

Extending the TIGER Query Language with Universal Quantification

Torsten Marek

Saarland University

Department of Computational Linguistics

66041 Saarbrücken, Germany

`tmarek@coli.uni-saarland.de`

Joakim Lundborg

Stockholm University

Department of Linguistics

106 91 Stockholm, Sweden

`joakim.lundborg@ling.su.se`

Martin Volk

University of Zürich

Institute of Computational Linguistics

8050 Zürich, Switzerland

`volk@cl.uzh.ch`

The query language in TIGERSearch is limited due to its lack of universal quantification. This restriction makes it impossible to ask simple queries like “Find sentences that do not include a certain word”. We propose an easy way to formulate such queries. We have implemented this extension to the query language in a tool that allows querying parallel treebanks including their alignment constraints. Our implementation of universal quantification relies on the view of node sets rather than single node unification. Our query tool is freely available.

1 Introduction

In the last ten years, many languages and tools for querying syntactically annotated corpora have been developed. Most of these tools and query languages have been designed for a specific corpus and a specific annotation format. TGrep2¹ was originally designed for the Penn Treebank and its tree format. TIGERSearch (Lezius, 2002; König and Lezius, 2003) was designed for the NEGRA and TIGER corpora (Brants et al., 2002) and uses the TIGER-XML format.

TIGERSearch is a tool that allows the user to search a treebank by formulating queries in feature-value pairs. For example, one may search for the word 'can' with the part-of-speech 'noun' by querying

(1) `[word="can" & pos="NN"]`

In addition to constraints over words, the TIGER query language allows the user to also use constraints over dominance relations (search for a node n1 that dominates a node n2 in the tree), precedence relations and node predicates (like arity, discontinuous, and root).

In general, the design of the input format influences the design of the query language to a large extent, since it defines what can be queried. For instance, TIGER-XML supports crossing branches, leading to non-terminal nodes whose terminals are not a proper substring of the sentence (discontinuous nodes). The TIGER query language thus has special functions for dealing with discontinuous nodes. In contrast, the Penn Treebank formalism does not support crossing branches directly, and thus TGrep2 has no means for this notion.

While certain limitations of query languages are therefore due to the original design and could only be approximated, other interesting queries may simply be missing from the query language. Lai and Bird (2004) list seven exemplary queries, named Q1-Q7, which each formalism should support, regardless of the annotation formalism².

In this paper, we will deal with queries that require universal quantification, i.e. selecting a tree by stating constraints over possibly many nodes rather than individual nodes. The sample queries contain two examples where universal quantification is needed (Lai and Bird, 2004, p. 2):

Q2. Find sentences that do not include the word 'saw'.

Q5. Find the first common ancestor of sequences of a noun phrase followed by a verb phrase.

With the TIGER query language and its implementation TIGERSearch, these queries can only be approximated. The result set generated for the approximated queries will likely contain errors.

¹<http://tedlab.mit.edu/~dr/TGrep2/>

²with one exception: Q6 assumes multiple layers of annotation.

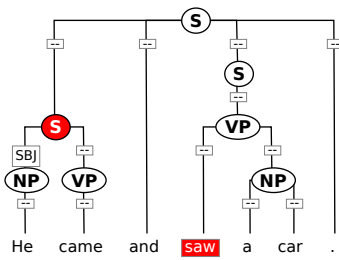


Figure 1: Wrong highlights.

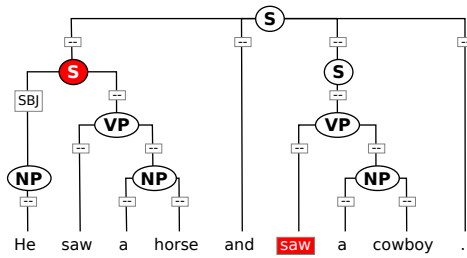


Figure 2: A false positive.

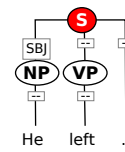


Figure 3: A false negative.

In section 2, we will analyze the limitations of TIGER and TIGERSearch and look at the solution proposed in Lezius (2002). In section 3, we will develop an extension to the TIGER query formalism that is able to deal with Q2 and similar queries, and we will analyze the different requirements of Q2 and Q5 in section 5.

In section 4, we will talk about the implementation of the extensions from the previous section in our own implementation of the TIGER query language, which is part of the Stockholm TreeAligner³ (Volk et al., 2007).

In the remainder of the paper, we will speak of syntax *graphs* rather than trees. These graphs are directed, acyclic and do not contain structure sharing (i.e. each node has exactly one direct ancestor). However, due to the existence of crossing branches, TIGER trees cannot be stored as nested lists or XML DOM trees directly, which is the usual understanding of trees.

Node descriptions are boolean expressions of feature constraints of the form "[feature=value]". They are used for finding nodes (assignments) in the corpus which are then used for the constraint resolution in TIGER queries.

2 Limitations of the TIGER Query Language

Figures 1 to 3 show three example graphs for Q2 (“Find sentences that do not include the word ‘saw’”). If the query were evaluated correctly, the result set would only contain the first coordinated sentence from figure 1 (“He came”) and the main sentence of figure 3 (“He left”).

In the TIGER query language however, every node variable is implicitly existentially quantified. For example, the query

$$(2) \quad \#s:[cat="S"] !>^* \#w:[word="saw"]$$

states that some node of category S does not dominate some word with the surface string ‘saw’. This query is interpreted by TIGERSearch so that it returns all combinations of two nodes $\#s$, $\#w$ in all graphs, such that $\#s$ does not dominate $\#w$. From the graphs that were actually meant to be matched by Q2, it will only return those that have a terminal ‘saw’ outside of any S node. Thus, the result set contains the graph from figure 1 (with an extra highlighting, which is distracting, but not fatal), but also, for instance, the graph from figure 2, which should definitely not be matched by the query. But all graphs that do not contain any ‘saw’, like the one in figure 3, will not show up in the result set.

Another attempt to formulate Q2 in TIGERSearch is the query

$$(3) \quad \#s:[cat="S"] >^* \#w:[word!="saw"]$$

which states that some node of category S dominates some terminal that does not have the surface string ‘saw’. Looking at the example graphs, it immediately becomes clear that there are many combinations of nodes that satisfy this query.

In general, there is no correct way to formulate the desired query in TIGERSearch. This limitation is acknowledged by the developers:

The use of the universal quantifier causes computational overhead since universal quantification usually means that a possibly large number of copies of logical expressions have to be produced. For the sake of computational simplicity and tractability, the universal quantifier is (currently) not part of the TIGER language. (TIGERSearch Help, section 10.3)

Section 5.7 of Lezius (2002) contains a proposal for an extension of the TIGER query formalism that combines universal quantification and the implication operator. Using this syntax, Q2 can now be formulated as

$$(4) \quad \forall \#w ([cat="S"] >^* \#w) \Rightarrow \#w:[word!="saw"]$$

³<http://dev.ling.su.se/treealigner>

which means that every node dominated by an S node must not be the surface string 'saw'.

Expressing the queries using an implication operator is natural given the unification-based evaluation of queries in TIGERSearch. As already mentioned, an actual implementation comes at great computational cost. For each \forall clause in the query, all nodes in the graph have to be iterated to find out if they satisfy $l \Rightarrow r$. In some cases, using the logical equivalent $!r \vee l$ can be used to speed up the queries. For instance, Q2 can simply be evaluated by

1. retrieving all S nodes in all graphs
2. retrieving all nodes where [word="saw"].
3. for each S node, check if it dominates none of the nodes from step 2.

Apart from runtime complexity considerations, the syntax is extended with a construct that is conceptually hard to grasp and that makes the grammar of the query language much more complex. We therefore decided to explore a different path.

3 Design

The Lezius solution presented in section 2 builds on the query calculus that is at the core of TIGERSearch's query evaluation engine. This calculus is based on unifying the partial graph description given by the query with any of the graph definitions in the corpus. If the unification succeeds, the graph matches the query and is returned in the result set.

In contrast, the query engine in the Stockholm TreeAligner is based on node sets, and combinations of nodes from the different sets to satisfy the constraints given in a query⁴.

In the previous analysis of Q2, we showed that it is possible to rephrase the query using logical equivalents. Therefore, the query "get all S nodes that do not contain the word saw" can be rephrased into "get all graphs where all instances of 'saw', if any, are not dominated by a specific S node".

This is essentially expressed in query 2. But as we already showed, the usage of the existential quantification will not lead to the expected results. However, if one of the two operands is not understood to be a single node from the graph, but a *set* of nodes, the result will be correct. Therefore, we introduce a new type into the query language, the node set which we indicate with the % symbol. Bare node descriptions are still bound by an implicit existential quantifier as before. A node set is only bound to a variable that starts with a percentage symbol:

$$(5) \quad \#s:[cat="S"] !>* \%w:[word="saw"]$$

If one operand in a constraint is a node set instead of a node, the semantics of the constraint are changed. Here, only those assignments to #s are returned where the

⁴cf. section 4 for a brief outline of the evaluation strategy.

constraint ' $!>^*$ ' holds for each node in the node set $\%w$, which contains all terminals with the surface string 'saw'. When applied to the graphs from the small sample corpus, this query now does not yield any false positives like the graph in figure 2.

The semantics of node predicates as defined in the TIGER query language do not change, they still operate at the node level. In the query

$$(6) \quad \%np:[cat="NP"] \ \& \ tokenarity(\%np, 2)$$

the node set $\%np$ will contain all NPs whose token arity is 2 (which means that each NP dominates exactly two tokens). See section 5 for a further discussion.

3.1 Node Sets

If each variable is bound by an existential quantifier, evaluation of a query can terminate as soon as one node description does not yield any results. Graphs that do not contain matching nodes for any of the descriptions will also be disregarded, which is why the graph from figure 3 will still not be matched by the query. To produce correct results, the semantics of node descriptions bound to node sets have to be changed. In contrast to existentially quantified nodes, which may not be undefined, a set can be the empty set \emptyset . If this is the case, a constraint is trivially true.

With this change in place, TIGER is in Cantor's paradise, and no one shall expel it from there. The basic semantics of set types are defined and new set predicates could be introduced to make set queries more powerful. As an example, consider the query "Return all NPs that do not contain any PP, but only if the graph contains PPs". Given that empty node sets are now allowed, the query has to be written as

$$(7) \quad [cat="NP"] \ !>^* \ \%pp:[cat="PP"] \ \& \ [cat="PP"].$$

The last term ensures that at least one PP exists in the graph. As a side effect, the result set will contain one entry for each combination of NP and PP in a matching graph, which is slightly more than what the query was supposed to yield.

To express constraints on set cardinality, a syntax for set algebra operations could be added to TIGER. As an example, to make sure that the set $\%pp$ contains elements, one could think of something like

$$(8) \quad [cat="NP"] \ !>^* \ \%pp:[cat="PP"] \ \& \ \{size(\%pp) > 0\}$$

Any expression enclosed in curly brackets is evaluated as an operation on sets. This addition would add a lot of power to the query language, but would make it much harder to use. Also, it requires a nontrivial amount of implementation effort and makes the query grammar more cumbersome. In our opinion, node set operations go beyond the scope of what can conveniently be handled inside of a single expression, too. If these expressions were added, it would be desirable to store node sets, reuse them in later queries or combine node sets from different queries.

Instead of adding full support for set operations, we introduce two new predicates that operate exclusively on node sets: *empty* and *nonempty*. The semantics of the predicates can be inferred from the names, and the previous query can be written in a straightforward manner:

(9) `[cat="NP"] !>* %pp:[cat="PP"] & nonempty(%pp)`

It is now also possible to search for graphs that do not contain a specific kind of node by using *empty*. The query

(10) `%w:[pos="det"] & empty(%w)`

returns all graphs that do not contain any determiner.

4 Implementation of the Query Language within the Stockholm TreeAligner

We have re-implemented the TIGER query language in a tool for the creation and exploration of parallel treebanks, the Stockholm TreeAligner. The TreeAligner allows the user to load two treebanks, typically with parallel (i.e. translated) sentences. For example, we have used the TreeAligner to work on a German treebank and its parallel English treebank. We have aligned the two treebanks first on the sentence level to get corresponding tree pairs and then on the word level and phrase level. Figure 4 shows a tree pair from our parallel treebank.

We currently distinguish between exact translation correspondence represented by green lines and approximate ('fuzzy') translation correspondence represented by red lines. Although it is useful in principle to make this distinction, it is often difficult in practice to consistently make this decision.

The TreeAligner is an editor which allows the user to create and modify alignments. It can also be used to browse and search parallel treebanks and this is where the query language comes in. We have re-implemented and extended the TIGER query language within the TreeAligner to search both treebanks in parallel. We also have added alignment constraints in order to combine the queries. Let us consider the following example query:

(11) `treebank1 #node1:[cat="NP"] > [cat="PP"]`
`treebank2 #node2:[cat="NP"] > [cat="PP"]`
`alignment #node1 * #node2`

Here, `#node1` is a variable that identifies a node of category NP in a graph in treebank1. And `#node2` identifies a node of category NP in treebank2. These variables correspond

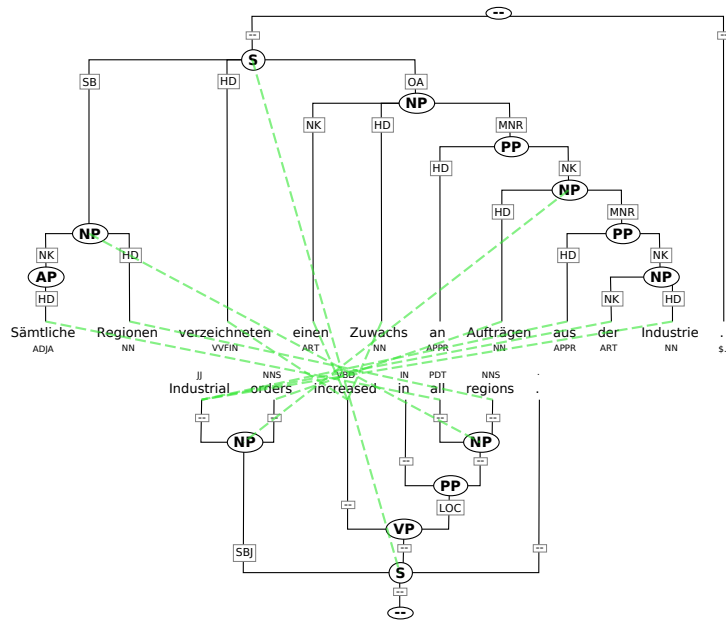


Figure 4: An example of two aligned sentences from the SMULTRON corpus

exactly to the syntax in the TIGERSearch query language. We then use these variables in the alignment query. The general alignment relation is indicated by the '*' operator. For a detailed description of the alignment query syntax, see Volk et al. (2007). This query searches through both treebank1 and treebank2 for noun phrases that dominate a prepositional phrase and returns all matches where these noun phrases are aligned with each other.

The implementation of the TIGER query language in the Stockholm TreeAligner is based on sets of nodes and constrained Cartesian products over these sets. Because of that, it was possible to implement the extensions described in section 3 with little programming effort. In contrast to existentially quantified nodes, node sets are subject to some restrictions. In a constraint, at most one operand may be a node set. Constraints that have two node sets as operands will lead to a runtime error⁵. In the result display, in contrast to existentially quantified nodes, nodes from sets are not highlighted. While there are no technical reasons for this, having the node sets stand out from the rest is not helpful, simply because too many nodes would be highlighted. This would lead to a confusing result display.

⁵This behavior may change in the future. However, to this date we do not have linguistically interesting queries that require a constraint with two node sets.

5 Beyond Node Descriptions

Using the implication operator, Q5 (finding the first common ancestor of a sequence NP VP) can be expressed in the following manner (where NT stands for any non-terminal):

$$\begin{aligned}
 (12) \quad & \#a:[NT] >^* \#np:[cat="NP"] \& \\
 & \#a >^* \#vp:[cat="VP"] \& \\
 & \#np .* \#vp \& \\
 & \forall \#b (\#a >^* \#b \& \#b >^* \#np) \Rightarrow (\#b !>^* \#vp)
 \end{aligned}$$

The query looks for a node $\#a$ that dominates a sequence of nodes $\#np \#vp$, but only if no node on the path from $\#a$ to $\#np$ (the condition of the \forall clause) also dominates $\#vp$. With this restriction, it is made sure that $\#a$ really is the first common ancestor.

It is clear that this query cannot be handled by the node set definition syntax, because, unlike in Q2, the right-hand side of the implication tests for a non-local feature of the node.

For these more complex queries, we propose an extension that is functionally equivalent to the implication syntax, but also relies on the notion of node sets. If the right-hand side of a node set definition is enclosed in curly braces, the contents of the definition are interpreted as a *sub-query* rather than a node description, which are enclosed in square brackets. In this sub-query, all existential node variables from the outer query may be referenced. In the example query, the node set of all nodes on the path between $\#a$ and $\#np$ would be defined as

$$(13) \quad \%b:\{\#a >^* \%:[NT] \& \% >^* \#np\}.$$

The free-standing percent sign is used as the placeholder for those nodes that are going to be elements of the set⁶. With the node set $\%b$ defined this way, query 12 can be rephrased as

$$\begin{aligned}
 (14) \quad & \#np:[cat="NP"] .* \#vp:[cat="VP"] \& \\
 & \#a:[NT] >^* \#np \& \\
 & \#a >^* \#vp \& \\
 & \%b:\{\#a >^* \%:[NT] \& \% >^* \#np\} !>^* \#vp
 \end{aligned}$$

While this query will need more time for evaluation, the evaluation order is not essentially changed:

1. get all nodes that can match $\#np$, $\#vp$, $\#a$ and $\%$
2. for a given assignment, execute the sub-query that defines $\%b$

⁶This syntax is still subject to change.

3. evaluate all other constraints.

While this approach is easy to implement, it is also very slow. Greater speeds can be achieved by early pruning of the search space, for instance by employing additional indices or intelligent constraint evaluation reordering and caching.

With sub-queries, the semantics of node predicates should be changed as well. A node predicate that is invoked on a set of nodes will only return true if all nodes in the set satisfy the predicate. To build the set of all nodes that fulfill a certain condition, a sub-query should be used. For instance, the set of all NP nodes with exactly two children can be created with

(15) `%x:{%:[cat="NP"] & arity(%, 2)}`

6 Related Research

The Stockholm TreeAligner is unique in its ability to display parallel syntax trees and their phrase alignment. It can be seen as an advancement over tools like Cairo (Smith and Jahr, 2000) or I*Link (Merkel et al., 2003) which were developed for creating and visualizing word alignments but are unable to display trees. The TreeAligner is also related to Yawat and Kwipc (Germann, 2007), two software tools for the creation and display of sub-sentential alignments in matrices.

Our implementation of the query language owes a lot to TIGERSearch, whose groundbreaking work and robust implementation we gratefully acknowledge. But our work is also related to Emdros (Petersen, 2004), who has built his own system for parallel corpus searches. In Petersen (2005), the author of Emdros shows that his system can handle Q2 correctly and is faster than TIGERSearch. Other related query systems are NXT search/NiteQL (Heid et al., 2004; Evert and Voormann, 2003) and LPath (Bird et al., 2005) none of which is specifically geared towards parallel treebanks.

LPath is an extension of XPath to support querying syntactic trees that are encoded in XML info sets. All syntactical relations (dominance, precedence) are directly encoded in the XML structure and thus inherit the restrictions that apply to XML trees. Secondary layers of annotation, even simple ones like the secondary edge feature in TIGER-XML, are not supported. LPath also provides an extension to the XPath syntax allowing `not` modifiers on predicates, making it possible to express Q2 in a straight-forward manner.

NiteQL has been designed for corpora with multiple layers of annotation on one and the same text, but also heavily relies on XML to encode syntactical structure, with the same restrictions. According to Lai and Bird (2004), NiteQL cannot handle queries that need universal quantification. The authors of NiteQL also say

TigerSearch (sic) has a nice graphical interface and again supports structural operators missing in NQL. (Nite XML Toolkit Documentation, section 8.4)

but do not give any examples for the differences.

7 Conclusion and Outlook

We have shown how the TIGERSearch query language can be integrated in a parallel treebank exploration tool, the Stockholm TreeAligner. This tool allows for parallel queries over two treebanks combined with their alignment constraints and is thus a useful tool for cross-language comparisons and translation studies (cf. Lundborg et al. (2007) for an introduction).

The TIGERSearch query language lacks universal quantification and is thus unable to handle certain queries. We have added this functionality to the query language and implemented it in the TreeAligner. In this paper, we have also sketched a possible syntax for queries that contain more complex universal quantification.

The extensions described in section 3 will be part of the next release of the Stockholm TreeAligner later this year, and are already implemented in the development version. This release will also contain a browser and query interface for monolingual treebanks, similar to the one provided by TIGERSearch. While the functionality for monolingual queries is already (and has always been) provided by the underlying implementation, so far the focus of the TreeAligner has been on creating and querying parallel treebanks.

We will also follow up on the planned extensions sketched in section 5. However, the TIGER implementation in the Stockholm TreeAligner is still in a very early stage. Feature completeness and a faster query evaluation will become more important for the upcoming releases.

References

- Bird, S., Chen, Y., Davidson, S., Lee, H., and Zheng, Y. (2005). Extending XPath to support linguistic queries. In *Proc. of Programming Language Technologies for XML (PLANX)*, pages 35–46, Long Beach, California.
- Brants, S., Dipper, S., Hansen, S., Lezius, W., and Smith, G. (2002). The TIGER treebank. In *Proceedings of the Workshop on Treebanks and Linguistic Theories*, Sozopol.
- Evert, S. and Voormann, H. (2003). NQL – A query language for multi-modal language data. Technical report, IMS, University of Stuttgart, Stuttgart.
- Germann, U. (2007). Two tools for creating and visualizing sub-sentential alignments of parallel text. In *Proc. of The Linguistic Annotation Workshop at ACL 2007*, pages 121–124, Prague.
- Heid, U., Voormann, H., Milde, J.-T., Gut, U., Erk, K., and Padó, S. (2004). Querying both time-aligned and hierarchical corpora with NXT search. In *Proc. of The Fourth Language Resources and Evaluation Conference*, Lisbon.
- König, E. and Lezius, W. (2003). The TIGER language – A description language for syntax graphs, formal definition. Technical report, IMS, University of Stuttgart.

- Lai, C. and Bird, S. (2004). Querying and updating treebanks: A critical survey and requirements analysis. In *Proceedings of the Australasian Language Technology Workshop*.
- Lezius, W. (2002). *Ein Suchwerkzeug für syntaktisch annotierte Korpora*. PhD thesis, IMS, University of Stuttgart.
- Lundborg, J., Marek, T., Mettler, M., and Volk, M. (2007). Using the Stockholm TreeAligner. In *Proc. of The 6th Workshop on Treebanks and Linguistic Theories*, Bergen.
- Merkel, M., Petterstedt, M., and Ahrenberg, L. (2003). Interactive word alignment for corpus linguistics. In *Proc. of Corpus Linguistics 2003*, Lancaster.
- Petersen, U. (2004). Emdros: A text database engine for analyzed or annotated text. In *COLING '04: Proceedings of the 20th international conference on Computational Linguistics*, page 1190, Morristown, NJ, USA. Association for Computational Linguistics.
- Petersen, U. (2005). Evaluating corpus query systems on functionality and speed: TIGERSearch and Emdros. In *Proc. of The International Conference on Recent Advances in NLP*, Borovets.
- Smith, N. A. and Jahr, M. E. (2000). Cairo: An alignment visualization tool. In *Proc. of LREC-2000*, Athens.
- Volk, M., Lundborg, J., and Mettler, M. (2007). A search tool for parallel treebanks. In *Proc. of The Linguistic Annotation Workshop (LAW) at ACL*, Prague.