k. (j_k > \text{split}) \text{ and } \exists k. (j_k \leq \text{split})\}$ 

 $\text{bitvecToList}([b_1 \dots b_n]) ::= [i \mid b_i = 1]$ 

```

**Figure 10.** The definition of trace.

```

1 backwardInference( $T_1, T_2$ )
2  $M \leftarrow \emptyset$ 
3 foreach  $i = 1$  to  $\text{rowNum}(T_2)$ :
4   foreach  $j = 1$  to  $\text{colNum}(T_2)$ :
5      $\omega \leftarrow (i, j)$ ;
6      $\Omega \leftarrow \{\omega' \mid \omega' \text{ in } T_1 \wedge \text{val}(\omega, T_2) = \text{val}(\omega', T_1)\}$ ;
7      $M \leftarrow M \cup \{\omega \mapsto \Omega\}$ ;
8 return  $M$ ;

```

**Figure 11.** Inferring a mapping instance from  $T_1$  to  $T_2$ .

$P_1 = \{([00110], [1, 2, 3])\}$ ). However, after Cartesian product, each row in input will have  $r$  duplicates in first 3 columns, and the number of filter-projector pairs between  $T_1 \times T'$  and  $T_{out}$  is increased to  $r^{r_{out}}$  (or  $2^{r_{out} * r}$  if we consider projection with 'Select Distinct',  $r_{out}$  is the row number of the output table).

As we will show in the next section, instead of directly instantiating all  $\langle \beta, \gamma \rangle$  pairs using  $\text{fullInstantiate}(\sigma)$ , the merge algorithm can use  $\rho$  to guide instantiate process of  $\sigma$  to avoid the problem listed above.

## 8. Merge Algorithm

As shown in Figure 6, the merge algorithm plays a crucial role in the synthesis pipeline: after obtaining a summary table  $\rho$  and a trace  $\sigma$  built from  $\text{coreTable}(\rho)$  and the output example  $T_{out}$ , the merge algorithm then heuristically searches for all  $(\bar{\beta}, \gamma)$  pairs on  $\rho$  with guide of  $\sigma$ .

Before presenting the merge algorithm, we first introduce several extra operations on both summary tables and mapping instances that are used in defining the merge algorithm.

**Auxiliary Functions** The two auxiliary functions defined on a summary table  $\rho$  are presented in Figure 12, including 1) pushing down a bv-filter in a compound summary table to a bv-filter in its left summary table, 2) heuristically enumerating  $(\bar{\beta}, \gamma)$  in  $\rho$  pairs with the guidance of a target filter-projector pair set  $P$ .

```

1 pushDownL( $\beta, T_1, T_2$ )
2    $\beta' \leftarrow []$ ;
3   for  $i = 1$  to  $\text{rowNum}(T_1)$ :
4      $b_i \leftarrow 0$ ;
5     forall  $j = 1$  to  $\text{rowNum}(T_2)$ :
6        $b_i \leftarrow b_i \vee \beta[i * \text{rowNum}(T_2) + j]$ ;
7      $\beta' \leftarrow \beta' \vee [b_i]$ ;
8   return  $\beta'$ ;

1 guidedEnumFilterCombST( $\rho, \text{target}$ )
2   if  $\rho = \text{ST}_p(T, B, Q)$ :
3      $\text{Combs} \leftarrow \text{guidedCombineST}(B, \text{target}, \text{maxLen})$ ;
4     return  $\{\bar{\beta} \mid \bar{\beta} \in \text{Combs}, (\bigoplus \bar{\beta}) \in \text{target}\}$ ;
5   if  $\rho = \text{ST}_c(\rho_1, \rho_2, B)$ :
6      $T_1 \leftarrow \text{baseTable}(\rho_1)$ ;
7      $T_2 \leftarrow \text{baseTable}(\rho_2)$ ;
8      $R_1 \leftarrow \{\bar{\beta}_1 \mid \bar{\beta}_1 \subseteq \text{bvSet}(\rho_1), \bar{\beta}_1 = \text{pullUpL}(\bar{\beta}, T_1, T_2),$ 
9        $\exists \beta_0 \in \text{target}.((\bigoplus \bar{\beta}_1) \oplus \beta_0 = \beta_0)\}$ ;
10     $R_2 \leftarrow \{\bar{\beta}_2 \mid \bar{\beta}_2 \subseteq \text{bvSet}(\rho_2), \bar{\beta}_2 = \text{pullUpR}(\bar{\beta}, T_1, T_2),$ 
11       $\exists \beta_0 \in \text{target}.((\bigoplus \bar{\beta}_2) \oplus \beta_0 = \beta_0)\}$ ;
12     $\text{Combs} \leftarrow \text{guidedCombineST}(B, \text{target}, \text{maxLen})$ ;
13     $R \leftarrow \{\bar{\beta}_1 \oplus \bar{\beta}_2 \mid \bar{\beta}_1 \in R_1, \bar{\beta}_2 \in R_2, \bar{\beta} \in \text{Combs},$ 
14       $\exists \beta_0 \in \text{target}.((\bigoplus \bar{\beta}_1) \oplus (\bigoplus \bar{\beta}_2) \oplus \beta_0 = \beta_0)\}$ ;
15    return  $R$ ;

1 guidedCombineST( $B, \text{target}, k$ )
2   if  $k = 1$ :
3     return  $\{[\beta] \mid \beta \in B \text{ and } \exists \beta_0 \in \text{target}.(\beta \oplus \beta_0 = \beta_0)\}$ ;
4    $\text{Combs} \leftarrow \text{guidedCombineST}(B, \text{target}, k - 1)$ ;
5    $R \leftarrow \{[\beta_1] \oplus [\beta_2] \mid \beta_1 \in B, \beta_2 \in \text{Combs},$ 
6      $\exists \beta_0 \in \text{target}.(\beta_1 \oplus \beta_2 = \beta_0)\}$ ;
7   return  $R$ ;

```

**Figure 12.** Extra operators on summary tables.

The first function  $\text{pushDownL}(\beta, T_1, T_2)$  can be treated as the inversion of the function  $\text{pullUpL}(\beta, T_1, T_2)$ . The input of the  $\text{pushDown}$  function include a bv-filter  $\beta$  (of size  $\text{rowNum}(T_1) * \text{rowNum}(T_2)$ ) and two tables  $T_1, T_2$ . The output of the function is a bv-filter  $\beta'$  (of size  $\text{rowNum}(T_1)$ ) on  $T_1$ , satisfying the property that  $\text{deDup}(\text{proj}(\text{filter}(T_1 \times T_2, \beta), \gamma)) = \text{deDup}(\text{filter}(T_1, \beta'))$ , where  $\gamma$  is a bitvector representing projection on all fields of  $T_1$  and  $\text{deDup}$  represents the function that removes all duplicated rows in a table.

The second operator is  $\text{guidedEnumFilterCombST}(\rho, P)$ : it takes in a summary table  $\rho$  and a set of target bv-filter set  $B$ , and searches for all bv-filter lists  $\bar{\beta}$  in  $\rho$ , such that  $\bigoplus \bar{\beta} \in B$ .

The key idea of the function is to prune all combinations  $\bar{\beta}$  such that  $\forall \beta_0 \in B.((\bigoplus \bar{\beta}) \oplus \beta_0 \neq \beta_0)$  (i.e. there exists some bit in  $\beta_0$  that is 1 but the corresponding bit in  $\bigoplus \bar{\beta}$  is 0) during the search process: since  $\forall \beta_0 \in B.((\bigoplus \bar{\beta}) \oplus \beta_0 \neq \beta_0)$  indicates that we cannot find other filters to combine with  $\bar{\beta}$  such that the combination result is in  $B$ , and we can safely prune out such combinations.

### 8.1 Merge Algorithm

With the auxiliary functions above, we now define the merge algorithm. As shown in Figure 13, the merge algorithm takes in a summary table  $\rho$  and the trace  $\sigma$  inferred from  $\text{coreTable}(\rho)$  and  $T_{out}$ , and returns a set of pairs of form  $(\bar{\beta}, \gamma)$ , where the first element is a list of bv-filters on  $\rho$  and each  $\beta$  in the list is a primitive filter from  $\text{bvSet}(\rho)$ , while the second element is a projector on  $\text{coreTable}(\rho)$ . Particularly, each pair  $(\bar{\beta}, \gamma)$  in the result set has the property:  $\text{deDup}(\text{proj}(\text{filter}(\text{coreTable}(\rho), \bigoplus \bar{\beta}), \gamma)) = \text{deDup}(T_{out})$ , and these pairs are passed to the next stage (decoding) to obtain all SynthSQL queries that are consistent with input-output examples within the subspace of  $\rho$ .

The merge algorithm takes the following steps to obtain the result. When the summary table is a primitive summary table, all

$(\beta, \gamma)$  pairs are directly instantiated from the program trace  $\sigma$  (line 3), and are directly used to guide searching candidate filter lists in  $\rho$  (lines 4-5) to obtain all candidate  $(\bar{\beta}, \gamma)$  pairs in  $\rho$ . When the summary table is a compound summary table of form  $\text{ST}_c(\rho_1, \rho_2, B)$ , the merge algorithm uses two steps to overcome the challenges listed in problem 2:

- *Step 1.* Instantiates only  $(\beta, \gamma)$  pairs from  $\sigma$  such that  $\gamma$  involves columns from both  $T_1$  and  $T_2$ , using the function  $\text{partialInstantiate}(\sigma, \text{bound})$  with  $\text{bound}$  set to  $\text{colNum}(T_1)$ . With this set  $P$  instantiated from  $\sigma$ , we use  $P$  to guide search filter combinations in  $\rho$  as in primitive summary tables (lines 10-11). This step instantiates all  $(\bar{\beta}, \gamma)$  pairs such that  $\gamma$  involves columns from both sub summary tables and  $\text{coreTable}(\rho)$  can be filtered with  $\bigoplus \bar{\beta}$  and projected with  $\gamma$  to obtain  $T_{out}$  (line 12).
- *Step 2.* In the second step, our goal is to find all  $(\bar{\beta}, \gamma)$  pairs in  $\rho$  such that  $\gamma$  involves only columns in the core table of the left summary table (we name these pairs type-2 pairs on  $\rho$ ) to ensure completeness of the search.

Since directly instantiating all  $(\beta, \gamma)$  pairs whose  $\gamma$  involves only columns of the left summary table is prohibitively expensive, instead of directly searching them in  $\rho$ , we construct a new summary table  $\rho'$  such that each type-2  $(\bar{\beta}, \gamma)$  pairs in  $\rho$  can be mapped to a  $(\bar{\beta}, \gamma)$  pair in  $\rho'$ , and we can achieve the goal of the search by searching over all candidate  $(\bar{\beta}, \gamma)$  pair in  $\rho'$  and mapping back to type-2 pairs in  $\rho$ .

This encoding and decoding process consists of the following steps: 1) first, we enumerate combination of filters that do not use filters in  $\rho_1$  (lines 14,15), 2) then we push down these filters to filters on  $T_1$  using  $\text{pushDownL}$  function, (we also use the mapping for decoding purpose) (line 16), and build  $\rho'_1$  by extending  $\rho_1$  with these pushed down filters, and 3) finally, we recursively call the merge function on  $\rho'_1$  to enumerate filter combinations (which contains much less pairs comparing to those in  $\rho$ ) (lines 18-20) and use the  $\text{map}_1, \text{map}_2$  to map each  $(\bar{\beta}, \gamma)$  pair found in  $\rho'_1$  to pairs in  $\rho$  using  $\text{mapBack}$  (line 22).

The correctness of the algorithm (in terms of completeness) is based on the property that there exists a map between type-2  $(\bar{\beta}, \gamma)$  pairs in  $\rho$  and  $(\bar{\beta}, \gamma)$  pairs in  $\rho'_1$ . As a result, the result in line 22 contains complete  $(\bar{\beta}, \gamma)$  pairs in  $\rho$  that can transform  $\text{coreTable}(\rho)$  to  $T_{out}$ .

After invoking the merge function between each pair of summary table and trace, we obtain all  $(\bar{\beta}, \gamma)$  on the summary tables that can generate  $T_{out}$  within the search space. As shown in Figure 6, the next step is to decode and rank the queries to recommend to the user.

**Lemma 3.** (Soundness of type-2 operations) *Given the two summary tables  $\rho = \text{ST}_c(\rho_1, \rho_2, B)$  and  $\rho'_1$  in the merge algorithm, if a filter-projector pair  $(\bar{\beta}, \gamma)$  is a candidate pair in  $\rho$  whose  $\gamma$  involves only columns in  $\text{coreTable}(\rho_1)$ , then there exists a pair  $(\bar{\beta}', \gamma')$  in  $\rho'_1$ , satisfying the property below and  $(\bar{\beta}, \gamma)$  is in the coding result set of  $(\bar{\beta}', \gamma')$ .*

$$\text{deDup}(T_{out}) = \text{deDup}(\text{proj}(\text{filter}(\text{coreTable}(\rho'_1), \bigoplus \bar{\beta}'), \gamma'))$$

*Proof.* For each  $\beta \in \bar{\beta}$ , if  $\beta$  is a filter generated from pulling up filters in  $\rho_1$  ( $\beta = \text{pullUpL}(\beta_0, T_1, T_2)$ ), where  $\beta_0$  is a filter from  $\rho_1$  and  $T_1$  is the core table of  $\rho_1$ , we have  $\beta_0 \in \text{bvSet}(\rho_1)$ . With this relations, we can decompose  $\beta$  into  $\beta_a$  and  $\beta_b$  where the former are those satisfying the above descriptions and the latter are the rest.

```

1 merge( $\rho, \sigma$ )
  // Input: summary table, mapping instance  $M$  built from  $\rho$  and  $T_{out}$ .
2  if  $\rho = ST_p(T, B, Q)$ :
3     $P \leftarrow \text{fullInstantiate}(\sigma)$ 
4     $C \leftarrow \text{guidedEnumFilterCombST}(\rho, \{\beta \mid \exists (\beta, \gamma) \in P\})$ 
5    return  $\{(\beta, \gamma) \mid \beta \in C \text{ and } (\bigoplus \beta, \gamma) \in P\}$ 
6  if  $\rho = ST_c(\rho_1, \rho_2, B)$ :
7     $T_1 \leftarrow \text{coreTable}(\rho_1)$ 
8     $T_2 \leftarrow \text{coreTable}(\rho_2)$ 
9    // only instantiate pairs whose projector  $\gamma$  involves columns in both sub
    // summary table.
10    $P \leftarrow \text{partialInstantiate}(\sigma, \text{colNum}(T_1))$ 
11    $C \leftarrow \text{guidedEnumFilterCombST}(\rho, \{\beta \mid \exists (\beta, \gamma) \in P\})$ 
12   result  $\leftarrow \{(\beta, \gamma) \mid \beta \in C \text{ and } (\bigoplus \beta, \gamma) \in P\}$ 
13   // Build a new summary table  $\rho'_1$  by pushing down filters from  $\text{bvSet}(\rho_2)$ 
    // and  $\text{bvSet}(\rho)$  to  $\rho_1$ , in a way that by searching  $\rho'_1$ , we can decode the search
    // result to obtain all candidate  $(\beta, \gamma)$  pairs in  $\rho$  whose  $\gamma$  involves only columns
    // in  $\rho_1$ .
14    $R_2 \leftarrow \{[\text{pullUpR}(\beta_i, T_1, T_2) \mid \beta_i \in \bar{\beta}] \mid \bar{\beta} \subseteq \text{bvSet}(\rho_2)\}$ 
15   Combs  $\leftarrow \{\bar{\beta} \mid \bar{\beta} \subseteq B, |\bar{\beta}| = \text{maxLen}\}$ 
16   map1  $\leftarrow \{\beta' \mapsto \bar{\beta}_2 + \bar{\beta} \mid \bar{\beta}_2 \in R_2, \bar{\beta} \in \text{Combs},$ 
     $\beta' = \text{pushDownL}(\bigoplus \bar{\beta}_2 \oplus \bigoplus \bar{\beta}, T_1, T_2)\}$ 
17   map2  $\leftarrow \{\beta \mapsto [\beta'] \mid \beta \in \text{bvSet}(\rho_1),$ 
     $\beta' = \text{promoteL}(\beta, T_1, T_2)\}$ 
18    $\rho'_1 \leftarrow \text{extendPrimBV}(\rho_1, \text{dom}(\text{map}_1))$ 
19    $\sigma'_1 \leftarrow \text{backwardInference}(\text{coreTable}(\rho'_1, T_{out}))$ 
20    $R'_1 \leftarrow \text{merge}(\rho'_1, \sigma'_1)$ 
21   result  $\leftarrow \text{result} \cup \{(\beta', \gamma) \mid (\bar{\beta}, \gamma) \in R'_1,$ 
     $\beta' \in \text{mapBack}(\bar{\beta}, \text{map}_1 \cup \text{map}_2),$ 
     $\text{proj}(\text{filter}(T, \bigoplus \bar{\beta}), \gamma) = T_{out}\}$ 
22   return result

1 extendPrimBV( $\rho, B'$ )
2  if  $\rho = ST_p(T, B, Q)$ :
3    return  $ST_p(T, B \cup B', Q)$ 
4  if  $\rho = ST_c(\rho_1, \rho_2, B)$ :
5    return  $ST_c(\rho_1, \rho_2, B \cup B')$ 

1 mapBack( $[\beta_1, \dots, \beta_n], \text{map}$ )
2  return  $\{\bar{\beta}'_1 + \dots + \bar{\beta}'_n \mid \forall i = 1, \dots, n. (\exists \beta_i \mapsto \bar{\beta}'_i \in \text{map})\}$ 

```

**Figure 13.** Merge Algorithm.

According to the merge algorithm, we have a  $\beta_x \in \text{dom}(\text{map}_1)$  such that filtering  $T_1$  with  $\beta_x$  is equal to projection the result of filtering  $T_1 \times T_2$  with  $\bigoplus \bar{\beta}_b$  by including only columns in  $T_1$ .

Thus, let  $\beta' = \bar{\beta}_a + \{\beta_x\}$ , the formula is then satisfied.  $\square$

**Lemma 4.** (Soundness of Merge) Given a summary tables  $\rho = ST_c(\rho_1, \rho_2, B)$  and the trace  $\sigma$  inferred between  $\text{coreTable}(\rho)$  and  $T_{out}$ , suppose a filter-projector pair  $(\beta, \gamma)$  satisfies the property that  $(\bigoplus \beta, \gamma) \in \text{fullInstantiate}(\sigma)$ , then  $(\beta, \gamma)$  is contained in the result of  $\text{merge}(\rho, \gamma)$ .

*Proof.* This is obvious given Lemma 3: for projections that involves both columns from  $\rho_1$  and  $\rho_2$ , filter-projector pairs are directly included in the result, and for pairs that involves only columns from  $\rho_1$ , according to Lemma 3, they can be discovered by  $R'$  in the algorithm and mapping them back into the compound summary table makes the set complete.  $\square$

## 9. Decoding and Ranking

**Decoding** For each summary table  $\rho$  after the merge algorithm, we obtain the set of all  $(\beta, \gamma)$  pairs from the merge algorithm such that  $\text{proj}(\text{filter}(\text{coreTable}(\rho), \bigoplus \beta), \gamma) = T_{out}$ . Thus, by decoding each  $(\beta, \gamma)$  pairs on  $\rho$  into queries, we obtain the set of all candidate queries within the query space defined by  $\rho$ , whose evaluation results are all  $T_{out}$ .

The decoding algorithm is presented in Figure 14: for a primitive summary table  $\rho$  (lines 3-9), we re-enumerate primitive filters

```

1 decode( $\rho, \bar{\beta}, \gamma$ )
2   $\bar{c} \leftarrow [c_i \mid c_i \in \text{schema}(\text{coreTable}(\rho)) \text{ and } i \in \gamma]$ 
3  if  $\rho = ST_p(T, B, Q)$ :
4     $n \leftarrow |\bar{\beta}|$ 
5     $F_1, \dots, F_n \leftarrow \emptyset$ 
6    foreach  $i = 1, \dots, n$ :
7       $F_i \leftarrow \{f \mid \text{eval}(f, T) = \beta_i\}$ 
8      //  $\beta_i$  is the  $i$ -th element bv-filter in  $\bar{\beta}$ 
9    return  $\{q \mid q = \text{"Select } \bar{c}$ 
     $\text{From } Q \text{ Where } f_1 \text{ And...And } f_n",$ 
     $f_i \in F_i (i = 1..n) \}$ 
10 if  $\rho = ST_c(\rho_1, \rho_2, B)$ :
11   // Decoding two sub summary tables into queries.
12    $\bar{\beta}_1 \leftarrow [\beta \mid \exists \beta' \in \bar{\beta}. (\beta \in \text{bvSet}(\rho_1) \wedge \beta' = \text{pullUpL}(\beta))]$ 
13    $\bar{\beta}_2 \leftarrow [\beta \mid \exists \beta' \in \bar{\beta}. (\beta \in \text{bvSet}(\rho_2) \wedge \beta' = \text{pullUpR}(\beta))]$ 
14    $S_1 \leftarrow \text{decode}(\rho_1, (\bar{\beta}_1, *))$ 
15    $S_2 \leftarrow \text{decode}(\rho_2, (\bar{\beta}_2, *))$ 
16    $\bar{\beta} \leftarrow [\beta \mid \beta \in \bar{\beta}_1 \wedge \beta \in \bar{\beta}_2]$ 
17    $n \leftarrow |\bar{\beta}|$ 
18   // Decoding join predicates.
19    $F_1, \dots, F_n \leftarrow \emptyset$ 
20   foreach  $i = 1, \dots, n$ :
21      $F_i \leftarrow \{f \mid \text{eval}(f, T) = \beta_i\}$ 
22     //  $\beta_i$  is the  $i$ -th element bv-filter in  $\bar{\beta}$ 
23   return  $\{q \mid q = \text{"Select } \bar{c} \text{ From } q_1 \text{ Join } q_2$ 
     $\text{On } f_1 \text{ And } \dots \text{ And } f_n",$ 
     $\forall i = 1..n. (f_i \in F_i \wedge q_1 \in S_1 \wedge q_2 \in S_2)\}$ 

```

**Figure 14.** Decoding Algorithm, where we use “\*” to refer to the projector “Select \*”.

on the  $\text{coreTable}(\rho)$  to decode each bv-filter to a primitive filter predicate on  $T$ , and further obtain a set of candidate queries by filling the holes of the query skeleton with these primitive filter predicates (line 9). For compound summary tables, the decoding process is similar: we first recursively invoke decoding algorithm on each sub summary tables (lines 12-15), then re-enumerate primitive join predicates to decode each bv-vector to a primitive join predicate<sup>§</sup> (lines 16-21), and the result is a set of join queries generated by filling the holes with the decoded sub summary tables and join predicates, as presented in line 23.

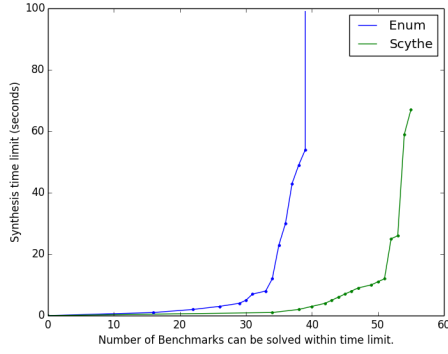
**Lemma 5.** (Completeness of the synthesis algorithm) Suppose  $Q$  is a query expressible in SynthSQL using the given input examples and constraints, such that  $Q$  is consistent with the input output examples, i.e.,  $Q(I) = T_{out}$ , then after decoding phase in the synthesis algorithm,  $Q$  is in the candidates set in Figure 6.

*Proof sketch.* This lemma is the immediate result from combining all of the previous lemmas.

**Ranking** We use a heuristic ranking strategy to rank the queries generated from the decoding algorithm before presenting the synthesis result to the user, by providing a score for each query based on the following three strategies.

- **Simplicity:** Queries with simpler structural and simpler queries are scored higher.
- **Coverage:** Queries with better coverage on the provided constraints (constants and aggregators) are awarded higher score.
- **Naturalness:** More natural join predicates are awarded with higher score, e.g., join on same key are awarded with higher score.

<sup>§</sup> In our implementation, each bv-filter is stored with its metadata (i.e., from which summary table was the bv-filter constructed), and thus bv-filters with same value but in different summary tables can be distinguished in lines 12,13,16.



**Figure 15.** The number of benchmarks can be solved with the increasing time limit specified by  $y$ -axis.

With these three ranking rules, each query in the candidate query list obtains a score, and top rated queries are suggested to the user as the synthesis result.

**User Interaction** When none of the recommended queries are correct (i.e., although these queries are consistent with the input-output examples, they return wrong answer when a user evaluates them on larger input tables), the user can either 1) reject the answer and ask the system to continue searching for more structurally complex queries (as shown in Figure 6), 2) provide a new input-output example from scratch to the system to restart synthesis or 3) modify the existing input-output tables (add, delete some rows in the input tables or modify some cells in the original table and revise the output example correspondingly) to rerun the synthesis process.

As we will show in the evaluation section, with our basic ranking algorithm and interaction model, users can typically obtain a correct query among the top 3 ranked query with at most 2 rounds of interaction, and our algorithm is very practical in practice for solving end-user querying problem.

## 10. Evaluation

To evaluate the effectiveness of our approach in practice, we have implemented our algorithm as a prototype PBE system named Scythe<sup>h</sup> and evaluated it on a set of 60 benchmarks collected from online help forums. Our evaluation aims at answering the following four questions.

- *Q1*: Does the hypothesis that users do provide constants and aggregators when asking questions hold widely in reality?
- *Q2*: Is SynthSQL grammar expressive enough to capture many interesting real-world use case?
- *Q3*: What is the performance of Scythe comparing against existing algorithms?
- *Q4*: Is the ranking algorithm and interaction model effective enough for finding correct users among query candidates?

**Evaluation Setup** Besides implementing Scythe, we have also implemented a PBE system Enum that uses enumerative search approach with equivalent-class partitioning [26] as the core algorithm for performance comparison purpose. Both systems have the same target language SynthSQL, same ranking algorithm and interaction model. Scythe and Enum are implemented in Java and are evaluated on a laptop with Intel Core i5 2.7GHz CPU, 8GB of RAM, and the JVM maximum heap size is set to 4GB.

<sup>h</sup> Our project and evaluation benchmarks are open-sourced and its website has been redacted due to the double-blinded review process.

**Benchmarks** Our benchmarks consist of 60 questions collected from Stack Overflow, under “sql” tags<sup>i</sup>. We choose Stack Overflow posts as our evaluation source since 1) these posts are real-world questions, 2) accepted expert answers are included in the posts, so that we can use them to check the correctness of our synthesized queries, and 3) problems in Stack Overflow are more diverse, benefiting from the none-duplicate policy in Stack Overflow. We collect these questions by looking for highly-voted (with more than 50 upvotes) and recent (asked within 6 months) posts that contain input-output examples. Particularly, we exclude posts where 1) input-output examples are incorrect or incomplete, 2) questions are downvoted by forum experts (scored less than -5) due to vagueness and 3) questions that are targeting for system specific functions (e.g., ‘FOR XML PATH’ in SQLServer).

From these benchmarks, 47 out of 60 questions involve using subqueries: 36 questions require using subqueries in the join-clause, 11 requires using subqueries in exists-clause. Comparing to SqlSynthesizer [37] that does not support subqueries (and is only able to express 13 out of 60 cases), our system captures a wider range of queries in practice. The number of input tables provides in these posts varies between 1-4, the row sizes of the tables vary between 3-12 and the columns sizes vary between 2-9.

**Answer to Q1:** All but two of the question posts contain constants or aggregators that are used in solution query appearing in the posts. This validates our assumption that we can obtain constants and aggregators from end users. In these two cases, the users indirectly provide these the aggregator ‘count’ and constant ‘1’ using the English word ‘duplicates’ or ‘top’ (we manually add the constants and aggregators for these two benchmark for performance evaluation).

**Answer to Q2:** Among these 60 benchmarks, SynthSQL grammar is able to expressive correct answers of 55 benchmarks. The five inexpressible cases requires using advanced operators in SQL that cannot be desugared to SynthSQL grammar, including arithmetic expression (2/5), partitioning (1/5), disjunctive filter predicates (1/5) and pivoting (1/5). Although imperfect, SynthSQL grammar is able to capture many real world cases including argmax, moving average, duplicate processing etc.

**Performance** We ran the 55 benchmarks that are expressible by SynthSQL on both Scythe and Enum. Figure 15 shows the number of benchmarks that can be solved by each algorithm ( $x$ -axis) with increasing time limit ( $y$ -axis). As shown in the figure, Scythe can solve all 55 benchmarks within 90 (and 51/55 within 15 seconds) seconds while Enum can totally only solve 38 of them and it is not able to solve more benchmarks even after increased time limits. Furthermore, we computed the average speedup of Scythe comparing to Enum on the 38 benchmarks that both algorithms can solve, and Scythe has an average speedup of  $17.1\times$  (with minimum  $1.12\times$  and maximum  $422\times$ ).

**User Interaction** We evaluated the effectiveness of our heuristic-based ranking algorithm based on how many rounds does the user need to interact with the database system before a correct answer (comparing against the forum expert answer) is among top 3 of the candidate list. Among 55 benchmarks, 47 benchmarks can be solved without providing new input-output examples, and the rest 8 of them cannot be solved with only one input-output example pair due to underfit. Particularly, these 8 benchmarks can be solved by modifying the original input-output examples (i.e. modify the examples so that new examples can better specify the original problem), we measure the difficulty of this interaction by calculating the number of rows needed to be modified (including addition a new row, deleting an existing row, and modifying values of cells in

<sup>i</sup> <http://stackoverflow.com/questions/tagged/sql>

the row) in the input-output examples. For these 8 examples, the average number of rows needed to be modified is 4.6, which validates that our interaction model remains very effective for end-users.

## 11. Related Work

**Program Synthesis Algorithms** Inductive program synthesis algorithms [3; 1; 11] can be roughly classified into the following five categories: 1) enumeration algorithms [26; 25; 9], 2) constraint-solver aided synthesis [34; 33], 3) type-directed synthesis [24; 8], 4) version space algebra based inference algorithms [17; 27; 12], and 5) stochastic search [28].

Among them, our synthesis algorithm is most closely related to the enumeration approach and version space algebra synthesis algorithms. The former technique [26; 9] typically enumerates programs in a given program space specified by language syntax, and memorize enumerated results as values to avoid re-visiting visited search subspace. The latter [27; 18] uses version space algebra to efficiently maintain programs inferred from output examples using inverse semantics of the language.

Our algorithm differs from these two approaches as we use abstract data structures in both forward enumeration and backward inference phases, and we designed a merge algorithm to efficiently instantiate programs from these abstract data structures to obtain queries that are consistent with input-output examples. The abstract representation and efficient merge algorithm play a key role in addressing the challenge brought by the structure-rich nature of the SQL language, which cannot be solved efficiently by existing synthesis algorithm.

**Programming by Examples** Our system is inspired by Programming by Examples (PBE) applications where users specify a program using input-output examples, and the system synthesizes a program that is consistent with the examples. Different PBE systems have been developed to solve problems in data transformation [12; 30; 31; 14], data extraction [18], map-reduce program synthesis [32], data structure transformation [36] and also SQL query synthesis [37; 35]. Particularly, the two previous programming by examples systems for SQL queries developed by Zhang et al. [37] and Tran et al. [35] restrict SQL syntax by disallowing subqueries due to algorithm efficiency issues. With our new synthesis algorithm, we are able to synthesize a richer set of SQL queries that are used in practice while maintaining synthesis efficiency.

**SQL Query Usability** Besides PBE systems, different approaches [5; 6; 21; 2] have been proposed to increase the usability of the SQL language. One technique to increase SQL usability is by designing graphics query interfaces [21; 2; 7; 15; 20]. These graphical interfaces are often good a certain domains benefiting from a set of predefined operators to support commonly used queries in a the domain. Our system are not domain-specific and are more general in increasing SQL usability for end-users. Another line of prior work is supporting natural language interface for database systems (NLIDB) [20; 13], where users specify a query with natural language sentences, and the system translates the sentences into a SQL query. NLIDB systems are good at generating structural complex lookup queries from well formed descriptions but less satisfying in analytical queries in which their descriptions are difficult to form using natural utterances.

**Interactive Refinement** Different approaches has been proposed to increase the ranking algorithm and user interaction interface to make PBE system easier to locate correct programs from the candidate set, including FlashProg [22], NaLIR [19] and BlinkFill [29]. Our system can be potentially integrated with these techniques to enhance the quality of the synthesized programs.

## 12. Conclusion

In this paper, we presented a new synthesis algorithm for synthesizing structurally rich SQL queries and built a practical PBE system Scythe based on the algorithm. Our algorithm solves the challenges in synthesizing SQL queries by using abstract representations to step-wisely reduce the original problem of searching SQL queries to simpler problem of searching bitvectors and integer lists. We have evaluated our system on 60 real-world benchmarks collected from online Q&A forums and the result indicates that our system is very practical and being able to solve most problems within seconds.

## References

- [1] Rajeev Alur, Rastislav Bodík, Garvit Juniwal, Milo M. K. Martin, Mukund Raghothaman, Sanjit A. Seshia, Rishabh Singh, Armando Solar-Lezama, Emina Torlak, and Abhishek Udupa. Syntax-guided synthesis. In *Formal Methods in Computer-Aided Design, FMCAD 2013, Portland, OR, USA, October 20-23, 2013*, pages 1–8, 2013.
- [2] Sonia Bergamaschi, Francesco Guerra, Matteo Interlandi, Raquel Trillo-Lado, and Yannis Velegrakis. Quest: a keyword search system for relational data based on semantic and machine learning techniques. *Proceedings of the VLDB Endowment*, 6(12):1222–1225, 2013.
- [3] Rastislav Bodík and Barbara Jobstmann. Algorithmic program synthesis: introduction. *STTT*, 15(5-6):397–411, 2013.
- [4] Surajit Chaudhuri. An overview of query optimization in relational systems. In *Proceedings of the seventeenth ACM SIGACT-SIGMOD-SIGART symposium on Principles of database systems*, pages 34–43. ACM, 1998.
- [5] Alvin Cheung and Armando Solar-Lezama. Computer-assisted query formulation. *Foundations and Trends in Programming Languages*, 3(1):1–94, 2016.
- [6] Alvin Cheung, Armando Solar-Lezama, and Samuel Madden. Optimizing database-backed applications with query synthesis. In *ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '13, Seattle, WA, USA, June 16-19, 2013*, pages 3–14, 2013.
- [7] Jonathan Danaparamita and Wolfgang Gatterbauer. Queryviz: helping users understand sql queries and their patterns. In *Proceedings of the 14th International Conference on Extending Database Technology*, pages 558–561. ACM, 2011.
- [8] Jonathan Frankle, Peter-Michael Osera, David Walker, and Steve Zdancewic. Example-directed synthesis: A type-theoretic interpretation. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2016*, pages 802–815, New York, NY, USA, 2016. ACM.
- [9] Joel Galenson, Philip Reames, Rastislav Bodík, Björn Hartmann, and Koushik Sen. Codehint: Dynamic and interactive synthesis of code snippets. In *Proceedings of the 36th International Conference on Software Engineering*, pages 653–663. ACM, 2014.
- [10] Richard A Ganski and Harry KT Wong. Optimization of nested sql queries revisited. In *ACM SIGMOD Record*, volume 16, pages 23–33. ACM, 1987.
- [11] Sumit Gulwani. Dimensions in program synthesis. In *Proceedings of the 12th international ACM SIGPLAN symposium on Principles and practice of declarative programming*, pages 13–24. ACM, 2010.
- [12] Sumit Gulwani. Automating string processing in spreadsheets using input-output examples. In *ACM SIGPLAN Notices*, volume 46, pages 317–330. ACM, 2011.
- [13] Sumit Gulwani and Mark Marron. Nlyze: Interactive programming by natural language for spreadsheet data analysis and manipulation. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 803–814. ACM, 2014.
- [14] William R Harris and Sumit Gulwani. Spreadsheet table transformations from examples. In *ACM SIGPLAN Notices*, volume 46, pages 317–328. ACM, 2011.

- [15] Sean Kandel, Jeffrey Heer, Catherine Plaisant, Jessie Kennedy, Frank van Ham, Nathalie Henry Riche, Chris Weaver, Bongshin Lee, Dominique Brodbeck, and Paolo Buono. Research directions in data wrangling: Visualizations and transformations for usable and credible data. *Information Visualization Journal*, 10(4):271–288, 2011.
- [16] Dileep Kini and Sumit Gulwani. Flashnormalize: Programming by examples for text normalization. In *Proceedings of the 24th International Conference on Artificial Intelligence, IJCAI’15*, pages 776–783. AAAI Press, 2015.
- [17] Tessa Lau, Steven A Wolfman, Pedro Domingos, and Daniel S Weld. Programming by demonstration using version space algebra. *Machine Learning*, 53(1-2):111–156, 2003.
- [18] Vu Le and Sumit Gulwani. Flashextract: A framework for data extraction by examples. In *ACM SIGPLAN Notices*, volume 49, pages 542–553. ACM, 2014.
- [19] Fei Li and Hosagrahar V Jagadish. Nalir: an interactive natural language interface for querying relational databases. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 709–712. ACM, 2014.
- [20] Fei Li and HV Jagadish. Constructing an interactive natural language interface for relational databases. *Proceedings of the VLDB Endowment*, 8(1):73–84, 2014.
- [21] Edward Lu and Ras Bodik. Quicksilver: Automatic synthesis of relational queries. Master’s thesis, EECS Department, University of California, Berkeley, May 2013.
- [22] Mikaël Mayer, Gustavo Soares, Maxim Grechkin, Vu Le, Mark Maron, Oleksandr Polozov, Rishabh Singh, Benjamin Zorn, and Sumit Gulwani. User interaction models for disambiguation in programming by example. In *Proceedings of the 28th Annual ACM Symposium on User Interface Software & Technology*, pages 291–301. ACM, 2015.
- [23] Mauro Negri, Giuseppe Pelagatti, and Licia Sbatella. Formal semantics of sql queries. *ACM Transactions on Database Systems (TODS)*, 16(3):513–534, 1991.
- [24] Peter-Michael Osera and Steve Zdancewic. Type-and-example-directed program synthesis. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2015*, pages 619–630, New York, NY, USA, 2015. ACM.
- [25] Phitchaya Mangpo Phothilimthana, Tikhon Jelvis, Rohin Shah, Nishant Totla, Sarah Chasins, and Rastislav Bodik. Chlorophyll: Synthesis-aided compiler for low-power spatial architectures. In *ACM SIGPLAN Notices*, volume 49, pages 396–407. ACM, 2014.
- [26] Phitchaya Mangpo Phothilimthana, Aditya Thakur, Rastislav Bodik, and Dinakar Dhurjati. Scaling up superoptimization. In *Proceedings of the Twenty-First International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 297–310. ACM, 2016.
- [27] Oleksandr Polozov and Sumit Gulwani. Flashmeta: A framework for inductive program synthesis. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*, pages 107–126. ACM, 2015.
- [28] Eric Schkufza, Rahul Sharma, and Alex Aiken. Stochastic superoptimization. *ACM SIGPLAN Notices*, 48(4):305–316, 2013.
- [29] Rishabh Singh. Blinkfill: Semi-supervised programming by example for syntactic string transformations.
- [30] Rishabh Singh and Sumit Gulwani. Learning semantic string transformations from examples. *Proceedings of the VLDB Endowment*, 5(8):740–751, 2012.
- [31] Rishabh Singh and Sumit Gulwani. Transforming spreadsheet data types using examples. *ACM SIGPLAN Notices*, 51(1):343–356, 2016.
- [32] Calvin Smith and Aws Albarghouthi. Mapreduce program synthesis. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 326–340. ACM, 2016.
- [33] Armando Solar-Lezama. *Program Synthesis by Sketching*. PhD thesis, Berkeley, CA, USA, 2008. AAI3353225.
- [34] Emina Torlak and Rastislav Bodik. A lightweight symbolic virtual machine for solver-aided host languages. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI ’14*, pages 530–541, New York, NY, USA, 2014. ACM.
- [35] Quoc Trung Tran, Chee-Yong Chan, and Srinivasan Parthasarathy. Query by output. In *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data*, pages 535–548. ACM, 2009.
- [36] Navid Yaghmazadeh, Christian Klinger, Isil Dillig, and Swarat Chaudhuri. Synthesizing transformations on hierarchically structured data. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 508–521. ACM, 2016.
- [37] Sai Zhang and Yuyin Sun. Automatically synthesizing sql queries from input-output examples. In *Automated Software Engineering (ASE), 2013 IEEE/ACM 28th International Conference on*, pages 224–234. IEEE, 2013.