

Automatic Generation and Reranking of SQL-Derived Answers to NL Questions

Alessandra Giordani and Alessandro Moschitti¹

Abstract. In this paper, given a relational database, we automatically translate a factoid question in natural language to an SQL query retrieving the correct answer. We exploit the structure of the DB to generate a set of candidate SQL queries, which we rerank with a SVM-ranker based on tree kernels. In particular we use linguistic dependencies in the natural language question and the DB metadata to build a set of plausible SELECT, WHERE and FROM clauses enriched with meaningful joins. Then, we combine all the clauses to get the set of all possible SQL queries, producing candidate queries to answer the question. This approach can be recursively applied to deal with complex questions, requiring nested SELECT instructions. We sort the candidates in terms of scores of correctness using a weighting scheme applied to the query generation rules. Then, we use a SVM ranker trained with structural kernels to reorder the list of question and query pairs, where both members are represented as syntactic trees. The f-measure of our model on standard benchmarks is in line with the best models (85% on the first question), which use external and expensive hand-crafted resources such as the semantic interpretation. Moreover, we can provide a set of candidate answers with a Recall of the answer of about 92% and 96% on the first 2 and 5 candidates, respectively.

1 Introduction

In the last decade, a variety of approaches have been developed to automatically convert natural language questions into machine-readable instructions. In the area of databases, question answering (QA) systems are supposed to answer to natural language questions by executing one or more SQL queries. This is obviously a complex task as systems have to deal with the lexical gap between natural language expressions and database structure. In this paper, we will demonstrate that it is possible to fill such gap by relying on (i) the informative metadata embedded in all real databases, (ii) natural language processing methods, e.g., syntactic parsing, and (iii) advanced machine learning to build kernel-based rerankers.

When designing a database, domain experts are requested to organize entities and relationships naming tables and columns in a meaningful way (i.e. *state_name* or *capital* instead of *table_1* or *table_2*). Moreover the database schema also specifies constraints and data types. This metadata is stored in an underlying database that contains tables of each database. The latter, in turn, contain columns referring to table names and column names. Such logic organization is referred to as *catalog*, and in SQL systems it is stored in a database called INFORMATION_SCHEMA (IS for brevity). A fragment sample is shown in Figure 1. IS can be inspected as a normal database, posing SQL

queries to obtain useful fields to build a new SQL query. In practice, we can use the same technique and technology to generate an answer to a given question and retrieve the answer.

This approach can also deal with cross-domain questions, as long as IS embeds shared metadata between multiple databases. For example, if we have both GEOQUERY and SAKILA data in the same database systems, we can find an answer for cross-domain questions like “Which movies were recorded in major cities of Texas?”.

In addition instead of using tailored dictionaries, we can enrich our knowledge based on the metadata added by the domain expert, when designing the database. Of course, it will be essential to rely on WordNet and similarity measure to generalize such metadata. For example, an answer for the question “Which rivers run through New York” can be found in the GeoQuery corpus. This is associated with a spatial database whose structure is stored in IS as shown in Figure 1.

While we have a simple matching for the word *rivers* with table *river* and column *river_name*, there isn’t a direct mapping between the word *run* in the question and any of the columns in the metadata. However, the disambiguation of the term *run* can be easily performed by looking at the less semantically distant metadata entry, i.e., *traverse*. This matching is re-confirmed when investigating on all possible interpretations of *New York* in this database (i.e. *city_name*, *state_name*, etc.), by the existing reference between column *traverse* in table *river* and column *state_name* in table *state*.

However, a link between both words *New* and *York* is not so easy, since there is no evidence of relatedness between the two words in the metadata: this means that the whole database should be looked up for their stems. Words can be matched with lots of values (e.g., “New York” both as city and as state name, but also with “New Jersey”), as shown by Figure 2. We can generate all possible (even ambiguous) queries exploiting related metadata information (i.e. primary and foreign keys, constraints, datatypes, etc.) and select the most plausible one using a re-ranker.

Last but not least, we deal with complex natural language (NL) questions, containing subordinates, conjunctions and negations and nested SQL queries. In particular, we designed a mapping algorithm that matches dependencies between NL components and SQL structure that allows to build a set of possible queries that answers a given question. This question answering problem and our proposed solution are described in detail later on in the paper. Section 2 gives a formal description of the problem while Section 3 describes the basic steps of our algorithm used to build clause. Section 4 shows how we prune and weigh queries in their possible combinations to generate an ordered set of meaningful queries among which we find the answer. Section 5 describes tree kernels our kernel-based rerankers. Section 6 discusses the results obtained using a reranking algorithm, while Section 7 draws some conclusions.

¹ Department of Computer Science and Engineering, University of Trento, Italy, email: agiordani@disi.unitn.it

TABLES			COLUMNS				
TABLE_SCHEMA	TABLE_NAME	...	TABLE_SCHEMA	TABLE_NAME	COLUMN_NAME	DATA_TYPE	...
geoquery	state		geoquery	state	state_name	varchar	
geoquery	city		geoquery	state	population	float	
geoquery	river		geoquery	city	city_name	varchar	
geoquery	border		geoquery	city	state_name	varchar	
geoquery	highlow		geoquery	river	traverse	varchar	
...	
sakila	city		sakila	city	city	varchar	
...	

KEY_COLUMN_USAGE						
TABLE_SCHEMA	TABLE_NAME	COLUMN_NAME	REFERENCED_TABLE_SCHEMA	REFERENCED_TABLE_NAME	REFERENCED_COLUMN_NAME	...
geoquery	city	state_name	geoquery	state	state_name	
geoquery	river	traverse	geoquery	state	state_name	
...	
geoquery	city	city_name	sakila	city	city	
...	

Figure 1. A DBMS catalog containing GEOQUERY and SAKILA

CITY				RIVER		
CITY_NAME	STATE_NAME	POPULATION	...	RIVER_NAME	TRAVERSE	...
new york	new york	7071640		delaware	new york	
newark	new jersey	329248		delaware	new jersey	
...				alleggheny	new york	
				hudson	new york	
				hudson	new jersey	
				...		

STATE			
STATE_NAME	CAPITAL	POPULATION	...
new york	albany	17558000	
new jersey	trenton	7365000	
...			

Figure 2. GEOQUERY database fragment

2 The Problem

We will begin by introducing the notion of typed dependencies and how to obtain a collapsed list of dependencies starting from an NL sentence. Then we will introduce the subset of Structured Query Language that our system can deal with and, in order to formalize the problem, we will recall the notation of corresponding operations in relational algebra.

2.1 NL Questions and Dependencies List

To represent the textual relationships of the NL sentence we use typed dependency relations. The Stanford Dependencies representation [9] provides a simple and consistent description of the binary grammar relations existing between a governor and a dependent. As shown in the example below, each dependency is written as *abbreviated_relation_name* (governor, dependent). The governor and the dependent are words in the sentence associated with a number indicating the position of the word in the sentence.

In particular we refer to collapsed representation, where dependencies involving prepositions, conjuncts, as well as information about the referent of relative clauses are collapsed to get direct dependencies between content words.

For example, the Stanford Dependencies Collapsed (*SDC*) representation for the question, q_1 : “What are the capitals of the states that border the most populated state?” is the following:

$SDC_{q_1} = attr(are-2, what-1), root(ROOT-0, are-2), det(capitals-4, the-3), nsubj(are-2, capitals-4), nsubj(border-9, states-7), rcmmod(states-7, border-9), det(states-13, the-10), advmod(populated-12, most-11), amod(state-13, populated-12), dobj(borders-9, state-13)$

The current representation contains approximately 53 grammatical relations but for our purposes we only use the following: adverbial and adjectival modifier, agent, complement, object, subject, relative clause modifier, prepositional modifier, and root.

2.2 SQL queries and Relational Algebra

The general SQL query with which our system can deal has the following form:

$SELECT\ COLUMN\ FROM\ TABLE\ [WHERE\ CONDITION]$ (1)

The query is interpreted starting from the relation in the FROM clause, selecting tuples that satisfy the condition indicated in the WHERE clause (optional) and then projecting the attribute in the SELECT clause.

In relational algebra, selection and projection are performed by σ and π operators respectively. The meaning of the SQL query above is the same as that of the relational expression:

$\pi_{COLUMN}(\sigma_{CONDITION}(TABLE))$ (2)

It is worth noting that while relational algebra formally applies to sets of tuples (i.e. relations), in a DBMS relations are bags so it may contain duplicate tuples [4]. For our purposes the fact of having duplicates in the result adds noise; this is why we always delete multiple copies of a tuple by using the keyword DISTINCT in the COLUMN field. In our QA task we expect that questions can be answered with a single result set (e.g. we can deal with “Cities in Texas” and “Populations in Texas” but not with the combined query “Cities and their population in Texas”). That is, even if in general COLUMN could

be a - possibly empty - list of attributes, in our system it just contains one attribute. We can apply to this attribute aggregation operators that summarize it by means of SUM, AVG, MIN, MAX and COUNT, always combined with DISTINCT keyword (e.g. `SELECT COUNT(DISTINCT state.state_name)`).

Instead, *CONDITION* is a logical expression where basic conditions, in the form $e_L \text{ OP } e_R$, with $\text{OP}=\{<, >, \text{LIKE}, \text{IN}\}$, are combined with AND, OR, NOT operators. While e_L is always in the form `table.column`, e_R could be:

- numerical value (e.g. `city.population > 15000`) or
- string value (e.g. `city.state_name LIKE "Texas"`) or
- nested query (e.g. `city.city_name IN (SELECT state.capital FROM state)`)

An example of a complex WHERE condition could be the following: `city.population > 15000 AND city.city_name NOT IN (SELECT state.capital FROM state) AND NOT city.state_name LIKE "Texas"` (i.e. “*major non-capital cities excluding texas*”).

The meaning of *TABLE* is more straightforward, since it should contain table name(s) to which the other two clauses refer. This clause could just be a single relation or a join operation, which selectively pairs tuples of two relations. We only deal with theta-joins where we take the Cartesian product of two relations and exclusively select those tuples that satisfy a condition C . The notation for theta-joins of relations R and S based on condition C is $R \bowtie_C S$. We use the SQL keyword ON to keep this condition C separated from the other WHERE conditions since it reflects a database requirement and shouldn't match to anything of the NL question. (e.g. `city JOIN state ON city.city_name = state.capital`).

The complexity of generated queries is fairly high indeed, since we can deal with questions that require nesting, aggregation and negation in addition to basic projection, selection and joining (e.g. “*How many states have major non-capital cities excluding Texas*”).

2.3 Problem Definition

The question answering task of finding an SQL query that retrieves an answer for a given NL question reduces to the following problem.

Given a question q represented by means of one typed dependency collapsed list SDC_q , generate the three sets of clauses $\mathcal{S}, \mathcal{F}, \mathcal{W}$ (argument of SELECT, FROM and WHERE, respectively) such that:

$$\exists s \in \mathcal{S}, \exists f \in \mathcal{F}, \exists w \in \mathcal{W} \text{ s.t. } \pi_s(\sigma_w(f)) \text{ answers } q \quad (3)$$

The query $answer \pi_s(\sigma_w(f))$ is chosen among the set of all possible queries $\mathcal{A} = \{\text{SELECT } s \times \text{FROM } f \times \text{WHERE } w\}$ in a way that maximizes the probability of generating a result set answering question q .

3 Building Clauses Sets

In order to generate all possible queries for a question q we need to find their possible SELECT, FROM and WHERE clauses (\mathcal{S}, \mathcal{F} and \mathcal{W}). We start from a dependency list SDC_q and (a) prune and stem its components, (b) add synonyms, (c) create the set of stems used to build \mathcal{S} and \mathcal{W} and (d) keep only dependencies possibly used in the recursive step to generate nested queries. Building the set \mathcal{F} from \mathcal{S} and \mathcal{W} is straightforward.

We are now going to briefly discuss some examples to introduce the objective of individual steps and clarify how the entire process is

carried out. The first question we take into account is the simplest one: “*What is the capital of Texas?*”. Its answer can be retrieved executing the query: `SELECT capital FROM state WHERE state.state_name='Texas'`. We can see that they share only two stems, *capital* and *Texas*. The key of categorizing stems (Section 3.2) is to recognize that the first stem will be used in \mathcal{S} and the second one in \mathcal{W} . In particular, since the word *Texas* is not a value in the \mathcal{S} , it is used as a r-value in the WHERE expression, while the l-value is derived from the column name under where it appears (Section 3.4).

The fact of being respectively projection and selection oriented can be inferred looking at their grammar relations, i.e. inspecting the dependency list (e.g. root of the sentence together subject dependent are typically used for projections). This list needs to be preprocessed (section 3.1) to take into account only relevant relations between the *stems* of the question. Let us consider for example the question: “*What is the capital of the most populous state?*” and its associated answering query `SELECT capital FROM state WHERE population = (SELECT max(population) FROM state)`. The matching words are *capital* and *state*, while stemming also allows to find a mapping through *popul*. We can note that this stem is used both in the l-value and in the r-value of the WHERE expression. In fact, this query requires nesting and indeed the categorizing algorithm needs to be recursive. This stem is classified both as a selection oriented stem for the outer query, and as a projection oriented one for the inner query (note that it requires aggregation, handled when generating the SELECT clause set, see Section 3.3).

Finally we will introduce one last example to clarify Section 3.5. While with the other examples it is straightforward to compile the FROM clause, since the other clauses refer to the same table, when we deal with columns belonging to different tables things get complicated. Take question “*What are the capitals states bordering Texas?*” and its associated query `SELECT capital FROM ... WHERE border = 'Texas'`. How can we fill in the dots in the FROM clause? Fields *capital* and *border* belong respectively to tables *state* and *border.info*. From the database catalog, we learn that these two tables are connected via the foreign key *state_name* and so the final \mathcal{F} will include the following join: `state JOIN border.info on state.state_name = border.info.state_name`.

3.1 Optimizing the Dependency List

As introduced in Section 2.1, we don't need all grammatical relations provided in output by the Stanford Dependency parser. For this reason before preprocessing the list of dependencies we need to prune the useless ones and remove from *governors* and *dependents* the appended number (indicating the position of the word in question q). Then, *govs* and *deps* are reduced to stems (using the Porter stemmer²).

In order to disambiguate the sense of the stems that do not appear in metadata but could match with it, we create a list of synonyms using off-the-shelf resources (like Wordnet and similarity measures) combined with our internal knowledge (represented by database constraints). Using this list we can substitute certain stems with their stemmed synonyms.

The resulting SDC_q is optimized to be processed by the next step. An example showing $SDC_{q_1}^{opt}$ with respect to the original SDC_{q_1} introduced in Section 2.1 can be found in Table 1.

3.2 Categorizing Stems

Before building \mathcal{S} and \mathcal{W} sets we need to identify those stems that are projection and/or selection oriented. Those stems will be added re-

² <http://tartarus.org/martin/PorterStemmer/>

(1) *root*(*ROOT*, *are*),
(2) *nsubj*(*are*, *capital*),
(3) *prep_of*(*capital*, *state*),
(4) *nsubj*(*border*, *state*),
(5) *rmod*(*state*, *border*),
(6) *advmod*(*populat*, *most*),
(7) *amod*(*state*, *populat*),
(8) *doobj*(*border*, *state*)
 $\Pi = \{\text{capital, state}\}$
 $\Sigma = \{\text{are}\} \Rightarrow \Sigma = \phi$
 $\Pi' = \{\text{state, border}\}$
 $\Sigma' = \{\text{border, state}\}$
 $\Pi'' = \{\text{most, populat, state}\}$
 $\Sigma'' = \phi$

Table 1. Categorizing stems into projection and/or selection oriented sets

spectively to Π and/or Σ categories according to the following rules. For each grammatical relation $rel(gov, dep)$ in SDC_q^{opt} :

1. If it is *ROOT*, *dep* is the key to populate \mathcal{W} so add it to Σ and remove the relation from SDC_q^{opt} . This stem can be an auxiliary verb, e.g., *is*, *are*, *has*, *have* and so on. It is useless to build the arguments of the queries but it could be used transitively to add other stems³.
2. If it starts with *nsubj*, check if $gov \in \Sigma$. If not (because there isn't any *ROOT* relation) add *gov* to Σ . Then add *dep* to Π and remove *rel* from SDC_q^{opt} , otherwise keep it, since it could be a subject related to a subordinate (we will need it in the recursive steps).
3. If it starts with *prep* or it ends with *obj*, we used it to create conditions (possibly involving nesting):
 - check if $gov \in \Pi$. If not (because no *ROOT* or *nsubj* relations were found so far) add *gov* to Π .
 - Then add *dep* to Σ if there is not any *table.column* like ⁴ *gov.dep*. Otherwise, also add *dep* to Π and remove *rel* from SDC_q^{opt} .
4. If it ends with *mod*, it implies that *dep* is a modifier of *gov*, so they should be paired together: if $gov \in \Sigma$ add *dep* to Σ and if $gov \in \Pi$ add *dep* to Π and remove *rel* from SDC_q^{opt} . This should be done only if *dep* is not a superlative (i.e. doesn't end with -st). The non-removed relations will be taken into account in the recursive step, adding both *dep* and *gov* to Π .
5. If none of the above rules can be applied, iterate the algorithm recursively building Π' and Σ' , Π'' and Σ'' and so on, until SDC_q^{opt} is empty.

In order to show how these steps are used to build projection and/or selection oriented sets from which we generate \mathcal{S} and \mathcal{W} , let us consider the list of optimized dependencies SDC_q^{opt} in Table 1.

At the first iteration we use *ROOT* to add *are* to Σ , then we also exploit it to add *capital* and include *state* to Π as soon as we check that there is an occurrence *state.capital* in IS. At this point these three relations have been deleted from SDC_q^{opt} obtaining $SDC_q^{opt'}$ used in the next iteration. Note that since *are* is a short stem, it should be deleted from Σ .

³ Stems of 3 or less characters would introduce too much noise in retrieving matching strings, so they will be eliminated in an additional step 6. Useful words like *in*, *of*, *not*, *or*, and are embedded in relation abbreviations when collapsing dependencies.

⁴ We query metadata seeking for something similar to *gov* as a table and to *dep* as a column, i.e. we search for table names using $\pi_{table.name}(\sigma_{table.name \cong dep \wedge column.name \cong gov}(IS.Columns))$. For brevity we use the symbol $s_1 \cong s_2$ for s_2 substring of s_1 , i.e. s_1 LIKE "%s2%".

$\mathcal{S} = \{state.capital^3, state.state_name^2, border_info.state_name^1, \dots\}$
 $\mathcal{S}' = \{border_info.state_name^3, border_info.border^2, state.state_name^2, \dots\}$
 $\mathcal{S}'' = \{max(state.population)^4, max(city.population)^3, state.population^3, \dots\}$

Figure 3. A subset of SELECT clauses for q_1

At the second iteration (first recursion step) we don't have a *ROOT* relation so we use *nsubj* to add *border* to Σ' and *state* to Π' . Since with *rmod* we find an occurrence *border.state_name* in IS, *border* is added also to Π . At this point, seeking through the end of the list we discard *doobj* because even if $border \in \Pi'$ we do not find *state.border* in IS, so these other three relations are deleted from $SDC_q^{opt'}$ obtaining $SDC_q^{opt''}$ for the last iteration.

In the third iteration we have $SDC_q^{opt''}$ composed by two *mod* relations, so we add all stems to Π'' and delete their associated relations from the list.

3.3 Building the SELECT Clauses Set

Once we have identified the set Π of projection-oriented stems, we can use it to search in metadata all the fields that could match with them. The generation process for \mathcal{S} is described by the following generative grammar.

$\mathcal{S} \rightarrow \text{AGGR } '(' \text{ FIELD } ')'$ | FIELD
AGGR \rightarrow max | min | sum | count | avg
FIELD \rightarrow TAB.COL
TAB $\in \bigcup_{x \in \Pi} \pi_{table.name}(\sigma_{table.name \cong x}(IS.Tables))$
COL $\in \bigcup_{x \in \Pi} \pi_{column.name}(\sigma_{column.name \cong x}(IS.Columns))$

With each element of \mathcal{S} , we also associate a weight w_i , calculated according to the procedure described in Section 4.3 (we will discuss it later). For example, considering the IS scheme in Figure 1, the SELECT clauses originated from Π of Table 1 are shown in Fig. 3. Note that the superscript numbers indicate the weight associated with each statement.

3.4 Building the WHERE Clauses Set

Before generating WHERE clauses, the selection-oriented set of stems Σ should be divided into two distinct sets: Σ_L and Σ_R .

The set Σ_L contains stems that find their matching in IS and allow us to build the set of left-hand side expressions $\mathcal{W}_L \rightarrow \text{FIELD}^{w_i}$, where FIELD is defined above and computed with Σ_L in place of Π (w_i is its associated weight).

For the remaining stems $\Sigma_R = \Sigma - \Sigma_L$ we should look up in the database to find a match⁵: $\forall col \in IS.Columns, \forall tab \in IS.Tables$ we generate the set $\mathcal{W}_R = \{x | \pi_{count(*)}(\sigma_{col \cong x}(Geoquery.tab)) \geq 0\}$.

Then, in order to build the WHERE clause set, \mathcal{W} , $\forall e_L \in \mathcal{W}_L, \forall e_R \in \mathcal{W}_R$ we first generate basic expressions $expr = e_L$ OP e_R and combine them by means of conjunctions and negations (see Section 2.2), keeping only those expressions $expr$ such that the execution of $\pi_{count(*)}(\sigma_{expr}(table))$ does not lead to an error for at least a *table* in the database.

To understand how it works, let us introduce a new example question q_2 : "what are the capitals of states bordering New York?". The SDC_q^{opt} is similar to SDC_q^{opt} except for the last three relations. Row (6) disappears while rows (7) and

⁵ Non-matching stems may semantically match a whole condition and need to be handled carefully. For example, *major*, if associated with *city* is translated into '*city.population > 15000*' while when talking about river is associated with '*river.length > 750*' [2]

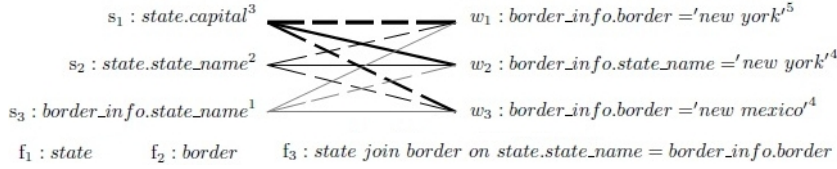


Figure 4. Possible pairing between clauses for q_2

(8) are replaced by $amod(york, new)$ and $doj(border, york)$, leading to $\Sigma' = \{border, new, york\}$. This set is split into $\Sigma'_L = \{border\}$ and $\Sigma'_R = \{new, york\}$. We build $\mathcal{W}'_L = \{border_info.border^3, border_info.state_name^2\}$ and $\mathcal{W}'_R = \{'new\ york'^2, 'new\ mexico'^1, 'new\ jersey'^1, 'newark'^1\}$. Finally we generate the set of possible valid conditions and their weights: $\mathcal{W} = \{border_info.border = 'new\ york'^5, border_info.state_name = 'new\ york'^4, \dots\}$.

Anyway, the set Σ_R could happen to be empty. For example, when the WHERE condition requires nesting: in this case e_R will be the whole subquery (e.g. Σ' in Table 1). It could be the case that also Σ_L is empty. In fact a query without a WHERE clause is valid (e.g. Σ'' in Table 1). In any case, even if there are no selection-based stems, \mathcal{W} may not be empty (e.g. Σ in Table 1). Taking into account all tables and columns we can get more conditions: $\mathcal{W}^*_R = \{tab.col\ such\ that\ tab \in \pi_{table.name}(IS.Columns)\ and\ col \in \pi_{column.name}(IS.Columns)\}$.

3.5 Building the FROM Clauses Set

The generation of the FROM clause \mathcal{F} is straightforward given \mathcal{S} and \mathcal{W} . This set will contain all tables to which clauses in \mathcal{S} and \mathcal{W} refer, enriched by pairwise joins.

As stated before, this information can be found running SQL queries over IS exploiting metadata stored in table KEY_COLUMN_USAGE (in short, K; see Figure 2). This table identifies all columns in the current databases that are restricted by some unique, primary key, or foreign key constraint. That is, for each usage of foreign key column in the table, we can determine how many aggregate table columns match that column usage.

First, we extract tables appearing in \mathcal{S} and \mathcal{W} (i.e. words ending with dot), creating a set F . At the beginning $\mathcal{F} = F$. Then $\forall t_1, t_2 \in F$ $\pi_{col.name, ref.col.name}(\sigma_{table.name=t_1 \wedge ref.table.name=t_2}(IS.K))$ retrieves c_1, c_2 to perform the: $join_{c_1=c_2}^{t_1 \triangleright t_2}$. In this way \mathcal{F} is enriched with the two-table join $t_1 \text{ join } t_2 \text{ on } t_1.c_1 = t_2.c_2$. In addition we can allow for more distant joins by finding an intermediate table useful to link two tables that are not directly referencing each other. This can be done performing a complex join between two instances of KEYS with multiple conditions, but due to for lack of space this can not be illustrated here.

With respect to our example with question q_1 and its SELECT clauses shown in Figure 3, the set of FROM clauses is:

$\mathcal{F}' = \{state, border, state\ join\ border\ on\ state.state_name = border_info.border, \dots\}$.

Note that there are no weights associated with FROM clauses because it is not possible to backtrack how many stems made each table appear in \mathcal{F} .

4 Generating Queries

In the previous section we saw how to create building blocks for queries starting from a question q . These elements should be paired together in a smart way to generate the set of queries that possibly answer q . This pairing is obtained by creating the Cartesian product

between clauses sets from which non-valid, redundant and meaningless clauses are deleted. We use a weighting scheme to order the most probable correct candidate queries.

4.1 Clause Cartesian Product

In order to find possible answering queries we generate the set $\mathcal{A} = \{\mathcal{S} \times \mathcal{F} \times \mathcal{W}\} \cup \{\mathcal{S} \times \mathcal{F}\}$. Given that at least one such query exists there should be one pairing $\langle s, f, w \rangle \in \mathcal{A}$, such that the execution of `SELECT s FROM f [WHERE w]` retrieve the correct answer. Given that each clause set contains on average up to ten items, this product can result in a very huge set. Thus, when generating all pairings some preliminary conditions are verified, e.g. tables appearing in SELECT and WHERE clauses should appear in the FROM clause as well, otherwise the execution of that query will fail. This avoids generating incorrect queries and wasting time trying to execute them.

To give a simple example, we illustrate in Figure 4 some generated clauses for the question q_2 , together with possible pairings. The pairing $\langle s_1, f_1, w_1 \rangle$ is not correct: it leads to the MySQL error `Unknown table: border_info`.

4.2 Pruning Useless Queries

Once the set \mathcal{A} of all valid pairings is built, we additionally prune some of them which are not useful. For example, meaningless queries project the same field compared to a value in the selection (e.g. the pairing $\langle s_3, f_2, w_2 \rangle$ answers the question “Which state is New York?” and is clearly useless).

Moreover there could be redundant queries that, if optimized, allow us to remove duplicates in the set, reducing its cardinality. For example, the pairing $\langle s_2, f_3, w_1 \rangle$ requires the columns $state.state_name$ and $border_info.border$ to be the same, so w_2 would select the same rows of w'_2 (i.e. $state.state_name = 'new\ york'$), but this means that table $border_info$ is no longer used and this pairing is equivalent to $\langle s_2, f_1, w'_2 \rangle$ which, as said above, is meaningless.

4.3 Weighting Scheme

As introduced in the previous sections, we weigh each clause in \mathcal{S} and \mathcal{W} by counting how many stems in the original question originated that clause.

In particular, for the SELECT clause, if there is a table that matches with a stem, its weight is +2 while the matching with columns weighs +1 (common stems between table and column are not valid). Superlatives matching with aggregation operators count as +1.

For the WHERE clause, a weight is computed in the same way as for the left-hand side of the conditions and a +1 is added for each matching value in the right-hand side. In addition when dealing with nested queries, the WHERE clause inherits also the weight of the nested query.

The FROM clauses are not associated with weights. However, we will take into account how many joins are involved when ordering queries with the same weight.

When pairing clauses the total weight is obtained just summing up the weight of its components, and it is used to order the final set $\bar{\mathcal{A}}$ of possible useful queries from the most to the least probable.

Figure 4 highlights this *probabilistic* score (obtained by the heuristic one by normalization) through the thickness of connection lines. Dashed lines illustrate pruned queries. The final ordered set answering q_2 is the following one:

$$\bar{\mathcal{A}} = \{ \langle s_1, f_3, w_2 \rangle^7, \langle s_3, f_2, w_1 \rangle^6, \langle s_2, f_3, w_2 \rangle^6, \langle s_1, f_1 \rangle^3, \langle s_2, f_1 \rangle^2, \langle s_3, f_2 \rangle^1 \}.$$

From the pairing with highest weight we derive the answering query, that is: `SELECT state.capital FROM state join border on state.state_name =border_info.border WHERE border_info.state_name='new york'`.

It is worth noting that more then a query can have the same weight. To deal with that, we implemented a comparator that privileges queries involving less joins and embed the most referenced table (e.g. `state` in the case of `GEOQUERY`). See, for example, the order of the second and third pairings in $\bar{\mathcal{A}}$: they have been swapped since f_3 contains a join while f_2 doesn't.

5 Kernel Methods for Ranking Question/Query Mapping

Once an initial rank of the candidate SQL queries has been derived, we can rely on machine learning methods to improve the probability of finding the correct answer in the top position. The need of designing suitable representations of the question and query pairs makes this operation quite complex. For this purpose, we rely on kernel methods.

5.1 Kernel Methods

In kernel-based machines, both learning and classification algorithms only depend on the inner product between instances. In several cases this can be efficiently and implicitly computed by kernel functions by exploiting the following dual formulation: $\sum_{i=1..l} y_i \alpha_i \phi(o_i) \phi(o) + b = 0$, where o_i and o are two objects, ϕ is a mapping from the objects to feature vectors \vec{x}_i and $\phi(o_i) \phi(o) = K(o_i, o)$ is a kernel function implicitly defining such mapping. In case of structural kernels, K determines the shape of the substructures that describe the objects above.

In the following section, we are going to first propose a structural representation of the question and query pairs, then we will illustrate the Syntactic Tree Kernel (STK) [3], which computes the number of syntactic tree fragments. In the last subsection we will show how to engineer new kernels from them, while the reranking kernel is presented in Sec. 5.5

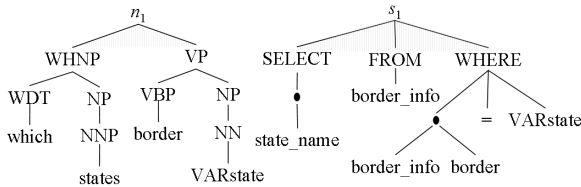


Figure 5. Question/Query Syntactic trees

5.2 Representing Question and Queries Pairs

In Data Mining and Information Retrieval the so-called bag-of-words (BOW) has been shown to be effective to represent textual documents, e.g. [13, 7]. However, in case of questions and queries, we

deal with small textual objects in which the semantic content is expressed by means of few words and poorly reliable probability distributions. In these conditions the use of syntactic representation improves BOW and should be always used.

Therefore, in addition to BOW, we represent questions and queries using their syntactic trees, as shown in Figure 5: for questions (a) we used the Charniak's syntactic parser [1] while for queries (b) we implemented an ad-hoc SQL parser. The latter builds a SQL parse tree for each query following its syntactic derivation according to MySQL grammar. The grammar has been slightly modified to accommodate the usage of the symbol \bullet for the production of *items* in the SELECT clause and in WHERE conditions. In such an SQL tree, the internal nodes are only the SQL keywords of the query plus the special symbol \bullet whereas the leaves are names of tables and columns of the database, category variables or operators. Note that, although we eliminated comma and dot from the grammar, it is still possible to obtain the original SQL query, by just performing a preorder traversal of the tree. The above structures can be represented in a learning algorithm using the kernel described in the next section.

5.3 Syntactic Tree Kernels (STK)

Convolution tree kernels [3] compute the similarity between two trees T_1 and T_2 by counting the common sub-trees, without enumerating the whole fragment space. In more detail, let N_1 and N_2 be the set of nodes in T_1 and T_2 , respectively. Moreover, let $I_i(n)$ be an indicator variable that is 1 if subtree i is rooted at n and 0 otherwise. Then the convolution kernel K over T_1 and T_2 is computed as:

$$STK(T_1, T_2) = \sum_{n_1 \in N_1, n_2 \in N_2} \Delta(n_1, n_2) \quad (4)$$

where

$$\Delta(n_1, n_2) = \sum_{n_1 \in N_1} \sum_{n_2 \in N_2} \sum_i I_i(n_1) I_i(n_2)$$

is computed efficiently using the following recursive definition:

- If the production rules⁶ at n_1 and n_2 are different, then $\Delta(n_1, n_2) = 0$.
- If the production rules at n_1 and n_2 are the same and n_1 and n_2 are pre-terminals, then $\Delta(n_1, n_2) = \lambda$.
- If the production rules at n_1 and n_2 are the same and n_1 and n_2 are not pre-terminals, then:

$$\Delta(n_1, n_2) = \lambda \prod_{j=1}^{nc(n_1)} (1 + \Delta(ch(n_1, j), ch(n_2, j)))$$

where $nc(n_1)$ is the number of children of n_1 in the tree and the j -th children of node n_i is denoted by $ch(n_i, j)$ (note that $nc(n_1) = nc(n_2)$ since the production rule is the same). λ ($0 < \lambda < 1$) is a decay factor to make the kernel less variable with respect to tree-fragment sizes.

5.4 Kernel Combination for Pairs

We need to represent the members of a pair and their interdependencies. For this purpose, given two kernel functions, $k_1(\cdot, \cdot)$ and $k_2(\cdot, \cdot)$, and two pairs, $p_1 = \langle n_1, s_1 \rangle$ and $p_2 = \langle n_2, s_2 \rangle$, a first approximation is given by summing the kernels applied to the components: $K(p_1, p_2) = k_1(n_1, n_2) + k_2(s_1, s_2)$. This kernel will produce the union of the feature spaces of questions and queries. A more effective kernel is the product $k(n_1, n_2) \times k(s_1, s_2)$, since it generates pairs of fragments, which are member of the Cartesian product of kernel spaces of the questions and queries. As additional feature

⁶ In a syntactic tree a node with its children correspond to a production rule of the grammar that generated it.

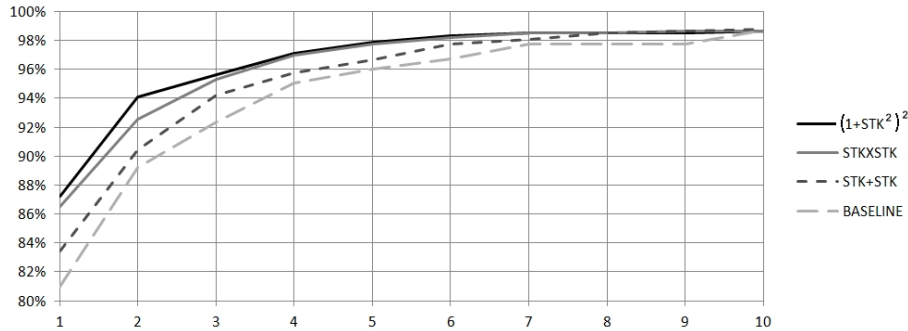


Figure 6. Recall of the correct answer within different k positions of the system rank

and kernel engineering, we also exploit the ability of the polynomial kernel to add feature conjunctions. By simply applying the function $(1 + K(p_1, p_2))^d$, we can generate conjunction up to d features. Thus, we can obtain tree fragment conjunctions and conjunctions of pairs of tree fragments.

The next section will show how to use such kernels for an SVM-based reranker.

5.5 Preference reranker

Our reranking model consists in learning to select the best candidate from a given candidate set. In order to use SVMs for training a reranker, we applied the Preference Kernel Method [14]. In the Preference Kernel approach, the reranking problem – learning to pick the correct candidate h_1 from a candidate set $\{h_1, \dots, h_k\}$ – is reduced to a binary classification problem by creating *pairs*: positive training instances $\langle h_1, h_2 \rangle, \dots, \langle h_1, h_k \rangle$ and negative instances $\langle h_2, h_1 \rangle, \dots, \langle h_k, h_1 \rangle$. This training set can then be used to train a binary classifier. At classification time, pairs are not formed (since the correct candidate is not known), while, the standard one-versus-all binarization method is still applied.

The kernels are then engineered to implicitly represent the *differences* between the objects in the pairs. If we have a valid kernel K over the candidate space \mathcal{T} , we can construct a preference kernel P_K over the space of pairs $\mathcal{T} \times \mathcal{T}$ as follows: $P_K(x, y) =$

$$\begin{aligned} P_K(\langle x_1, x_2 \rangle, \langle y_1, y_2 \rangle) &= K(x_1, y_1) + \\ &K(x_2, y_2) - K(x_1, y_2) - K(x_2, y_1), \end{aligned} \quad (5)$$

where $x, y \in \mathcal{T} \times \mathcal{T}$. It is easy to show that P_K is also a valid Mercer’s kernel. This makes it possible to use kernel methods to train the reranker. The several kernels defined in the previous section can be used in place of K^7 in Eq. 5.

6 The Experiments

We ran several experiments to evaluate the accuracy of our approach for automatic generation and selection of correct SQL queries from NL questions. We experimented with a well-known dataset GeoQuery developed in order to study semantic parsing.

6.1 Setup

To learn the reranker, we used SVM-Light-TK⁸, which extends the SVM-Light optimizer [7] with tree kernels. i.e. Syntactic Tree Kernel (STK) as described in Section 5. We modeled many different combinations described in the next section. We used the default parameters,

⁷ More precisely, we also multiply K for the inverse of rank position.

⁸ <http://disi.unitn.it/~moschitt/Tree-Kernel.htm>

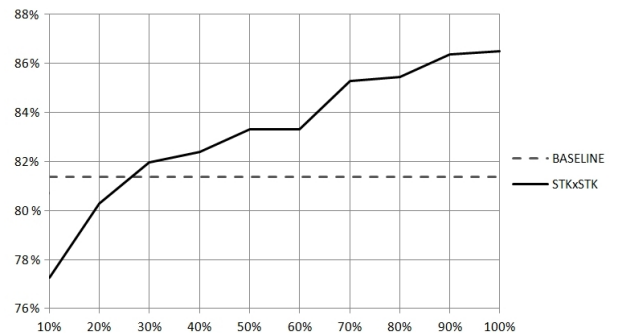


Figure 7. Learning curve comparison between simple answer generator and the reranking model using the $STK \times STK$ kernel.

i.e. the cost and trade-off parameters = 1 (for normalized kernels) and $\lambda = 0.4$ (see Sec. 4).

To generate the set of possible SQL queries we applied our algorithm described in Section 3 to GEOQUERIES⁹ corpus. We started from a set of 700 NL questions¹⁰. Thanks to our generative algorithm we discovered and fixed all errors and inconsistencies in SQL queries, except for 3 cases that still lead to a MySQL error. Indeed, since we can’t test the correctness of our generated query (without a result set to compare with) we considered a subset of 697 pairs.

6.2 Generative Results

Given a question from GeoQuery, our algorithm was able to generate a correct SQL query in the first 25 in 95.3% of the cases. This also means that our system cannot answer to 33 questions. This is due to (1) empty clauses set \mathcal{S} and/or \mathcal{W} , for example, “How many square kilometers in the us?” does not contain any useful stem; and (2) from mismatching nested queries, for example, “Count the states which have elevations lower than what alabama has” contains an implicit reference to a missing piece of question. In addition there are ambiguous questions like “Which states does the colorado?” from which we retrieve an incomplete dependency set.

For all remaining questions from which we succeed in generating an ordered list of possible queries, we find that the query on top of the list retrieves the correct result set in 82% of the cases. For the other questions, it can be found within the first 10 generated answers

⁹ Available at <http://www.cs.utexas.edu/ml/geo.html>

¹⁰ This are the first 700 questions of the 880 ones that Mooney’s group [15] paired with logical formulas in Prolog and that Popescu et al. [11] manually converted into SQL.

Table 2. Kernel combination recall (\pm Std. Dev) for GEO dataset

Combination	Rec@1	Rec@2	Rec@3	Rec@4	Rec@5
NO RERANKING	81.4 \pm 5.8	87.6 \pm 3.8	90.8 \pm 3.1	94.0 \pm 2.4	95.0 \pm 2.0
STK + STK	83.5 \pm 3.6	90.4 \pm 3.5	94.2 \pm 2.9	95.8 \pm 2.0	96.7 \pm 1.7
STK \times STK	86.5 \pm 4.0	92.6 \pm 3.7	95.3 \pm 3.2	97.0 \pm 1.8	97.7 \pm 1.4
(1+STK ²) ²	87.2 \pm 3.9	94.1 \pm 3.4	95.6 \pm 2.7	97.1 \pm 1.9	97.9 \pm 1.4
BOW \times STK	86.7 \pm 4.1	92.1 \pm 3.2	95.6 \pm 2.5	97.1 \pm 1.4	97.6 \pm 1.2

for 99% of the questions (once the 33 questions above have been removed). This can be observed in Figure 6, which plots the Recall (of the correct question) curve of the generative approach, i.e., the baseline. As pointed out in the graphic, the right query is found among the first three in 93% of the cases.

6.3 Reranking Results

Figure 6 also shows the plot for different rerankers using the following kernels: STK+STK, STK \times STK and (1+STK \times STK)², which provide better rankings (the first STK is applied to the question parse trees whereas the second STK is applied to the query derivation tree). For example, the latter kernel retrieved the correct answers 94% of times by only using the first two answers.

To better evaluate the results of our rerankers, we applied standard 10-fold cross validation and measure the average Recall and Std Dev. of selecting a query for each question. The results for different kernel models for reranking are reported in Table 2. The first column of Table 2 lists kernel combination by means of product and sum between pairs of basic kernels used for the question and the query, respectively. The other columns show the percentage of questions for which we found at least 1 correct answer in the top @X positions (average Recall@X over 10 folds \pm Std. Dev).

The results are rather exciting since they compare favorably with the state-of-the-art. The best system on this datasets was designed in [16] and shows a Precision of 96.3% and a Recall of 79.3%, for an f-measure of 86.9%, while our system shows a Precision of 82.8% and a Recall of 87.2%, for an f-measure of 85.0% (when we include the 33 missing questions in the evaluation). Two main facts should be noted:

- our system performs just 2 points less than the system designed in [16] but it does not need any hand-crafted manual resource, i.e., the semantic trees manually designed in [16] for each question, and it is very simple to implement.
- unlike it has been done in previous work, we can also provide multiple ranked answers. If we select the first n candidates, we highly increasing the Recall of the correct answers, e.g., within the first 2 we have a f-measure of 90% (considering the 33 missing questions).

Other closely related work, e.g., [5], suggests that lower results than ours can be obtained using different approaches. These rely either on semantic grammar specified by an expert user [10], or on enriching the information contained in the pairs [11] and implementing ad-hoc rules in a semantic parser [8, 12]. Our system instead, requires no intervention since the database metadata already contains all the needed data.

Finally, we report the learning curve of one basic reranker in Figure 7, showing how recall of STK \times STK increases for larger training sets. The plot reveals that as soon as we provide a reasonable percentage of training data (25% of the available data corresponding to 9 folds of 700 questions – one fold is used for testing) for reranking, the model improves on the baseline.

The main contribution of this research consist in the fact that given a NL question we can generate a set of mapping SQL queries. Moreover if we can rely on a relatively small set of correct pairs of questions and queries to train a SVM classifier, we are able to re-rank the set of generated pairs to select the correct one with a fairly high accuracy.

7 Conclusions and Future Work

In this paper, we have approached the question answering task of implementing a NL interface to databases by automatically generating SQL queries based on grammatical relations and matching metadata. To our knowledge, the underlying idea that we have proposed to build and combine clauses sets is novelty. Additionally, we are firstly experimented with a preference reranking kernel, which is able to boost the accuracy of our generative model.

Given the high accuracy, the simplicity and the practical usefulness of our approach, (e.g., we can generate the correct question in the first 5 candidates in 95% of the cases), we believe that our methods can be successfully used in the future for real-world applications.

In the future we plan to experiment with datasets in different domains (e.g. ATIS corpus). Moreover, given that current challenges in Semantic Web tackle similar problem [6] (scaling question answering approaches to Linked Data, i.e. Question Answering over Linked Data), it would be interesting to apply our algorithms to semantic search and question answering over RDF data.

ACKNOWLEDGEMENTS

The research described in this paper has been partially supported by the European Community’s Seventh Framework Programme (FP7/2007-2013) under the grants #247758: ETERNALS – Trustworthy Eternal Systems via Evolving Software, Data and Knowledge, and #288024: LIMOSINE – Linguistically Motivated Semantic aggregation engines

REFERENCES

- [1] E. Charniak, ‘A maximum-entropy-inspired parser’, in *Proceedings of NAACL’00*, (2000).
- [2] Philipp Cimiano and Michael Minock, ‘Natural language interfaces: What is the problem? - a data-driven quantitative analysis’, in *NLDB*, pp. 192–206, (2009).
- [3] M. Collins and N. Duffy, ‘New ranking algorithms for parsing and tagging: Kernels over discrete structures, and the voted perceptron’, in *Proceedings of ACL’02*, (2002).
- [4] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom, *Database Systems: The Complete Book*, Prentice Hall Press, Upper Saddle River, NJ, USA, 2 edn., 2008.
- [5] Alessandra Giordani and Alessandro Moschitti, ‘Corpora for automatically learning to map natural language questions into sql queries’, in *Proceedings of LREC’10*, Valletta, Malta, (may 2010). European Language Resources Association (ELRA).

- [6] Johan Granberg and Michael Minock, 'A natural language interface over the musicbrainz database', in *Proceedings of the 1st Workshop on Question Answering over Linked Data (QALD-1) : Co-located with the 8th Extended Semantic Web Conference*, pp. 38–43, (2011). QC 20120413.
- [7] T. Joachims, 'Making large-scale SVM learning practical', in *Advances in Kernel Methods*, eds., B. Schölkopf, C. Burges, and A. Smola, (1999).
- [8] Rohit J. Kate and Raymond J. Mooney, 'Using string-kernels for learning semantic parsers', in *Proceedings of the 21st ICCL and 44th Annual Meeting of the ACL*, pp. 913–920, Sydney, Australia, (July 2006). Association for Computational Linguistics.
- [9] Bill MacCartney Marie-Catherine de Marneffe and Christopher D. Manning, 'Generating typed dependency parses from phrase structure parses', in *Proceedings LREC 2006*, (2006).
- [10] Michael Minock, Peter Olofsson, and Alexander Näslund, 'Towards building robust natural language interfaces to databases', in *NLDB '08: Proceedings of Natural Language and Information Systems*, Berlin, Heidelberg, (2008).
- [11] Ana-Maria Popescu, Oren A Etzioni, and Henry A Kautz, 'Towards a theory of natural language interfaces to databases', in *Proceedings of the 2003 International Conference on Intelligent User Interfaces*, Miami, (2003). Association for Computational Linguistics.
- [12] S Ruwanpura, 'Sq-hal: Natural language to sql translator'.
- [13] Gerard Salton, 'Recent trends in automatic information retrieval', in *SIGIR '86, Proceedings of the 9th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Pisa, Italy, September 8-10, 1986*, pp. 1–10. ACM, (1986).
- [14] Libin Shen and Aravind K. Joshi, 'An SVM-based voting algorithm with application to parse reranking', in *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pp. 9–16, (2003).
- [15] L. R. Tang and Raymond J. Mooney, 'Using multiple clause constructors in inductive logic programming for semantic parsing', in *Proceedings of the 12th European Conference on Machine Learning*, pp. 466–477, Freiburg, Germany, (2001).
- [16] Luke S. Zettlemoyer and Michael Collins, 'Learning to map sentences to logical form: Structured classification with probabilistic categorial grammars', in *UAI*, pp. 658–666, (2005).