

Natural Language Explanations for Query Results

Daniel Deutch
Tel Aviv University
danielde@post.tau.ac.il

Nave Frost
Tel Aviv University
navefrost@mail.tau.ac.il

Amir Gilad
Tel Aviv University
amirgilad@mail.tau.ac.il

ABSTRACT

Multiple lines of research have developed Natural Language (NL) interfaces for formulating database queries. We build upon this work, but focus on presenting a highly detailed form of the *answers* in NL. The answers that we present are importantly based on the *provenance* of tuples in the query result, detailing not only the results but also their *explanations*. We develop a novel method for transforming provenance information to NL, by leveraging the original NL query structure. Furthermore, since provenance information is typically large and complex, we present two solutions for its effective presentation as NL text: one that is based on provenance factorization, with novel desiderata relevant to the NL case, and one that is based on summarization.

1. INTRODUCTION

Developing Natural Language (NL) interfaces to database systems has been the focus of multiple lines of research (see e.g. [17, 2, 21]). In this work we complement these efforts by providing *NL explanations to query answers*. The explanations that we provide elaborate upon answers with additional important information, and are helpful for understanding *why* does each answer qualify to the query criteria.

As an example, consider the Microsoft Academic Search database (<http://academic.research.microsoft.com>) and consider the NL query in Figure 1a. A state-of-the-art NL query engine, NaLIR [17], is able to transform this NL query into the SQL query also shown (as a Conjunctive Query, which is the fragment that we focus on in this paper) in Figure 1b. When evaluated using a standard database engine, the query returns the expected list of organizations. However, the answers (organizations) in the query result lack *justification*, which in this case would include the authors affiliated with each organization and details of the papers they have published (their titles, their publication venues and publication years). Such additional information, corresponding to the notion of *provenance* (e.g. [12, 14, 6]) can lead to a richer answer than simply providing the names of organizations: it allows users to also see relevant details of the qualifying organizations. Provenance information is also valuable for validation of answers: a user who sees an organization name as an answer is likely to have a harder time

```
return the organization of authors who published papers
in database conferences after 2005
```

(a) NL Query

```
query(oname) :- org(oid, oname), conf(cid, cname),
pub(wid, cid, ptitle, pyear), author(aid, aname, oid),
domainConf(cid, did), domain(did, dname),
writes(aid, wid), dname = 'Databases', pyear > 2005
```

(b) CQ Q

Figure 1: NL Query and CQ Q

```
TAU is the organization of Tova M. who published
'OASSIS...' in SIGMOD in 2014
```

Figure 2: Answer For a Single Assignment

validating that this organization qualifies as an answer, than if she was presented with the full details of publications.

We propose a novel approach of presenting *provenance information for answers of NL queries, again as sentences in Natural Language*. Continuing our running example, Figure 2 shows one of the answers outputted by our system in response to the NL query in Figure 1a.

Our solution consists of the following key contributions.

Provenance Tracking Based on the NL Query Structure.

A first key idea in our solution is to leverage the *NL query structure* in constructing NL provenance. In particular, we modify NaLIR so that we store exactly which parts of the NL query translate to which parts of the formal query. Then, we evaluate the formal query using a provenance-aware engine (we use SelP [7]), further modified so that it stores which parts of the query “contribute” to which parts of the provenance. By composing these two “mappings” (text-to-query-parts and query-parts-to-provenance) we infer which parts of the NL query text are related to which provenance parts. Finally, we use the latter information in an “inverse” manner, to translate the provenance to NL text.

Factorization. A second key idea is related to the provenance size. In typical scenarios, a single answer may have multiple explanations (multiple authors, papers, venues and years in our example). A naïve solution is to formulate and present a separate sentence corresponding to each explana-

© VLDB Endowment 2017. This is a minor revision of the paper entitled “Provenance for Natural Language Queries”, published in the Proceedings of the VLDB Endowment, Vol. 10, No. 5, 2150-8097/17/01. DOI: <https://doi.org/10.14778/3055540.3055550>

We are extremely grateful to Fei Li and H.V. Jagadish for generously sharing with us the source code of NaLIR, and providing invaluable support.

tion. The result will however be, in many cases, very long and repetitive. As observed already in previous work [4, 18], different assignments (explanations) may have significant parts in common, and this can be leveraged in a *factorization* that groups together multiple occurrences. In our example, we can e.g. factorize explanations based on author, paper name, conference name or year. Importantly, we impose a novel constraint on the factorizations that we look for (which we call *compatibility*), intuitively capturing that their structure is consistent with a partial order defined by the parse tree of the question. This constraint is needed so that we can translate the factorization back to an NL answer whose structure is similar to that of the question. Even with this constraint, there may still be exponentially many (in the size of the provenance expression) compatible factorizations, and we look for the factorization with minimal size out of the compatible ones; for comparison, previous work looks for the minimal factorization with no such “compatibility constraint”. The corresponding decision problem remains coNP-hard (again in the provenance size), but we devise an effective and simple greedy solution. We further translate factorized representations to concise NL sentences, again leveraging the structure of the NL query.

Summarization. We propose *summarized* explanations by replacing details of different parts of the explanation by their synopsis, e.g. presenting only the number of papers published by each author, the number of authors, or the overall number of papers published by authors of each organization. Such summarizations incur by nature a loss of information but are typically much more concise and easier for users to follow. Here again, while provenance summarization has been studied before (e.g. [1, 18]), the desiderata of a summarization needed for NL sentence generation are different, rendering previous solutions inapplicable here. We observe a tight correspondence between factorization and summarization: every factorization gives rise to multiple possible summarizations, each obtained by counting the number of sub-explanations that are “factorized together”. We provide a robust solution, allowing to compute NL summarizations of the provenance, of varying levels of granularity.

2. PRELIMINARIES

We provide here the necessary preliminaries on Natural Language Processing, conjunctive queries and provenance.

2.1 From NL to Formal Queries

We start by recalling some basic notions from NLP, as they pertain to the translation process of NL queries to a formal query language. A key notion that we will use is that of the *syntactic dependency tree* of NL queries. This is essentially a node-labeled tree where labels consist of two components, as follows: (1) Part of Speech (*POS*): the syntactic role of the word; (2) Relationship (*REL*): the grammatical relationship between the word and its parent in the dependency tree.

We focus on a sub-class of queries handled by NaLIR, namely that of Conjunctive Queries, possibly with comparison operators ($=, >, <$) (NaLIR further supports nested queries and aggregation). The corresponding NL queries in NaLIR follow one of the two (very general) abstract forms described in Figure 3: an object (noun) is sought for, that satisfies some properties, possibly described through a complex sub-

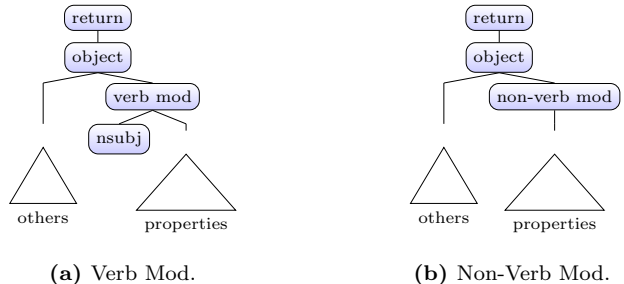


Figure 3: Abstract Dependency Trees

sentence rooted by a *modifier* (which may or may not be a verb, a distinction whose importance will be made clear later).

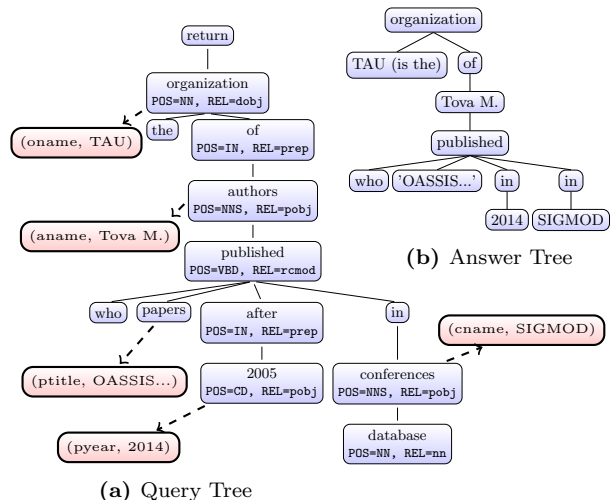


Figure 4: Question and Answer Trees

EXAMPLE 2.1. Reconsider the NL query in Figure 1a; its dependency tree is depicted in Figure 4a (ignore for now the arrows). The part-of-speech (*POS*) tag of each node reflects its syntactic role in the sentence – for instance, “organization” is a noun (denoted “NN”), and “published” is a verb in past tense (denoted “VBD”). The relation (*REL*) tag of each node reflects the semantic relation of its sub-tree with its parent. For instance, the *REL* of “of” is *prep* (“prepositional modifier”) meaning that the sub-tree rooted at “of” describes a property of “organization” while forming a complex sub-sentence. The tree in Figure 4a matches the abstract tree in Figure 3b since “organization” is the object and “of” is a non-verb modifier (its *POS* tag is *IN*, meaning “preposition or subordinating conjunction”) rooting a sub-sentence describing “organization”.

The dependency tree is transformed by NaLIR, based also on schema knowledge, to SQL. We focus in this work on NL queries that are compiled into Conjunctive Queries (CQs).

EXAMPLE 2.2. Reconsider our running example NL query in Figure 1a; a counterpart Conjunctive Query is shown in Figure 1b. Some words of the NL query have been mapped by NaLIR to variables in the query, e.g., the word “organization” corresponds to the head variable (*oname*). Additionally, some parts of the sentence have been compiled

$(\text{name,TAU}) \cdot (\text{name,Tova M.}) \cdot (\text{ptitle,OASSIS...}) \cdot$
 $(\text{cname,SIGMOD}) \cdot (\text{pyear,14'}) +$
 $(\text{name,TAU}) \cdot (\text{name,Tova M.}) \cdot (\text{ptitle,Querying...}) \cdot$
 $(\text{cname,VLDB}) \cdot (\text{pyear,06'}) +$
 $(\text{name,TAU}) \cdot (\text{name,Tova M.}) \cdot (\text{ptitle,Monitoring...}) \cdot$
 $(\text{cname,VLDB}) \cdot (\text{pyear,07'}) +$
 $(\text{name,TAU}) \cdot (\text{name,Slava N.}) \cdot (\text{ptitle,OASSIS...}) \cdot$
 $(\text{cname,SIGMOD}) \cdot (\text{pyear,14'}) +$
 $(\text{name,TAU}) \cdot (\text{name,Tova M.}) \cdot (\text{ptitle,A sample...}) \cdot$
 $(\text{cname,SIGMOD}) \cdot (\text{pyear,14'}) +$
 $(\text{name,UPENN}) \cdot (\text{name,Susan D.}) \cdot (\text{ptitle,OASSIS...}) \cdot$
 $(\text{cname,SIGMOD}) \cdot (\text{pyear,14'})$

Figure 5: Value-level Provenance

to boolean conditions based on the MAS schema, e.g., the part “in database conferences” was translated to $\text{dname} = \text{'Databases'}$ in the CQ depicted in Figure 1b. Figure 4a shows the mapping of some of the nodes in the NL query dependency tree to variables of Q (ignore for now the values next to these variables).

The translation performed by NaLIR from an NL query to a formal one can be captured by a *mapping* from (some) parts of the sentence to parts of the formal query. It can also be defined as a partial function from the nodes of the dependency tree to the variables of the query. We denote it by dependency-to-query-mapping.

2.2 Provenance

After compiling a formal query corresponding to the user’s NL query, we evaluate it and keep track of *provenance*, to be used in explanations. To define provenance, we first exemplify the standard notion of *assignments* for CQs.

Assignments allow for defining the semantics of CQs: a tuple t is said to appear in the query output if there exists an assignment α s.t. $t = \alpha(\text{head}(Q))$. They will also be useful in defining provenance below.

EXAMPLE 2.3. Consider again the query Q in Figure 1b and the database in Figure 6. The tuple (TAU) is an output of Q when assigning the highlighted tuples to the atoms of Q . As part of this assignment, the tuple (2, TAU) (the second tuple in the org table) and (4, Tova M., 2) (the second tuple of the author table) are assigned to the first and second atom of Q , respectively. In addition to this assignment, there are 4 more assignments that produce the tuple (TAU) and one assignment that produces the tuple (UPENN).

We next leverage assignments in defining provenance, introducing a simple value-level model. The idea is that assignments capture the *reasons* for a tuple to appear in the query result, with each assignment serving as an *alternative* such reason (indeed, the existence of a single assignment yielding the tuple suffices, according to the semantics, for its inclusion in the query result). Within each assignment, we keep record of the value assigned to each variable, and note that the *conjunction* of these value assignments is required for the assignment to hold. Capturing alternatives through the symbolic “+” and conjunction through the symbolic “.”, we arrive at the following definition of provenance as sum of products.

DEFINITION 2.4. Let $A(Q, D)$ be the set of assignments for a CQ Q and a database instance D . We define the value-

level provenance of Q w.r.t. D as

$$\sum_{\alpha \in A(Q, D)} \prod_{\{x_i, a_i \mid \alpha(x_i) = a_i\}} (x_i, a_i)$$

oid	oname
1	UPENN
2	TAU

aid	aname	oid
3	Susan D.	1
4	Tova M.	2
5	Slava N.	2

wid	cid	ptitle	pyear
6	10	“OASSIS...”	2014
7	10	“A sample...”	2014
8	11	“Monitoring...”	2007
9	11	“Querying...”	2006

aid	wid
4	6
3	6
5	6
4	7
4	8
4	9

cid	cname
10	SIGMOD
11	VLDB

cid	did
10	18
11	18

did	name
18	Databases

Figure 6: DB Instance

EXAMPLE 2.5. Re-consider our running example query and consider the database in Figure 6. The value-level provenance is shown in Figure 5. Each of the 6 summands stands for a different assignment (i.e. an alternative reason for the tuple to appear in the result). Assignments are represented as multiplication of pairs of the form (var, val) so that var is assigned val in the particular assignment. We only show here variables to which a query word was mapped; these will be the relevant variables for formulating the answer.

By composing the dependency-to-query-mapping from the NL query’s dependency tree to query variables, and the assignments of query variables to values from the database, we associate different parts of the NL query with values. We will use this composition of mappings throughout the paper as a means of assembling the NL answer to the NL query.

EXAMPLE 2.6. Continuing our running example, consider the assignment represented by the first monomial of Figure 5. Further reconsider Figure 4a, and now note that each node is associated with a pair (var, val) of the variable to which the node was mapped, and the value that this variable was assigned in this particular assignment. For instance, the node “organization” was mapped to the variable oname which was assigned the value “TAU”.

3. FIRST STEP: A SINGLE ASSIGNMENT

We start describing our transformation of provenance to NL for a single assignment. The solution will serve as the basis for the general case of multiple assignments.

3.1 Basic Solution

We follow the structure of the NL query dependency tree and generate an answer tree with the same structure by replacing/modifying the words in the question with the values from the result and provenance that were mapped using the dependency-to-query-mapping and the assignment. Yet, note that simply replacing the values does not always result in a coherent sentence, as shown in the following example.

EXAMPLE 3.1. Re-consider the dependency tree depicted in Figure 4a. If we were to replace the value in the organization node to the value “TAU” mapped to it, the word “organization” will not appear in the answer although it is needed to produce the coherent answer depicted in Figure 2. Without this word, it is unclear how to deduce the information about the connection between “Tova M.” and “TAU”.

We next account for these difficulties and exemplify our approach that outputs the dependency tree of a detailed answer; We do so by augmenting the query dependency tree into an answer tree. we will further translate this tree to an NL sentence.

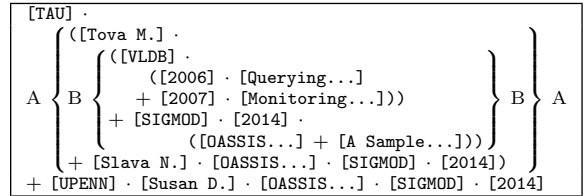
Recall that the dependency tree of the NL query follows one of the abstract forms in Figure 3. We distinguish between two cases based on nodes whose *REL* (relationship with parent node) is *modifier*; in the first case, the clause begins with a verb modifier (e.g., the node “published” in Fig. 4a is a verb modifier) and in the second, the clause begins with a non-verb modifier (e.g., the node “of” in Fig. 4a is a non-verb modifier). In short, the children of verb modifier nodes are replaced with the value mapped to them while the children of non-verb modifier nodes stay as part of the tree and the value mapped to them is added to the tree.

EXAMPLE 3.2. Re-consider Figure 4a, and note the mappings from the nodes to the variables and values as reflected in the boxes next to the nodes. To generate an answer, we follow the NL query structure, “plugging-in” mapped database values. We start with “organization”, which is the first node to be considered. Observe that “organization” has the child “of” which is a non-verb modifier, so we add “TAU” as its child. On the other hand, the node “authors” has the child “published” which is a verb modifier, so we replace “authors” with the value “Tova M.”, mapped to it. Another case is the handling of the nodes “after” and “in” which are modifiers as well. These nodes refer to times and locations, hence we replace the subtree rooted at these nodes with the node mapped to their child (in the case of “after” it is “2014” and in the case of “in” it is “SIGMOD”) and attach the node “in” as the parent of the node, in both cases as it is the suitable word for equality for years and locations.

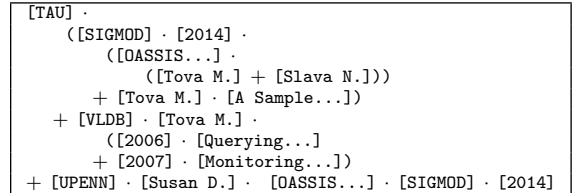
So far we have augmented the NL query dependency tree to obtain the dependency tree of the answer. The last step is to translate this tree to a sentence. To this end, we recall that the original query, in the form of a sentence, was translated by NaLIR to the NL query dependency tree. To translate the dependency tree to a sentence, we essentially “revert” this process, further using the mapping of NL query dependency tree nodes to (sets of) nodes of the answer.

4. THE GENERAL CASE

In general, as illustrated in Section 2, the provenance may include multiple assignments. We next generalize the construction to account for this. Note that a naïve solution in this respect is to generate a sentence for each individual assignment and concatenate the resulting sentences. However, already for the small-scale example presented here, this would result in a long and unreadable answer (recall Figure 5 consisting of six assignments). Instead, we propose two solutions: the first based on the idea of provenance factorization [18, 4], and the second leveraging factorization to provide a summarized form.



(a) f_1



(b) f_2

Figure 7: Provenance Factorizations

4.1 NL-Oriented Factorization

We start by defining the notion of factorization in a standard way (see e.g. [18, 8]).

DEFINITION 4.1. Let P be a provenance expression. We say that an expression f is a factorization of P if f may be obtained from P through (repeated) use of some of the following axioms: distributivity of summation over multiplication, associativity and commutativity of both summation and multiplication.

EXAMPLE 4.2. Re-consider the provenance expression in Figure 5. Two possible factorizations are shown in Figure 7, keeping only the values and omitting the variable names for brevity (ignore the A, B brackets for now). In both cases, the idea is to avoid repetitions in the provenance expression, by taking out a common factor that appears in multiple summands. Different choices of which common factor to take out lead to different factorizations.

How do we measure whether a possible factorization is suitable/preferable to others? Standard desiderata [18, 8] are that it should be short or that the maximal number of appearances of an atom is minimal. On the other hand, we factorize here as a step towards generating an NL answer; to this end, it will be highly useful if the (partial) order of nesting of value annotations in the factorization is consistent the (partial) order of corresponding words in the NL query. We will next formalize this intuition as a constraint over factorizations. We start by defining a partial order on nodes in a dependency tree:

DEFINITION 4.3. Given an dependency tree T , we define \leq_T as the descendant partial order of nodes in T : for each two nodes, $x, y \in V(T)$, we say that $x \leq_T y$ if x is a descendant of y in T .

EXAMPLE 4.4. In our running example (Figure 4a) it holds in particular that $authors \leq organization$, $2005 \leq authors$, $conferences \leq authors$ and $papers \leq authors$, but $papers$, 2005 and $conferences$ are incomparable.

Next we define a partial order over elements of a factorization, intuitively based on their nesting depth. To this end, we first consider the *circuit form* [3] of a given factorization:

EXAMPLE 4.5. Consider the partial circuit of f_1 in Figure 8. The root, \cdot , has two children; the left child is the leaf “TAU” and the right is a $+$ child whose subtree includes the part that is “deeper” than “TAU”.

Given a factorization f and an element n in it, we denote by $level_f(n)$ the distance of the node n from the root of the circuit induced by f multiplied by (-1) . Intuitively, $level_f(n)$ is bigger for a node n closer to the circuit root.

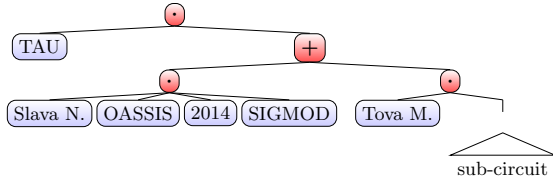


Figure 8: Sub-Circuit of f_1

Our goal here is to define the correspondence between the level of each node in the circuit and the level of its “source” node in the NL query’s dependency tree (note that each node in the query corresponds to possibly many nodes in the circuit: all values assigned to the variable in the different assignments). In the following definition we will omit the database instance for brevity and denote the provenance obtained for a query with dependency tree T by $prov_T$. Recall that dependency-to-query-mapping maps the nodes of the dependency tree to the query variables and the assignment maps these variables to values from the database.

DEFINITION 4.6. Let T be a query dependency tree, let $prov_T$ be a provenance expression, let f be a factorization of $prov_T$, let τ be a dependency-to-query-mapping and let $\{\alpha_1, \dots, \alpha_n\}$ be the set of assignments to the query. For each two nodes x, y in T we say that $x \leq_f y$ if $\forall i \in [n] : level_f(\alpha_i(\tau(x))) \leq level_f(\alpha_i(\tau(y)))$.

We say that f is T -compatible if each pair of nodes $x \neq y \in V(T)$ that satisfy $x \leq_T y$ also satisfy that $x \leq_f y$.

Essentially, T -compatibility means that the partial order of nesting between values, for each individual assignment, must be consistent the partial order defined by the structure of the question. Note that the compatibility requirement imposes constraints on the factorization, but it is in general far from dictating the factorization, since the order $x \leq_T y$ is only partial – and there is no constraint on the order of each two provenance nodes whose “origins” in the query are unordered. Among the T -compatible factorizations, we will prefer shorter ones.

DEFINITION 4.7. Let T be an NL query dependency tree and let $prov_T$ be a provenance expression for the answer. We say that a factorization f of $prov_T$ is optimal if f is T -compatible and there is no T -compatible factorization f' of $prov_T$ such that $|f'| < |f|$ ($|f|$ is the length of f).

The following example shows that the T -compatibility constraint still allows much freedom in constructing the factorization. In particular, different choices can (and sometimes should, to achieve minimal size) be made for different sub-expressions, including ones leading to different answers and ones leading to the same answer through different assignments.

EXAMPLE 4.8. Recall the partial order \leq_T imposed by our running example query, shown in part in Example 4.4. It implies that in every compatible factorization, the organization name must reside at the highest level, and indeed TAU was “pulled out” first in Figure 8; similarly the author name must be pulled out next. In contrast, since the query nodes corresponding to title, year and conference name are unordered, we may, within a single factorization, factor out e.g. the year in one part of the factorization and the conference name in another one. As an example, Tova M. has two papers published in VLDB (“Querying...” and “Monitoring”) so factorizing based on VLDB would be the best choice for that part. On the other hand, suppose that Slava N. had two paper published in 2014; then we could factorize them based on 2014. The factorization could, in that case, look like the following (where the parts taken out for Tova and Slava are shown in bold):

```
[TAU] ·
([Tova M.] ·
([VLDB] ·
  ([2006] · [Querying...]
  + [2007] · [Monitoring...]))
+ [SIGMOD] · [2014] ·
  ([OASSIS...] + [A Sample...]))
+ ([Slava N.] ·
  ([2014] ·
  ([SIGMOD] · [OASSIS...]
  + [VLDB] · [Ontology...])))
```

The following example shows that in some cases, requiring compatibility can come at the cost of compactness.

EXAMPLE 4.9. Consider the query tree T depicted in Figure 4a and the factorizations $prov_T$ (the identity factorization) depicted in Figure 5, f_1, f_2 presented in Figure 7. $prov_T$ is of length 30 and is 5-readable, i.e., the maximal number of appearances of a single variable is 5 (see [8]). f_1 is of length 20, while the length of f_2 is only 19. In addition, both f_1 and f_2 are 3-readable. Based on those measurements f_2 seems to be the best factorization, yet f_1 is T -compatible with the question and f_2 is not. For example, conferences \leq_T authors but “SIGMOD” appears higher than “Tova M.” in f_2 . Choosing a T -compatible factorization in f_1 will lead (as shown below) to an answer whose structure resembles that of the question, and thus translates to a more coherent and fitting NL answer.

Note that the identity factorization is always T -compatible, so we are guaranteed at least one optimal factorization (but it is not necessarily unique). We next study the problem of computing such a factorization.

4.2 Computing Factorizations

Recall that our notion of compatibility restricts the factorizations so that their structure resembles that of the question. Without this constraint, finding shortest factorizations is coNP-hard in the size of the provenance (i.e. a boolean expression) [13]. The compatibility constraint does not reduce the complexity since it only restricts choices relevant to part of the expression, while allowing freedom for arbitrarily many other elements of the provenance. Also recall (Example 4.8) that the choice of which element to “pull-out” needs in general to be done separately for each part of the provenance so as to optimize its size (which is the reason for the hardness in [13] as well). In general:

PROPOSITION 4.10. *Given a dependency tree T , a provenance expression $prov_T$ and an integer k , deciding whether there exists a T -compatible factorization of $prov_T$ of size $\leq k$ is coNP-hard.*

Greedy Algorithm. Despite the above result, the constraint of compatibility does help in practice, in that we can avoid examining choices that violate it. For other choices, we devise a simple algorithm that chooses greedily among them. More concretely, the input to Algorithm 1 is the query tree T_Q (with its partial order \leq_{T_Q}), and the provenance $prov_{T_Q}$. The algorithm output is a T_Q -compatible factorization f . Starting from $prov$, the progress of the algorithm is made in steps, where at each step, the algorithm traverses the circuit induced by $prov$ in a BFS manner from top to bottom and takes out a variable that would lead to a minimal expression out of the valid options that keep the current factorization T -compatible. Naturally, the algorithm does not guarantee an optimal factorization (in terms of length), but performs well in practice.

In more detail, we start by choosing the largest nodes according to \leq_{T_Q} which have not been processed yet (Line 2). Afterwards, we sort the corresponding variables in a greedy manner based on the number of appearances of each variable in the expression using the procedure *sortByFrequentVars* (Line 3). In Lines 4-5, we iterate over the sorted variables and extract them from their sub-expressions. This is done while preserving the \leq_{T_Q} order with the larger nodes, thus ensuring that the factorization will remain T_Q -compatible. We then add all the newly processed nodes to the set *Processed* which contains all nodes that have already been processed (Line 6). Lastly, we check whether there are no more nodes to be processed, i.e., if the set *Processed* includes all the nodes of T_Q (denoted $V(T_Q)$, see the condition in Line 7). If the answer is “yes”, we return the factorization. Otherwise, we make a recursive call. In each such call, the set *Processed* becomes larger until the condition in Line 7 holds.

Algorithm 1: GreedyFactorization

input : T_Q - the query tree, \leq_{T_Q} - the query partial order, $prov$ - the provenance, τ, α - dependency-to-query-mapping and assignment from nodes in T_Q to provenance variables, *Processed* - subset of nodes from $V(T_Q)$ which were already processed (initially, \emptyset)

output: f - T_Q -compatible factorization of $prov_{T_Q}$

- 1 $f \leftarrow prov$;
- 2 $Frontier \leftarrow \{x \in V(T_Q) \mid \forall (y \in V(T_Q) \setminus Processed) \text{ s.t. } x \not\leq_{T_Q} y\}$;
- 3 $vars \leftarrow sortByFrequentVars(\{\alpha(\tau(x)) \mid x \in Frontier\}, f)$;
- 4 **foreach** $var \in vars$ **do**
- 5 Take out var from sub-expressions in f not including variables from $\{x \mid \exists y \in Processed : x = \alpha(\tau(y))\}$;
- 6 $Processed \leftarrow Processed \cup Frontier$;
- 7 **if** $|Processed| = |V(T_Q)|$ **then**
- 8 **return** f ;
- 9 **else**
- 10 **return** $GreedyFactorization(T_Q, f, \tau, \alpha, Processed)$;

EXAMPLE 4.11. *Consider the query tree T_Q depicted in Figure 4a, and provenance $prov$ in Figure 5. As explained*

above, the largest node according to \leq_{T_Q} is organization, hence “TAU” will be taken out from the brackets multiplying all summands that contain it. Afterwards, the next node according to the order relation will be author, therefore we group by author, taking out “Tova M.”, “Slava N.” etc. The following choice (between conference, year and paper name) is then done greedily for each author, based on its number of occurrences. For instance, VLDB appears twice for Tova.M. whereas each paper title and year appears only once; so it will be pulled out. The polynomial [SlavaN.].[OASSIS...].[SIGMOD].[2014] will remain unfactorized as all values appear once. Eventually, the algorithm will return the factorization f_1 depicted in Figure 7, which is T_Q -compatible and much shorter than the initial provenance expression.

PROPOSITION 4.12. *Let f be the output of Algorithm 1 for the input dependency tree T_Q , then f is T_Q -compatible.*

Complexity. Denote the provenance size by n . The algorithm complexity is $O(n^2 \cdot \log n)$: at each recursive call, we sort all nodes in $O(n \cdot \log n)$ (Line 3) and then we handle (in *Frontier*) at least one node (in the case of a chain graph) or more. Hence, in the worst case we would have n recursive calls, each one costing $O(n \cdot \log n)$.

4.3 Factorization to Answer Tree

The final step is to turn the obtained factorization into an NL answer. Similarly to the case of a single assignment (Section 3), we leverage the mappings and assignments to convert the query dependency tree into an answer tree that reflects the factorization. Intuitively, we follow the structure of a single answer, replacing each node there by either a single node, standing for a single word of the factorized expression, or by subtree, standing for some brackets (sub-circuit) in the factorized expression.

EXAMPLE 4.13. *Consider the factorization f_1 depicted in Figure 7, and the structure of single assignment answer depicted in Figure 4b which was built based on an answer tree for a single assignment. Given this input, we will generate an answer tree corresponding to the following sentence:*

TAU is the organization of
Tova M. who published
in VLDB
‘Querying...’ in 2006 and
‘Monitoring...’ in 2007
and in SIGMOD in 2014
‘OASSIS...’ and ‘A sample...’
and Slava N. who published
‘OASSIS...’ in SIGMOD in 2014.
UPENN is the organization of Susan D. who published
‘OASSIS...’ in SIGMOD in 2014.

Note that the query has two results: “TAU” and “UPENN”. “UPENN” was produced with a single assignment, but there are five different assignments producing “TAU”. Focusing on the factorization part of the result “TAU”, notice that the authors were pulled out first, then the conferences, and then the years and papers, so this will be reflected in the factorized answer tree. For example, we replace the node authors with the values from the factorization that correspond to this word, i.e., Tova M. and Slava N. The answer tree can also be changed based on the hierarchy of the factorization. For instance, although the node paper is closer to the root of the tree then the nodes year and conference in the original answer tree, the order of these nodes in the new answer

tree will be reversed since f_1 extracted the values “VLDB”, “SIGMOD” and “2014”.

Why require compatibility? We conclude this part of the paper by revisiting our decision to require compatible factorizations, highlighting difficulties in generating NL answers using non-compatible factorizations.

EXAMPLE 4.14. Consider factorization f_2 from Figure 7. “TAU” should be at the beginning of the sentence and followed by the conference names “SIGMOD” and “VLDB”. The second and third layers of f_2 are composed of author names (“Tova M.”, “Slava N.”), paper titles (“OASSIS”, “A sample...”, “Monitoring...”) and publication years (2007, 2014). Changing the original order of the words such that the conference name “SIGMOD” and the publication year “2014” will appear before “Tova M.” breaks the sentence structure in a sense. It is unclear how to algorithmically translate this factorization into an NL answer, since we need to patch the broken structure by adding connecting phrases. One hypothetical option of patching f_2 and transforming it into an NL answer is depicted below. The bold parts of the sentence are not part of the factorization and it is not clear how to generate and incorporate them into the sentence algorithmically. Even if we could do so, it appears that the resulting sentence would be quite convoluted:

TAU is the organization of authors who published in SIGMOD 2014
‘OASSIS...’ which was published by
Tova M. and Slava N.
and Tova M. published ‘A sample...’
and Tova M. published in VLDB
‘Querying...’ in 2014
and ‘Monitoring...’ in 2007.
UPENN is the organization of Susan D. who published
‘OASSIS...’ in SIGMOD in 2014

Observe that the resulting sentence is much less clear than the one obtained through our approach, even though it was obtained from a shorter factorization f_2 ; the intuitive reason is that since f_2 is not T-compatible, it does not admit a structure that is similar to that of the question, thus is not guaranteed to admit a structure that is coherent in Natural Language. Interestingly, the sentence we would obtain in such a way also has an edit distance from the question [9] that is shorter than that of our answer, demonstrating that edit distance is not an adequate measure here.

4.4 From Factorizations to Summarizations

When there are many assignments and/or the assignments involve multiple distinct values, even an optimal factorized representation may be too long and convoluted for users to follow.

EXAMPLE 4.15. Reconsider Example 4.13; if there are many authors from TAU then even the compact representation of the result could be very long. In such cases we need to summarize the provenance in some way that will preserve the “essence” of all assignments without actually specifying them, for instance by providing only the number of authors/-papers for each institution.

To this end, we employ *summarization*, as follows. First, we note that a key to summarization is understanding which parts of the provenance may be grouped together. For that, we use again the mapping from nodes to query variables:

```
(A) [TAU] · Size([Tova M.], [Slava N.]) · Size([VLDB], [SIGMOD]) ·
    Size([Querying...], [Monitoring...],
        [OASSIS...], [A Sample...]) · Range([2006], [2007], [2014])
(B) [TAU] · (
    [Tova M.] ·
    Size([VLDB], [SIGMOD]) ·
    Size([Querying...], [Monitoring...],
        [OASSIS...], [A Sample...]) · Range([2006], [2007], [2014])
    [Slava N.] · [OASSIS...] · [SIGMOD] · [2014])
```

Figure 9: Summarized Factorizations

```
(A) TAU is the organization of 2 authors who published
4 papers in 2 conferences in 2006 - 2014.
(B) TAU is the organization of Tova M. who published
4 papers in 2 conferences in 2006 - 2014 and Slava N.
who published ‘OASSIS...’ in SIGMOD in 2014.
```

Figure 10: Summarized Sentences

we say that two nodes are of the same *type* if both were mapped to the same query variable. Now, let n be a node in the circuit form of a given factorization f . A summarization of the sub-circuit of n is obtained in two steps. First, we group the descendants of n according to their type. Then, we summarize each group separately. The latter is done in our implementation simply by either counting the number of distinct values in the group or by computing their range if the values are numeric. In general, one can easily adapt the solution to apply additional user-defined “summarization functions” such as “greater / smaller than X” (for numerical values) or “in continent Y” for countries.

EXAMPLE 4.16. Re-consider the factorization f_1 from Figure 7. We can summarize it in multiple levels: the highest level of authors (summarization “A”), or the level of papers for each particular author (summarization “B”), or the level of conferences, etc. Note that if we choose to summarize at some level, we must summarize its entire sub-circuit (e.g. if we summarize for Tova. M. at the level of conferences, we cannot specify the papers titles and publication years).

Figure 9 presents the summarizations of sub-trees for the “TAU” answer, where “size” is a summarization operator that counts the number of distinct values and “range” is an operator over numeric values, summarizing them as their range. The summarized factorizations are further converted to NL sentences which are shown in Figure 10. Summarizing at a higher level results in a shorter but less detailed summarization.

5. RELATED WORK

Multiple lines of work (e.g. [2, 17, 21]) have proposed NL interfaces to formulate database queries, and additional works [10] have focused on presenting the *answers* in NL, typically basing their translation on the schema of the output relation. Among these, works such as [2, 17] also harness the dependency tree in order to make the translation from NL to SQL by employing *mappings* from the NL query to formal terms. To our knowledge, no previous work has focused on formulating the *provenance* of output tuples in NL. This requires fundamentally different techniques (e.g. that of factorization and summarization, building the sentence based on the input question structure, etc.) and leads to answers of much greater detail.

The tracking, storage and presentation of provenance have been the subject of extensive research (see e.g. [12, 14, 6]) while the field of provenance applications has also been

Table 1: Sample use-cases and results

Query	Single Assignment	Multiple Assignments - Summarized
Return the homepage of SIGMOD	http://www.sigmod2011.org/ is the homepage of SIGMOD	
Return the authors who published papers in SIGMOD before 2015 and after 2005	Tova M. published "Auto-completion..." in SIGMOD in 2012	Tova M. published 10 papers in SIGMOD in 2006-2014
Return the authors from TAU who published papers in VLDB	Tova M. from TAU published "XML Repository..." in VLDB	Tova M. from TAU published 11 papers in VLDB
Return the authors who published papers in database conferences	Tova M. "published Auto-completion..." in SIGMOD	Tova M. published 96 papers in 18 conferences
Return the organization of authors who published papers in database conferences after 2005	TAU is the organization of Tova M. who published 'OASSIS...' in SIGMOD in 2014	TAU is the organization of 43 authors who published 170 papers in 31 conferences in 2006 - 2015

broadly studied (e.g. [7, 19]). A longstanding challenge in this context is the complexity of provenance expressions, leading to difficulties in presenting them in a user comprehensible manner. Approaches in this respect include showing the provenance in a graph form (see e.g. [20, 6]), allowing user control over the level of granularity ("zooming" in and out [5]), or otherwise presenting different ways of provenance visualization [14]. Other works have studied allowing users to query the provenance (e.g. [16, 15]) or to a-priori request that only parts of the provenance are tracked (see for example [7, 11]). Importantly, provenance factorization and summarization have been studied (e.g., [1, 18, 4]) as means for compact representation of the provenance. Usually, the solutions proposed in these works aim at reducing the size of the provenance but naturally do not account for its presentation in NL; we have highlighted the different considerations in context of factorization/summarization in our setting.

6. CONCLUSION AND LIMITATIONS

We have studied in this paper, for the first time to our knowledge, provenance for NL queries. We have devised means for presenting the provenance information again in Natural Language, in factorized or summarized form.

There are two main limitations to our work. First, a part of our solution was designed to fit NaLIR, and will need to be replaced if a different NL query engine is used. Specifically, the "sentence generation" module will need to be adapted to the way the query engine transforms NL queries into formal ones; our notions of factorization and summarization are expected to be easier to adapt to a different engine. Second, our solution is limited to Conjunctive Queries. One of the important challenges in supporting NL provenance for further constructs such as union and aggregates is the need to construct a concise presentation of the provenance in NL.

Acknowledgments. This research was partially supported by the Israeli Science Foundation (ISF, grant No. 1636/13), and by ICRC - The Blavatnik Interdisciplinary Cyber Research Center. The contribution of Amir Gilad is part of a Ph.D. thesis research conducted at Tel Aviv University.

7. REFERENCES

- [1] E. Ainy, P. Bourhis, S. B. Davidson, D. Deutch, and T. Milo. Approximated summarization of data provenance. In *CIKM*, pages 483–492, 2015.
- [2] Y. Amsterdamer, A. Kukliansky, and T. Milo. A natural language interface for querying general and individual knowledge. *VLDB*, pages 1430–1441, 2015.
- [3] P. Brgisser, M. Clausen, and M. A. Shokrollahi. *Algebraic Complexity Theory*. Springer Publishing Company, Incorporated, 2010.
- [4] A. P. Chapman, H. V. Jagadish, and P. Ramanan. Efficient provenance storage. In *SIGMOD*, pages 993–1006, 2008.
- [5] S. Cohen-Boulakia, O. Biton, S. Cohen, and S. Davidson. Addressing the provenance challenge using zoom. *Concurr. Comput. : Pract. Exper.*, pages 497–506, 2008.
- [6] S. B. Davidson and J. Freire. Provenance and scientific workflows: challenges and opportunities. In *SIGMOD*, pages 1345–1350, 2008.
- [7] D. Deutch, A. Gilad, and Y. Moskovitch. Selective provenance for datalog programs using top-k queries. *PVLDB*, pages 1394–1405, 2015.
- [8] K. Elbassioni, K. Makino, and I. Rauf. On the readability of monotone boolean formulae. *JoCO*, pages 293–304, 2011.
- [9] M. Emms. Variants of tree similarity in a question answering task. In *Proceedings of the Workshop on Linguistic Distances*, pages 100–108, 2006.
- [10] E. Franconi, C. Gardent, X. I. Juarez-Castro, and L. Perez-Beltrachini. Qelo Natural Language Interface: Generating queries and answer descriptions. In *Natural Language Interfaces for Web of Data*, 2014.
- [11] B. Glavic. Big data provenance: Challenges and implications for benchmarking. In *Specifying Big Data Benchmarks - First Workshop, WBDB*, pages 72–80, 2012.
- [12] T. Green, G. Karvounarakis, and V. Tannen. Provenance semirings. In *PODS*, pages 31–40, 2007.
- [13] E. Hemaspaandra and H. Schnoor. Minimization for generalized boolean formulas. In *IJCAI*, pages 566–571, 2011.
- [14] M. Herschel and M. Hlawatsch. Provenance: On and behind the screens. In *SIGMOD*, pages 2213–2217, 2016.
- [15] Z. G. Ives, A. Haeberlen, T. Feng, and W. Gatterbauer. Querying provenance for ranking and recommending. In *TaPP*, pages 9–9, 2012.
- [16] G. Karvounarakis, Z. G. Ives, and V. Tannen. Querying data provenance. In *SIGMOD*, pages 951–962, 2010.
- [17] F. Li and H. V. Jagadish. Constructing an interactive natural language interface for relational databases. *Proc. VLDB Endow.*, pages 73–84, 2014.
- [18] D. Olteanu and J. Závodný. Factorised representations of query results: Size bounds and readability. In *ICDT*, pages 285–298, 2012.
- [19] S. Roy and D. Suci. A formal approach to finding explanations for database queries. In *SIGMOD*, pages 1579–1590, 2014.
- [20] Y. L. Simmhan, B. Plale, and D. Gannon. Karma2: Provenance management for data-driven workflows. *Int. J. Web Service Res.*, pages 1–22, 2008.
- [21] D. Song, F. Schilder, C. Smiley, C. Brew, T. Zielund, H. Bretz, R. Martin, C. Dale, J. Duprey, T. Miller, and J. Harrison. TR discover: A natural language interface for querying and analyzing interlinked datasets. In *ISWC*, pages 21–37, 2015.