

Automatically Synthesizing SQL Queries from Input-Output Examples

Sai Zhang Yuyin Sun
Computer Science & Engineering
University of Washington, USA
{szhang, sunyuyin}@cs.washington.edu

Abstract—Many computer end-users, such as research scientists and business analysts, need to frequently query a database, yet lack the programming knowledge to do such tasks smoothly. To alleviate this problem, we present a *programming by example* technique (and its tool implementation, called SQLSynthesizer) that permits end-users to automate such querying tasks. Our technique takes from users an input and output example of how the database should be queried, and then synthesizes a SQL query that reproduces the example output from the example input. Later, when the synthesized SQL query is applied to another, potentially larger, database with a similar schema as the example input, the synthesized SQL query produces a corresponding result that is similar to the example output.

[[may delay this paragraph later]] Our technique has several notable features: it only needs small input-output examples to infer a desirable SQL query, it is fully automated and does not require users to provide annotations/hints of any form, and it can rank multiple possible solutions to provide to users the most likely result.

[[num is missing here]] We evaluated SQLSynthesizer on XXX SQL exercises from a classic database textbook and XXX non-trivial SQL questions raised by real-world users from popular online forums. As a result, SQLSynthesizer infers correct SQL queries for XXX textbook exercises and XXX forum questions, including those that receive no human replies.

I. INTRODUCTION

The big data revolution over the past few years has resulted in significant advances in digitization of massive amounts of data and accessibility of computational devices to massive proportions of the population. A key challenge faced by many enterprise or computer end-users nowadays is the management of their increasingly large and complex databases.

A. End-users' Difficulties in Writing SQL Queries

Although the relational database management system (RDBMS) and the de facto language (SQL) are perfectly adequate for most end-users' needs [16], the costs associated with deployment and use of database software and SQL are prohibitive.

The problem is exacerbated by the fact that many end-users have myriad diverse backgrounds including research scientists, business analysts, commodity traders, human resource managers, finance professionals, and marketing managers. Those end-users who need to query databases to analyze the data are not professional programmers, but are experts in some other domains. They need to query a variety of information on their database and use the information to support their business

decisions. **[[move the following sentence to somewhere else]]** For example, as pointed out by [10], conventional RDBMS software remains underused in the long tail of science: the large number of users, such as the research scientists who are in relatively small labs and individual researchers, have limited IT training, staff and infrastructure yet collectively produce the bulk of scientific knowledge.

Most non-professional database end-users often learn SQL queries by first **[[x]]** consulting a textbook or browsing some online resources to get familiar basic idioms of SQL. Then, they may try to write some experimental SQL queries, execute them on a sample database to explore the data, and subsequently **[[xx]]** modify the queries to write new queries based on the query results. However, such practice is inefficient and time-consuming. **[[explain why]]**. Those end-users can clearly understand their goals but simply can not write a correct SQL query for their tasks, either due to the syntax complexity of the language itself, or the structure complexity of the underlying databases, or others. To assist them in performing database query tasks, a highly accessible tool that can be used to “describe” their needs and connect their intentions to executable SQL queries would be highly desirable.

B. Existing Solutions

Graphical User Interfaces (GUIs) and *general programming languages* are two state-of-the-art approaches in helping end-users perform database queries. However, both of them are far from satisfactory.

Many RDBMS come with a well-designed GUI with tons of features. However, a GUI is often fixed, and does not permit users to personalize a database's functionality for their specific tasks. On the other hand, as a GUI supports more and more customization features, users may struggle to discover those features, which can significantly degrade its usability.

General programming languages, such as SQL, Java (with JDBC), or other domain specific query languages, **[[add non-SQL]]** serve as a fully expressive medium for communicating a user's intention to a database. However, learning a practical programming language (even a simplified, high-level domain specific language) often requires a substantial amount of time and energy that a typical end-user would not prefer, and should not be expected, to invest.

C. Synthesizing SQL Queries from Input-Output Examples

In this paper, we present a technique (and its tool implementation, called SQLSynthesizer) to automatically synthesize SQL queries from input-output examples. SQLSynthesizer takes example input table(s) and output table from the end-users, and then automatically infers a SQL query (or multiple queries, if exist) that queries against the input tables and returns the output example. If the inferred SQL query is applied to the example input, then it produces the example output; and if the SQL query is applied to other similar inputs (potentially much larger tables), then the SQL query produces a corresponding output.

[[explain why input-output examples is well-suited?]]

Although input-output examples may lead to underspecification, they serve as a straightforward way of describing *what* the task is and a natural interface that a tool can provide assistance.

Our technique is designed to be used by non-professional **[[beginner]]** database end-users when they do not know how to write a desirable SQL query. It aims to replace the role of the forum expert, **[[also trainers?]]** which not only removes human from the loop, but also enables end-users to solve their problems in a few seconds as opposed to a few days (while waiting for an expert's reply). **[[right place?]]** We also envision this technique to be useful in an online education setting. Recently, several education initiatives such as EdX, Coursera, and Udacity are teaming up with experts to provide high quality online courses to several thousands of students worldwide. One challenge, which is not present in a traditional classroom setting, is to provide answers on questions raised by a large number of students. **[[explain why SQLSynthesier is useful.]]**

SQLSynthesizer uses three steps to link a user's intention to a desirable SQL query:

[[add to somewhere]] helping end-users writing a wide range of SQL queries.

- XXX
- XXX
- XXX

Compared to previous approaches [], SQLSynthesizer has several notable features:

- XXX
- XXX

D. Evaluation

We evaluated SQLSynthesizer on XXX exercises from a classic database textbook [24]. SQL exercises from a database textbook often cover the most widely-used SQL features that a course instructor wishes students to master. **[[expressiveness]]** Those exercises can be served as golden test to evaluate the expressiveness of our supported SQL subset. Second, we collected 5 non-trivial SQL problems raised by real end-users on online help forums, and tested whether our technique can effectively synthesize desirable SQL queries.

As a result, our technique successfully synthesized 5 out of 6 textbook exercises and solved all 5 forum problems, within

a very small amount of time (1 minute per exercise/problem, on average).

E. Contributions

This paper makes the following contributions:

- **Technique.** We present a technique automatically synthesize SQL queries from input-output examples (Section IV).
- **Implementation.** We implemented our technique in a tool, called SQLSynthesizer (Section V). It is available at: <http://sqlsynthesizer.googlecode.com>.
- **Evaluation.** We applied SQLSynthesizer to XXX textbook SQL exercises and XXX SQL-related problems from online forums. The results show the usefulness of the proposed technique (Section VI).

II. ILLUSTRATING EXAMPLE

We use a SQL exercise from a classic database textbook [24] (Chapter 5, Exercise **[[xx]]**) as an example to illustrate the use of SQLSynthesizer.

Find out the name and max score of the students whose level is senior and enrolled in more than 3 courses.¹

Although the problem description is quite clear and most end-users have a clear idea of what the desirable query should do, they **[[still have difficult ...]]**

As an alternative solution, end-users can use SQLSynthesizer to solve this problem. As depicted in Figure 1, a user provides SQLSynthesizer two input tables and one expected output table. SQLSynthesizer automatically inferred the desirable SQL query (Figure 1(b)) to produce the output from the input tables.

The synthesized query by SQLSynthesizer first joins the student, then ... **[[project, and aggregation]]**

[[add the following justification]] To the best of our knowledge, none of the existing techniques can do ...

III. SUPPORTED SQL SUBSET

This section presents the design of a supported SQL subset in SQLSynthesizer.

The first step is to ... **[[]]** However, no systematic study has ever been conducted to this end, and little empirical knowledge has ever been provided. Without such study or empirical knowledge, understanding XXX remains difficult.

To address this issue, we conduct an online survey and a series of follow-up email interviews.

A. Online Survey: Eliciting Design Requirements

[[Explain why a survey is needed (the language is too big). It will become intractable. And how]]

Decide the important **[[Took a different way than existing work in eliciting design requirement.]]**

Our online survey consists of 6 questions that can be divided into three parts. The first part includes simple demographic questions about participants. In the second part, participants

¹As described in the textbook, there is a student table having two columns: student_id and nameA, and an enrolled table having three columns: student_id, course_id, and score.

student_id	name	level
1	Adam	senior
2	Bob	junior
3	Erin	senior
4	Rob	junior
5	Dan	senior
6	Peter	senior
7	Sai	senior

student_id	course_id	score
1	1	4
1	2	2
2	1	3
2	2	2
2	3	3
3	2	1
4	1	4
4	3	4
5	2	5
5	3	2
5	4	1
6	2	4
6	4	5
7	1	2
7	3	3
7	4	4

```

select student.name, max(enrolled.score)
from student, enrolled
where student.student_id = enrolled.student_id
and student.level = 'senior'
group by student.student_id
having count(enrolled.course_id) > 3

```

name	max_score
Dan	5
Sai	5

(a) Two input tables: student (Left) and enrolled (Right)

(b) A SQL query inferred by SQLSynthesizer

(c) The output table

Fig. 1. Example input-output tables and the SQL query synthesized by SQLSynthesizer. In this example, users provide SQLSynthesizer with two input tables (shown in (a)) and an output table (shown in (c)); SQLSynthesizer automatically infers a SQL query (shown in (b)) that transforms the two input tables into the output table.

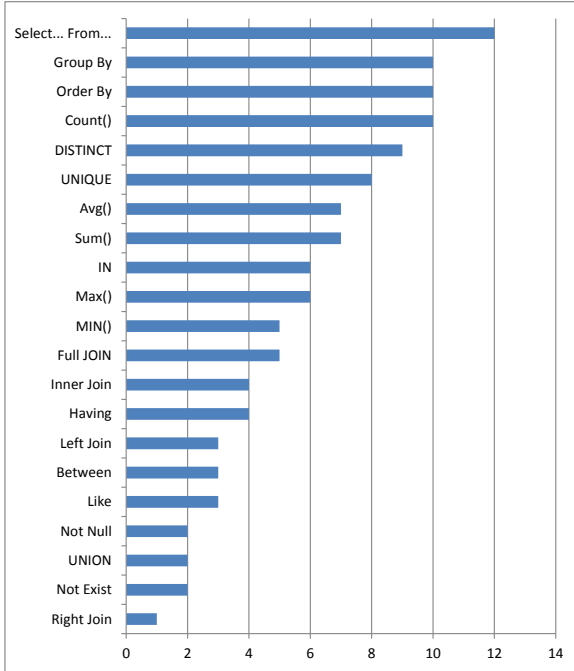


Fig. 2. Survey results...

were asked to select their most-often **[[most important?]]** used SQL language features. Instead of directly asking participants about the SQL features, which might be vague and difficult to respond, we presented them a list of *all* standard SQL features in writing a query. Additionally, participants were asked to report their own experience in writing SQL queries in the third part of the survey.

Before distributing our survey, we conducted pilot interviews with three graduate students with XXX experience at University of Washington. We ran the survey with them and made notes of their comments. According to their feedback, we refined

```

<query> ::= SELECT <expr>+ FROM <table>+
          WHERE <cond>+
          GROUP BY <column>+ HAVING <cond>+
          ORDER BY <column>+

<table> ::= atom
<column> ::= <table>.atom
<cond> ::= <cond> && <cond>
          | <cond> || <cond>
          | ( <cond> )
          | <cexpr> <op> <cexpr>
          | NOT EXIST ( <query> )

<op> ::= = | > | <
<cexpr> ::= const | <column>
<expr> ::= <cexpr> | count( <column> )
          | sum( <column> ) | max( <column> ) | min( <column> )

```

Fig. 3. Syntax of the supported SQL subset.

the survey questions and adjusted the wording to make sure that the questions are relevant and clear.

We sent out invitation to graduate mailing lists at University of Washington, and posted our survey on professional online forums (e.g., StackOverflow). As of April 2013, we received XXX responses. On average, the respondents have XXX years of experience in Computer Science, and XXX years of experience in using database.

Figure ??

B. Syntax

Figure 3 defines the syntax of the supported language, which is a subset of the standard SQL 93 language. This SQL subset supports common query operations across multiple tables (i.e., `select ... from ... where..`), and conjunction of predicates. The language subset shares the same semantics with the standard SQL language.

[[some design rationale]] 1. discard some vendor-specific grammars, like top in Microsoft,

2 discard fuzzy matching part for example, like " ... ", wildcards, makes undecidable, can use string manipulation.

3. exclude delete, update query statement, because we focus on query

4. Omit distinct, in (which can be replaced by NOT EXIST, or condition query)

5. omit between, discard Left Joint, Right Join, only remaining Full Join

6. Not Null (can be simply replaced by != NULL)

7. Nested query can always be re-written using conditions

[[revise the following]] Differing from language subsets used in the existing work [5], our language subset supporting joining operations across multiple tables, and includes widely-used database operations such as group by, order by, and having, as well as a few common aggregation functions such as count, sum, max, and min.

C. Follow-up Interviews: Feedback about the SQL Subset

After prototyping the SQL language in Figure ??, we performed follow-up email interviews to gain participants' feedback about the tailored SQL subset. Participants were first asked to rate the expressiveness of the SQL subset in Figure ?? on a 6-point scale (5-sufficient; 0-not sufficient at all), and then to provide their comments. **[[xxx]]**.

IV. TECHNIQUE

This section explains SQLSynthesizer's 3 steps, as introduced in Section I and illustrated in Figure 4.

[[add some transition sentence]]

A. Query Skeleton Creation

In this step, SQLSynthesizer performs a simple scan over the provided input-output examples to capture the basic structure of the desirable SQL query.

A query skeleton is an incomplete SQL query, consisting of three basic parts: the table set used in the query, the table columns used to join tables, and the table columns used to project the query results. Besides, other query parts, such as query conditions and aggregation operators, are left as unknown and will be determined later.

We next describe how SQLSynthesizer determines the table set, joining columns, and projection columns used in a query.

Determining the Table Set. End-users are unwilling to provide more than enough inputs. Thus, every table in the input example is expected to be used in the desirable SQL query. Based on this observation, we assume that every input table should be used *as least* once in the result query. By default, the table set contains all given input tables. However, it is possible that one input table will be used for multiple times. SQLSynthesizer does not forbid this case, rather, it uses a heuristic to estimate the table set: if the *same* column from an input table appears more than once in the example output, we add the input table to the table set the same number of times.

```
select      Student.Student_name, <Aggregation>
from        student, enrolled
where       student.Student_key = enrolled.Student_key
           and <Conditions>
group by    Student.Student_name
having      <Conditions>
```

Fig. 5. The SQL skeleton created for the motivating example in Figure 1.

[[we view it as a strong indicator that this table will be joined multiple times and add it to our table set T using an alias]]

Determining Joining Columns. Given a set of tables, there are many ways to join them. Enumerating all possibilities can introduce a huge number of joining conditions and would quickly become intractable. To prune the search space, we use three simple but effective effective rules. First, tables are often joined on their primary keys with the same data types, such as joining the student table with the enrolled table on the student_id column (Figure 1). By contrast, it is unlikely to join two tables with an Integer column and a String column. Second, tables are often joined on columns with the same name, such as joining the student table with the enrolled table on the student_name column. Third, it is only meaningful to join two tables on columns that have the same data type and some overlapped values.

SQLSynthesizer restricts the search space in uses the above three rules **[[need to revise]]**

[[need to implement above.]] [[mention how many skeletons will be created]]

Determining Output Columns. For each column in the output table, SQLSynthesizer checks whether its name appears in any of the input table. If so, SQLSynthesizer uses the matched column from the input table as the output column. to the output column set. Otherwise, the output column must be produced by using aggregation operators. **[[same names]]** Consider the example in Figure 1, SQLSynthesizer determines that column name comes from the student table, while column max_Score must be created by using an aggregation operator. **[[If there is no column name]] [[check the values in the output column]]**

[[It is possible that multiple skeleton can be created. add an algorithm here.]]

Figure 5 shows the created query skeleton for the motivating example in Figure 1. In this skeleton, three unknown structures represented by <Aggregation> or <Conditions> are in red, and will be filled in the next phase. **[[revise text]]**

[[how to create group by]]

B. SQL Query Completion

In this step, SQLSynthesizer takes as input a query skeleton, infers the missing parts, namely query conditions (Section IV-B1) and aggregation operators (Section IV-B2); and then produces a set of syntactically-correct SQL queries.

1) *Inferring Query Conditions:* The problem of inferring query conditions can be cast as finding appropriate rules that can perfectly divide the whole searching space into positive

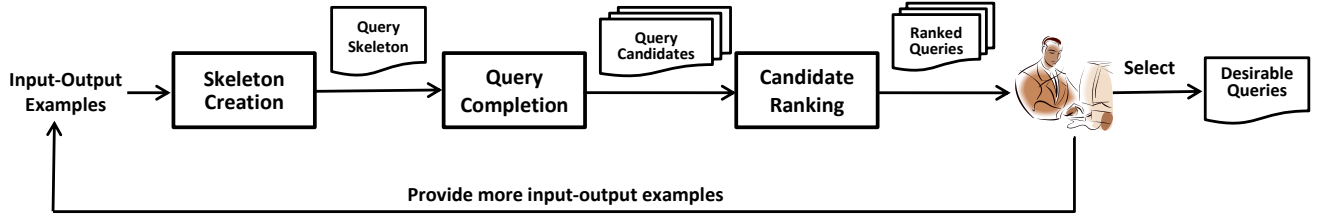


Fig. 4. Illustration of SQLSynthesizer’s workflow in inferring SQL queries from input-output examples. SQLSynthesizer consists of 3 steps: (1) The “Skeleton Creation” step (Section IV-A) infers a partially-complete SQL query as a skeleton from the given examples, (2) The “Query Completion” step (Section IV-B) uses a rule-learning algorithm [] and type-directed search to produces a list of syntactically-valid query candidates that satisfy the provided input-output examples. , and (3) The “Candidate Ranking” step (Section IV-C) ranks all inferred SQL query candidates and provides users a ranked list of SQL queries with the likely ones on the top. Users select the desirable SQL queries from the produced ranked query list. On some examples, if SQLSynthesizer produces SQL queries that satisfy the input and output examples, but does not address the intention that the user wants; SQLSynthesizer can be used interactively by asking users to provide more informative examples and then refine the SQL queries.

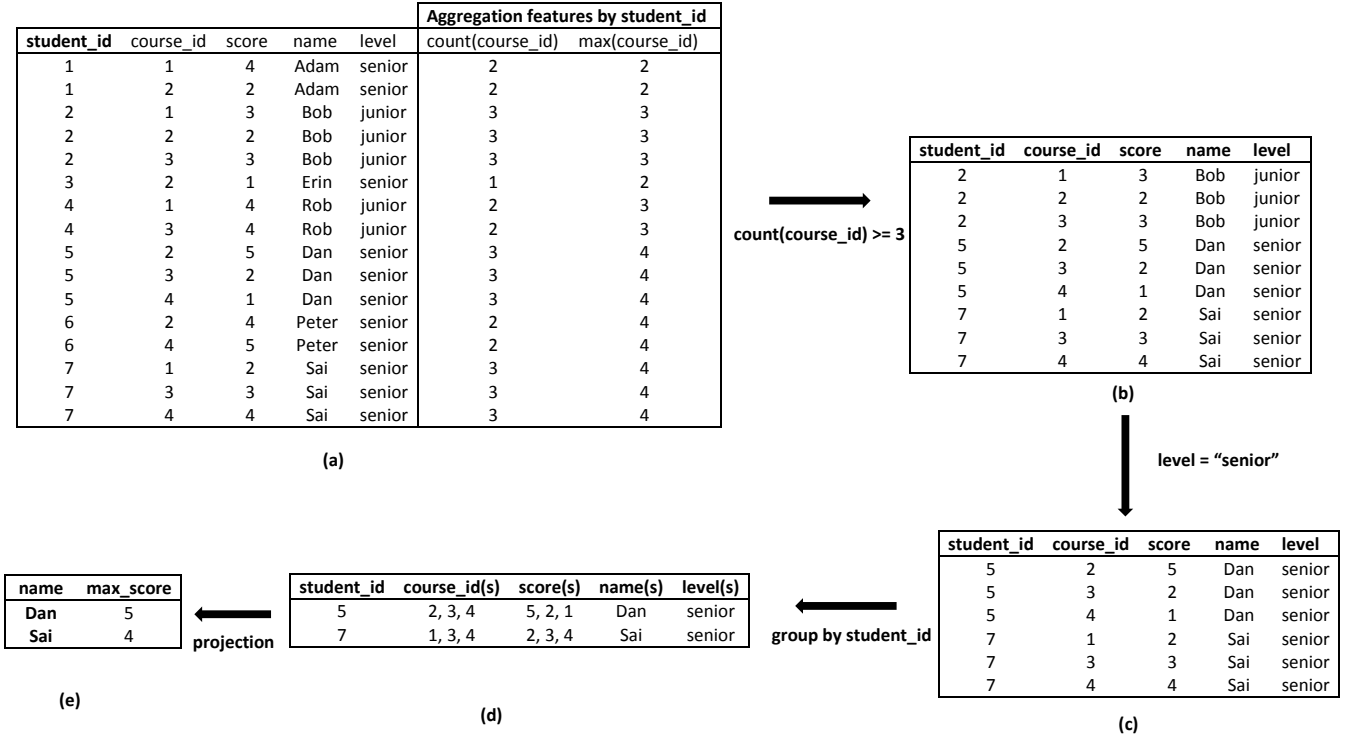


Fig. 6. Illustration of xxx

part and negative part. In our context, the searching space contains all tuples generated by joining the input tables, the positive part are all tuples in the output table, and the negative part are the rest tuples.

SQLSynthesizer uses the PART learning algorithm [9] from the machine learning community to infer a set of accurate and expressive rule set. Compared to other well-known learning algorithms like decision-tree learning or **[unclear]**, PART has two notable features. **[unclear]** First, it utilizes the “divide-and-conquer” strategy to repeatedly build rules and remove the instances it covers until no examples are left. When creating each rule, it employs a pruned decision tree built from current set of instance and only makes the leaf with the largest coverage into the resulting rules, without keeping the whole learned tree in memory.

A key challenge in using the PART algorithm is how to

devise meaningful features. Existing approaches [] simply uses concrete values in a tuple as features. However, doing so loses much useful structure information needed in a SQL query. **[not enough, no support for max, min, sum]** Thus, we must transform tuples into appropriate feature representation.

SQLSynthesizer addresses this problem by encoding two types of additional features for each tuple:

- **Aggregation Features.** Aggregation features are the aggregation results grouped by each String type column over every other columns. Table 7 shows an example. (1) Aggregation, including COUNT, MAX, MIN and AVG, whose results might be used in query condition.
- **Comparison Features.** Comparison feature is the result of comparing two comparable columns. We consider two possible values: {1,0}, which represents means whether two columns under comparison satisfy the predicate or not,

group by	aggregation
C_1	COUNT(C_2), MAX(C_3), MIN(C_3), AVG(C_3)
C_2	COUNT(C_1), MAX(C_3), MIN(C_3), AVG(C_3)

Fig. 7. The generated aggregation features for a table with 3 columns: C_1 , C_2 , and C_3 , in which columns C_1 and C_2 are String type and column C_3 is Integer type.

predicate	comparison result
$C_1 = C_2$	0
$C_1 < C_2$	1
$C_1 > C_2$	0
$C_3 = C_4$	1
$C_3 < C_4$	0
$C_3 > C_4$	0

Fig. 8. The generated comparison features for a table with 4 columns: C_1 , C_2 , C_3 , and C_4 , in which columns C_1 and C_2 are String type, and columns C_3 and C_4 are Integer type. Columns with the same type are comparable, such as C_1 and C_2 , and C_3 and C_4 .

respectively. Comparison results between two comparable columns. The above two additional knowledge encoding permits our technique to make use of correlations between columns, rather than only values from each isolated and sequential columns. Table 8 shows an example.

Combining using concrete tuple values, aggregation features, and comparison features, our technique is able to extract expressive feature representation for tuples, and permits users to encode domain knowledge and structural information about a SQL query.

[[should give an example here of why such feature is useful]]

The PART algorithm returns rules representing three types of features. **[[which 3 types?]]** **[[move below text to the example figure]]** For example, it returns the following rules for our motivating example in Section II:

```
COUNT(enrolled.course_id) > 2
&& student.level = 'senior'.
```

SQLSynthesizer next splits the returned rules into two parts: the query selection condition, and the having condition. SQLSynthesizer treats rules derived from the value features as the query conditions (COUNT(enrolled.course_id) > 2 for the example of Figure 1), and rules derived from the aggregation features as the having condition (&& student.level = 'senior' for the example of Figure 1). **[[explain why?]]**

2) *Searching for Aggregation Operators:* For each column produced by an aggregation operator, the whole search space includes all possible combinations of table columns and the five supported aggregation operators (see Figure 3). SQLSynthesizer leverages the following two observations to further reduce the search space:

- The data type of an output column must be compatible with the aggregation operator's return type. For instance, if an output column has the String type, it must not use aggregation operators (e.g., count and sum) that returns an Integer.

name	score
Bob	4
Dan	5
Jim	2

(a) The input table: student

name
Bob
Dan

(b) The output table

1. `select name from student where score > 3`
2. `select name from student where name = 'Bob' or name = 'Dan'`

Fig. 9. Illustration of SQLSynthesizer's query candidate ranking strategy. SQLSynthesizer produces two queries for the given input-output examples. Based on our heuristic (Section IV-C), the first query differs from the second query by using simpler conditions, and thus ranks higher (and is less likely to overfit the given examples).

- If an arithmetic aggregation operator, such as max and min, is used, each value in the output column must have appeared in the input table.

[[Order by structure, relatively each to add]]

C. Query Candidate Ranking

It is possible that multiple SQL queries satisfying the given input-output examples will be returned. This may adversely impact end-users who want to write simple query tasks but now may require to provide more bits for disambiguation of their intent (which manifests in the need to provide more examples and more rounds of interaction). To alleviate this problem, we employ the Occam's razor principle, which states that the simplest explanation is usually the correct one, to rank the more likely queries higher in the output list. We define a comparison scheme between different SQL queries by defining a partial order between them. Some of these choices are subjective, but have been observed to work well.

A SQL query is simpler than another one if it uses smaller number of query conditions (including having and order by) or the expressions in each query condition are pairwise simpler. Simpler query conditions suggests the extraction logics are more common and general.

In our implementation, SQLSynthesizer computes a cost for each query, and prefers queries with lower costs. The cost for a SQL query is computing by summarizing the text length of each query condition used. Figure 9 shows an example.

[[write an algorithm here?]]

D. Discussions

[[limitations]]

[[soundness and completeness]]

End-users can use our approach to obtain SQL query to transform multiple, huge database tables by constructing small, representative input and output example tables. On some examples, we speculate that our approach may produce a SQL query that satisfies the input and output examples given by the user, but does not address the intention that the user wants. To address this issue, we adapt a simple interaction model from [15] to ask users to investigate the results of an output SQL query and report any discrepancy. In this case, the user can

refine the inferred SQL query by providing a more informative input-output example (or multiple input-output examples that together describe the required behavior) that demonstrate the behavior on which the originally-inferred SQL query behaves incorrectly.

V. IMPLEMENTATION

We implemented the proposed technique in a tool, called SQLSynthesizer. SQLSynthesizer uses the Weka toolkit [14] to implement the rule-learning algorithm described in Section IV-B. **[[inspect the above]]**

SQLSynthesizer uses MySQL [21] as the backend database to check the correctness of each inferred SQL query. Specifically, SQLSynthesizer populates the backend database with the given input tables. When a SQL query is synthesized, SQLSynthesizer executes the query on the database to observe whether the output matches the given output.

VI. EVALUATION

Our evaluation evaluates SQLSynthesizer’s effectiveness in the following our aspects:

- Is the supported SQL subset expressive enough to describe a variety of queries?
- the time cost of SQLSynthesizer in synthesizing SQL queries (Section ??).
- c

A. Benchmarks

To answer these two questions, we collected a set of SQL query synthesis scenarios in the following two ways. First, we picked up 6 SQL exercises from a classic database textbook [24]. Those 6 SQL exercises cover many commonly-used SQL features, which a database course instructor think would be the most useful parts, and expect her students to get familiar with. Second, we collected 5 SQL problems raised by real end-users from popular online help forums [6], [26], [28].

For a given exercise or problem, if it has already been associated with input-output examples (e.g., those provided by users from online forums), we directly apply our tool on the existing examples. Otherwise, we manually provide small concise and representative input-output examples for our tool. If for a given exercise or problem, our tool inferred a SQL query that does not behave as we expected when applied to other inputs, we will manually find an input on which the SQL query misbehaves and reapplied our tool to the new input. We repeat this process until our tool infers a desirable SQL query.

B. Evaluation Procedure

[[describe how to perform evaluation]]

Our experiments were run on a 2.67GHz Intel Core PC with 4GB physical memory (2GB was allocated for the JVM), running Windows 7.

C. Results

1) *Success Ratio*: For each SQL synthesis scenario, we measured two results. The time cost in producing a desirable SQL query by our tool, and the time we spent in writing input-output examples. Table 10 summarizes our experimental results.

[[need to edit]] As shown in Table 10, our technique synthesized expected SQL queries for 10 out of 11 scenarios within a very small amount of time (less than 1 minute each). The human effort spent in providing input-output examples is also very limited: it took us less than 5 minutes for one scenario.

2) *A Real Example*: We next describe a real SQL problem picked up from an online SQL help forum². The thread was started by a novice user, who needed help to write a SQL query to get result from three input tables. In the help thread, the novice user described his required query in a few paragraphs of English, but also include several small, representative input-output examples as shown in Figure ??, to better express his intention. This thread receives no replies as of May 2012, and we speculated that writing a SQL query to join three tables to produce certain output results can be non-trivial.

We ran our tool on the input-output examples in Figure ?. The tool produced 6 valid answers in less than 1 minutes, all of which satisfy the given examples. The highest ranked SQL query is shown in Figure ??, which is quite unintuitive to write. The SQL query in Figure ?? first joins three input tables on columns T1.Column2, T2.Column2, T2.Column1, and T3.Column1 using some selected columns, and then aggregates the results based on column Table2.Column3’s value. Finally, it returns the minimal values of columns T1.Column1, T1.Column4, and T3.Column2 from each aggregated group as the results.

3) *Performance*: We measured SQLSynthesizer’s performance by recording the average time cost in producing a ranked list of SQL queries. As shown in Figure ??, the performance of SQLSynthesizer is reasonable. On average, it uses less than XX minutes to produce the results in one iterative round. Most of the time is spent querying the backend database to validate the correctness of each synthesized SQL query. **[[caching might be helpful]]**

4) *Number of Interactive Rounds*: This is a measure of the generalization power of the conditional learning part of the algorithm and the ranking scheme. We observed that the tool typically requires just **[[XXX]]** round of interaction, when the user is smart enough to give an example for each input format (which typically range from 1 to 3) to start with. This was indeed the case for most scenarios in our benchmarks, even though our algorithm can function robustly without this assumption. The maximum number of interactive rounds required in any scenario was **[[XXX]]** (with 2 to 3 being a more typical number). **[[the largest table]]** The maximum number of examples required in any scenario over all possible interactions was 10.

²<http://forums.tutorialized.com/sql-basics-113/join-problem-147856.html>

ID	Source	Short description
1	Textbook	Find all students that registrate
2	Textbook	
3	Textbook	
4	Textbook	
5	Textbook	
6	Textbook	
7	Textbook	
8	Textbook	
9	Textbook	
10	Textbook	
11	Textbook	
12	Textbook	
13	Textbook	
14	Textbook	
15	Textbook	
16	Textbook	
17	Textbook	
18	Textbook	
19	Textbook	
20	Textbook	
21	Textbook	
22	Textbook	
23	Textbook	
24	Textbook	
25	Textbook	

Fig. 10. lack a title.

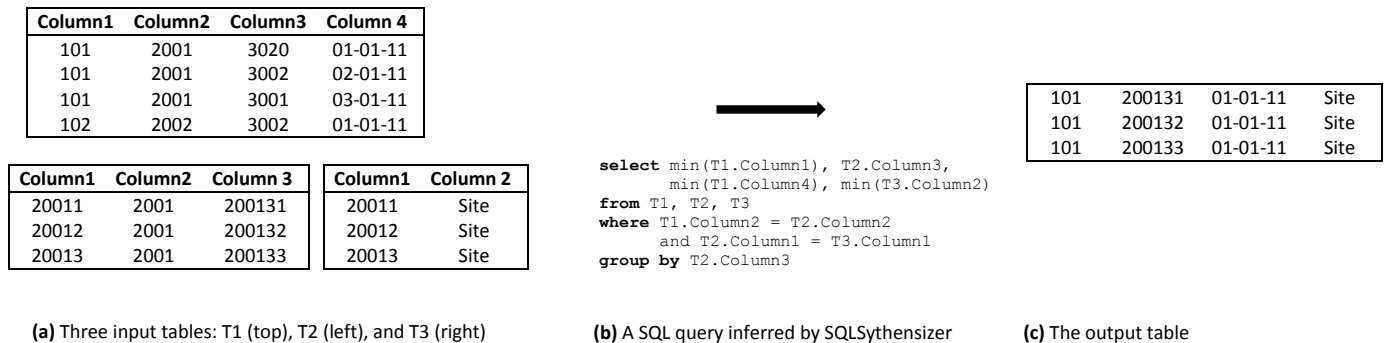


Fig. 11. Input-output examples ((a) and (c)) taken from an online SQL help forum thread. SQLSynthesizer automatically synthesizes 6 SQL queries that can produce the output table from the three input tables. (b) shows the highest ranked SQL query.

5) *Comparison with a Previous Technique:* . **[[not ready yet]]**

D. Initial Experience

We deployed the beta-test version of SQLSynthesizer to a small number of developers and have been using it ourselves, and refining it, since early May 2012. Designing and deploying SQLSynthesizer, along with feedback from the handful of users, has helped us to better understand the issues and to improve the tool's design.

Here is one example piece of feedback from an external user, via private communication:

[[some quote from users]]

Prior to developing SQLSynthesizer, we studied over 100 forum posts about common problems about writing correct SQL queries. **[[most of them use examples]]** We anticipate that future user studies will identify additional strengths and weaknesses that will allow us to further improve SQLSynthesizer.

E. Experimental Discussion

Limitations. The experiments indicate several limitations of our technique. First, XXX Second, XXX.. **[[fill above]]**

Threats to Validity. There are three major threats to validity in our evaluation. First, the XXX textbook exercises and XXX forum problems may not be representative. Thus, we can not

claim the results can be generalized to an arbitrary scenario. Second, we only compared SQLSynthesizer with **[[comparison results]]**. Using other query inference and recommendation techniques **[]** might achieve different results. Third, **[[about user study]]**

Experimental Conclusions. We have three chief findings: (1) The supported SQL subset in SQLSynthesizer is expressive enough to describe a variety of database queries. (2) SQLSynthesizer efficiently synthesizes desirable SQL queries on **[[how many]]** for both textbook exercises and online forum problems, with a small amount of human effort small input-output examples. (3) SQLSynthesizer infers more SQL queries than existing techniques **[]**. **[[edit]]**

VII. RELATED WORK

This section discusses two categories of closely-related work on on reverse query processing methods and automated program synthesis techniques.

A. Reverse Query Processing

In the database community, there is a rich body of work [5], [27], [30] that share the same broad principle of “reverse query processing”. Zloof’s work on *Query by Example* (QBE) [30] provided an intuitive form-based Graphical User Interface for writing database queries, which is completely different than ours. **[[add why it is different]]**

[[mention the difference]] Tran et al. [27] studied the problem of *Query By Output* (QBO). QBO shares a similar goal with QBE, but is significantly different than SQLSynthesizer. In QBO, we are given a view: V and n relations: R_1, R_2, \dots, R_n and asked to return a query Q such that $Q(R_1, R_2, \dots, R_n) = V$. **[[what’s the difference between?]]** However, QBO is cast as a classification problem (into two classes: tuples in V , and tuples not in V) and decision trees are used to find a compact query. A decision tree is constructed in a greedy fashion by determining a “good” predicate to split the tuples into two classes, which form the root nodes of two decisions trees (each of which is constructed in a recursive fashion). The QBO technique presented in [27] infers select-project-join queries. Such queries can only be applied to satisfy a limited number of examples, and cannot be applied to many of the other real world cases that we studied (e.g., queries with aggregation operators such the examples in Figure 1 and 11).

Recently, Sarma et al. [5] proposed a technique to address the *View Definitions Problem* (VDP). **[[explain more here]]** With a rather different goal than SQLSynthesizer, VDP aims to find the most succinct and accurate view definition, when the view query is restricted to a specific family of queries. VDP can be treated as a special case of QBO where $n = 1$ and there are no joins nor projections. Compared to SQLSynthesizer, there are three key differences: (1) VDP aims to infer a relation from a large representative example view, while SQLSynthesizer attempts to infer a SQL query from a set of few input-output examples **[[revise below]]** (which is a critical usability aspect for end-users who lack adequate SQL experience). (2) the technique proposed by Sarma et al. [5] does not consider join or

projection operations and assumes the output view to be a strict subset of the input database with the same schema**[[revise]]**, while our work eliminates these assumptions, and considers both join and projection operations.**[[revise above]]**

In the data mining community, a slightly related problem, called the *redescription mining problem* [23], has been recently studied by researchers. **[[explain the problem]]** However, SQLSynthesizer is significantly different from existing redescription mining techniques **[]** in several ways. First, the goal of redescription mining is to find different subsets of data that afford multiple descriptions across multiple vocabularies covering the same dataset, while our goal is to infer a SQL query to achieve the desirable output. Second, we are concerned with a fixed subset of the data (i.e., the output example table). **[[explain]]** Third, none of the existing redescription mining techniques account for structural (relational) information in the data. **[[revise]]**

[[what is the limitation of non program synthesis, why sythensis would be useful.]]

B. Automated Program Synthesis

Program synthesis aims to create an executable program in some underlying language from specifications that can range from logical declarative specifications to examples or demonstrations [12], [15], [25]. This topic has been studied intensively in Artificial Intelligence and Programming Language communities with the goal of easing the burden of algorithm designers, software developers, or end-users. Interest in it has grown rapidly in recent years [11], and many approaches have been proposed. **[[need to inspect]]** Closely-related to SQLSynthesizer, example-based program synthesis techniques [1], [2], [7], [8], [17], [19], [20] take textual input-output examples and infer programs that transform text. However, existing techniques focused on synthesizing programs for string expressions [12], [25], spreadsheet layout transformation [15], , geometry constructions [13], thus can not be directly applied to our SQL query synthesis domain, due to different abstractions, granularities, and tasks. **[[explain the difference between excel work]]**

There is few work for inferring SQL queries for database end-users. Query recommendation systems, such as SQLShare **[]** and SQLSnippet **[]**, **[[explain]]** may achieve a similar goal of helping end-users write better SQL queries, but they rely on information that we cannot assume access to in many realistic scenarios: a query log [18], or user history and preferences [16]. **[[inferring SQL from code, for performance]]** Add this as related work [3]

VIII. CONCLUSIONS AND FUTURE WORK

[[move to somewhere else]] General purpose programming platforms have never been easy for end-users who are not professional programmers. Users are often stucked with the process of *how* to accomplish a certain task by receiving step-by-step, detailed, and syntactically correct instructions, instead of simply describing *what* the task is.

This paper studied the problem of automated SQL query synthesis from simple input-output examples, and presented a practical technique (and its tool implementation, called SQLSynthesizer) that can infer desirable SQL queries that fulfill end-users' intentions. Our experimental results show that SQLSynthesizer is effective in automating a variety of database query tasks, and it does so with small input-output examples.

The source code of our tool implementation is available at: <http://sqlsynthesizer.googlecode.com>

Our future work will concentrate on the following topics:

Enrich the supported SQL subset. SQLSynthesizer supports a SQL subset containing some widely-used language features. **[[miss some corner cases]]** We plan to enrich the supported SQL subset, and design a corresponding algorithm to synthesize more general queries.

Illustration of synthesis steps. Besides producing a final result, end-users may also be interested in knowing how a SQL query is inferred step by step. Showing detailed inference steps not only makes SQLSynthesizer more usable, but also permits end-users to better understand the whole process and spot possible errors earlier. We plan to apply recent advance in data visualization [17] to the context of program synthesis.

Noise detection and tolerance in users' inputs. The current technique requires users to provide noise-free input-output examples. Even in the presence of a small amount of user-input noises (e.g., a typo), the inference algorithm will declare failure when it fails to learn a valid SQL query. To overcome this limitation, we plan to design a more robust inference algorithm that can attempt to identify and tolerate user-input noises, and even suggest a fix to the noisy example.

REFERENCES

- [1] A. Arasu, S. Chaudhuri, and R. Kaushik. Learning string transformations from examples. *Proc. VLDB Endow.*, 2(1):514–525, Aug. 2009.
- [2] D. M. Barbosa, J. Cretin, N. Foster, M. Greenberg, and B. C. Pierce. Matching lenses: alignment and view update. In *Proceedings of the 15th ACM SIGPLAN international conference on Functional programming*, ICFP '10, pages 193–204, New York, NY, USA, 2010. ACM.
- [3] A. Cheung, A. Solar-Lezama, and S. Madden. Inferring sql queries using program synthesis. *CoRR*, abs/1208.2013, 2012.
- [4] W. W. Cohen. Fast effective rule induction. In *Proceedings of the 12th International Conference on Machine Learning*, ICML'95, pages 115–123. Morgan Kaufmann, 1995.
- [5] A. Das Sarma, A. Parameswaran, H. Garcia-Molina, and J. Widom. Synthesizing view definitions from data. In *Proceedings of the 13th International Conference on Database Theory*, ICDT '10, pages 89–103, New York, NY, USA, 2010. ACM.
- [6] Database Journal. <http://forums.databasejournal.com/>.
- [7] K. Fisher. Learnpads: Automatic tool generation from ad hoc data. In *In SIGMOD*, 2008.
- [8] K. Fisher, D. Walker, K. Q. Zhu, and P. White. From dirt to shovels: fully automatic tool generation from ad hoc data. In *Proceedings of the 35th annual ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, POPL '08, pages 421–434, New York, NY, USA, 2008. ACM.
- [9] E. Frank and I. H. Witten. Generating accurate rule sets without global optimization. In *Proceedings of the 15th International Conference on Machine Learning*, ICML'98, pages 144–151. Morgan Kaufmann, 1998.
- [10] J. Gray, D. T. Liu, M. Nieto-Santisteban, A. Szalay, D. J. DeWitt, and G. Heber. Scientific data management in the coming decade. *SIGMOD Rec.*, 34(4):34–41, Dec. 2005.
- [11] S. Gulwani. Dimensions in program synthesis. In *Proceedings of the 12th international ACM SIGPLAN symposium on Principles and practice of declarative programming*, PPDP '10, pages 13–24, New York, NY, USA, 2010. ACM.
- [12] S. Gulwani. Automating string processing in spreadsheets using input-output examples. In *Proceedings of the 38th annual ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, POPL '11, pages 317–330, New York, NY, USA, 2011. ACM.
- [13] S. Gulwani, V. A. Korthikanti, and A. Tiwari. Synthesizing geometry constructions. In *Proceedings of the 32nd ACM SIGPLAN conference on Programming language design and implementation*, PLDI '11, pages 50–61, New York, NY, USA, 2011. ACM.
- [14] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. The weka data mining software: an update. *SIGKDD Explor. Newsl.*, 11(1):10–18, Nov. 2009.
- [15] W. R. Harris and S. Gulwani. Spreadsheet table transformations from examples. In *Proceedings of the 32nd ACM SIGPLAN conference on Programming language design and implementation*, PLDI '11, pages 317–328, New York, NY, USA, 2011. ACM.
- [16] B. Howe, G. Cole, N. Khoussainova, and L. Battle. Automatic example queries for ad hoc databases. In *Proceedings of the 2011 international conference on Management of data*, SIGMOD '11, pages 1319–1322, New York, NY, USA, 2011. ACM.
- [17] S. Kandel, A. Paepcke, J. Hellerstein, and J. Heer. Wrangler: interactive visual specification of data transformation scripts. In *Proceedings of the 2011 annual conference on Human factors in computing systems*, CHI '11, pages 3363–3372, New York, NY, USA, 2011. ACM.
- [18] N. Khoussainova, Y. Kwon, M. Balazinska, and D. Suciu. Snipsuggest: context-aware autocompletion for sql. *Proc. VLDB Endow.*, 4(1):22–33, Oct. 2010.
- [19] T. Lau, S. A. Wolfman, P. Domingos, and D. S. Weld. Programming by demonstration using version space algebra. *Mach. Learn.*, 53(1-2):111–156, Oct. 2003.
- [20] T. A. Lau, P. Domingos, and D. S. Weld. Version space algebra and its application to programming by demonstration. In *Proceedings of the Seventeenth International Conference on Machine Learning*, ICML '00, pages 527–534, San Francisco, CA, USA, 2000. Morgan Kaufmann Publishers Inc.
- [21] MySQL. <http://www.mysql.com>.
- [22] J. R. Quinlan. *C4.5: programs for machine learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1993.
- [23] N. Ramakrishnan, D. Kumar, B. Mishra, M. Potts, and R. F. Helm. Turning cartwheels: an alternating algorithm for mining redescrptions. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '04, pages 266–275, New York, NY, USA, 2004. ACM.
- [24] R. Ramakrishnan and J. Gehrke. *Database Management Systems*. Addison-Wesley (3rd Edition), 2007.
- [25] R. Singh and S. Gulwani. Learning semantic string transformations from examples. In *Proceedings of the 37th International Conference on Very Large Data Bases*, VLDB '2012, New York, NY, USA, 2012. ACM.
- [26] StackOverflow. <http://www.stackoverflow.com>.
- [27] Q. T. Tran, C.-Y. Chan, and S. Parthasarathy. Query by output. In *Proceedings of the 35th SIGMOD international conference on Management of data*, SIGMOD '09, pages 535–548, New York, NY, USA, 2009. ACM.
- [28] Tutorialized Forums. <http://forums.tutorialized.com>.
- [29] y. Giulia Pagallo, D. Haussler, and P. Rosenbloom. Boolean feature discovery in empirical learning. *Machine Learning.*, 5(1):71–99, Mar 1990.
- [30] M. M. Zloof. Query-by-example: the invocation and definition of tables and forms. In *Proceedings of the 1st International Conference on Very Large Data Bases*, VLDB '75, pages 1–24, New York, NY, USA, 1975. ACM.