

**Hans Sandbu**

# Spatio-temporal Keyword search on RDF graphs

Master's Thesis in Computer Science, Spring 2020

Norwegian University of Science and Technology





## Abstract

This paper provides a template for writing a Master's Thesis (parts of it can also be used when writing a Fall Project Report). The template does not form a compulsory style that you are obliged to use, but rather provides a common starting point for all students. For a given thesis, tuning of the template may still be required, depending on the nature of the thesis and the author's writing style. Such tuning might involve moving a chapter to a section or vice versa, or removing or adding sections and chapters.

[If you write a Fall Project Report, it should normally focus on the background, related work (i.e., your literature study), and future work sections — with the “future work” section containing the plan for the Master's Thesis work to be carried out in the Spring. Architectural and experimental sections can also be included, but in preliminary versions. All those sections should of course be updated in the Master's Thesis, and adapted to the actual work carried out.]

Note that the template contains a lot of examples of how to write different parts of the thesis as well as how to cite authors and how to use LaTeX and BibTeX. Some of those examples might only be clear if you actually look at the LaTeX source itself.

The abstract is your sales pitch which encourages people to read your work, but unlike sales it should be realistic with respect to the contributions of the work. It should include:

- the field of research,
- a brief motivation for the work,
- what the research topic is,
- the research approach(es) applied, and
- contributions.

The abstract length should be roughly half a page of text. Do not include lists, tables or figures. Avoid abbreviations and references.

## Sammendrag

Husk at hvis du er en norsk student og skriver masteren din på engelsk, så *må* du lage et sammendrag på norsk.

(If you are a non-Norwegian student, it is not obligatory to include an abstract in Norwegian.)

## **Preface**

Spatiotemporal keyword search on RDF graphs.

Hans Sandbu  
Trondheim, 4th May 2020



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background and Motivation . . . . .	2
1.2	Goals and Research Questions . . . . .	2
1.3	Research Method . . . . .	2
1.4	Contributions . . . . .	2
1.5	Thesis Structure . . . . .	2
<b>2</b>	<b>Background Theory</b>	<b>3</b>
2.1	Definitions . . . . .	3
2.1.1	Knowledge Base . . . . .	3
2.1.2	Ontology . . . . .	3
2.1.3	Knowledge Graphs . . . . .	4
	Subject . . . . .	4
	Predicate . . . . .	4
	Object . . . . .	4
2.1.4	Facts and entities . . . . .	4
2.2	Existing knowledge graphs and ontologies . . . . .	5
2.2.1	Uses for Knowledge graphs . . . . .	5
2.2.2	Yago . . . . .	5
	Temporal data . . . . .	5
	Spatial data . . . . .	6
2.2.3	DBPedia . . . . .	6
2.2.4	Common tools and technology . . . . .	6
2.2.5	Jena . . . . .	6
2.3	Introducing Figures . . . . .	6
2.4	Introducing Tables in the Report . . . . .	6
<b>3</b>	<b>Related Work</b>	<b>7</b>
<b>4</b>	<b>Architecture</b>	<b>9</b>
4.1	Search . . . . .	9
4.1.1	BFS search . . . . .	9
4.2	Pruning . . . . .	11
4.2.1	Predicate pruning . . . . .	11
4.2.2	Unqualified place pruning . . . . .	11
4.2.3	Bounding . . . . .	12

<b>5</b>	<b>Experiments and Results</b>	<b>13</b>
5.1	Experimental Plan . . . . .	13
5.1.1	Time . . . . .	13
5.1.2	Scoring and ranking . . . . .	13
	BFS . . . . .	13
	BFS with pruning . . . . .	13
5.2	Experimental Setup . . . . .	13
5.3	Data set . . . . .	13
5.4	Queries . . . . .	14
5.5	Pruning . . . . .	14
5.5.1	Predicate pruning . . . . .	15
5.6	Experimental Results . . . . .	15
5.7	BFS search . . . . .	15
5.7.1	Time . . . . .	15
5.7.2	Accuracy . . . . .	16
<b>6</b>	<b>Evaluation and Discussion</b>	<b>19</b>
6.1	Evaluation . . . . .	19
6.2	Discussion . . . . .	19
<b>7</b>	<b>Conclusion and Future Work</b>	<b>21</b>
7.1	Contributions . . . . .	21
7.2	Future Work . . . . .	21
	<b>Bibliography</b>	<b>23</b>
	<b>Appendices</b>	<b>25</b>



# List of Figures

2.1	Trondheim as a subject in YAGO . . . . .	6
-----	--	---



# List of Tables



# 1 Introduction

Creating efficient structures for computers to interpret information opens up many new areas for information retrieval. One method is structuring the information as a graph with vertices connected by edges ( $G(v, e)$ ). In such a structure each vertex ( $v$ ) can hold a small piece of information, and the edge ( $e$ ) represents a relation. When information is modeled in this way, it can be called a knowledge graph (KG).

Graph based models allows complex relations to be modeled. Finding relations in a graph is also computationally cheaper than other models. Using standardized models such as the Resource Description Framework (RDF), which is created specifically for graph based data, has made it possible to develop efficient storage and retrieval solutions for graph data. One such solution is the framework Jena Wilkinson et al. (2004) for Java.

Using RDF projects like YAGO (Yet Another Great Ontology) Suchanek et al. (2007) have created graphs containing spatial and temporal data (spatiotemporal). This makes it possible to apply reason to the dataset, and derive new information from what is already there. This is a feature that typically separates a KG from an ontology Färber et al. (2017). In a graph, spatiotemporal data will be described using at least 2 vertices, one for a start date, and one for a place. When modeling time intervals, multiple vertices are needed.

There are few applications using RDF as an information storage system. Some of the more well known are DBPedia, YAGO, and Creative Commons. These applications often consists of a large linked dataset, or in the case of Creative Commons, it is used for embedding licenses. The applications built with RDF usually don't have much utility. A common utility built on top of RDF data sets are tools used to describe social relations, though there have been developed applications such as MusicBrainz that use RDF for relations between songs, artists, albums and other tags given to entities.

Creating methods that allows for easy access to information in RDF data makes it possible to create applications with more utility. Keyword and text search in RDF data is still something being developed and improved. Using frameworks such as Jena, it is possible to build solutions for keyword search, and spatiotemporal search. In this thesis (paper?) methods for efficient keyword search on spatial and temporal data will be explored.

## 1.1 Background and Motivation

## 1.2 Goals and Research Questions

Research questions?

RQ1: How can spatio-temporal data be integrated into exiting keyword query methods for RDF data.

RQ2: What methods can be used to create more effective queries on RDF data.

RQ3: How do spatial and temporal RDF query methods differ from from other query methods.

## 1.3 Research Method

What methodology will you apply to address the goals: theoretic/analytic, model/abstraction or design/experiment? This section will describe the research methodology applied and the reason for this choice of research methodology.

## 1.4 Contributions

This section just provides a brief summary of the main contributions of the work. The main description of the contributions will come in Section 7.1, after the results are presented. (Hence Section 1.4 can also be left out, leaving the discussion completely to Section 7.1.)

The format of this section will generally be as follows: *Donec non turpis nec neque egestas faucibus nec id neque. Etiam consectetur, odio vitae gravida tempus, diam velit sagittis turpis, a molestie ligula tellus at nunc. Nam convallis consequat vestibulum. Proin dolor neque, dapibus a pellentesque a, commodo a nibh.*

1. *Lorem ipsum dolor sit amet, consectetur adipiscing elit.*
2. *Lorem ipsum dolor sit amet, consectetur adipiscing elit.*
3. *Lorem ipsum dolor sit amet, consectetur adipiscing elit.*

where the items on the list briefly describe the key contributions.

## 1.5 Thesis Structure

This section provides the reader with an overview of what is coming in the next chapters. You want to say more than what is explicit in the chapter name, if possible, but still keep the description short and to the point.

## 2 Background Theory

### 2.1 Definitions

- definitions (Find names, and some definitions)
- Result tree: Result  $r$  for a given query  $q$  with tokens  $Q_t$  on an RDF graph  $G\langle V, E \rangle$ , where the result forms a minimum spanning tree  $T_m\langle V', E' \rangle$ ,  $V'$  contains a set of tokens  $T_t$ ,  $T_m$  is rooted in a place vertex  $p$ , so that  $V' \subseteq V$ ,  $E' \subseteq E$ , and  $T_t \subseteq Q_t$ 
  - Possible results: Collection of all result trees  $r_t$  for a given query  $q$
  - Root nodes: Node with a `yagoGeoEntity` class, used as the starting node when traversing the graph.
  - Accuracy: Score given to a result tree, where score is based on nodes  $n$  in the result tree, node distance  $nd$  from root, query  $q$ , and node tokens  $nt$ , and number of query words hit  $h$  so that  $\text{sum}(nt/\text{len}(q))/\text{sum}(n^*(nd+1)) = \text{score}$  and  $nt (= q$
  - Queries: Any set of words from in English set of words.
  - Pruning: Removal of unnecessary nodes  $n$ , where the removed nodes tokens  $nt$  is not part of the query tokens  $qt$ ,  $nt \not\subseteq Q_t$ .

Knowledge base, ontology and knowledge graph are terms with multiple definitions, and are often used interchangeably. The definitions here are used to clarify what they mean in this paper, and is not a definitive definition.

#### 2.1.1 Knowledge Base

The term knowledge base can have many definitions, often because the terms knowledge base, ontology, and knowledge graph is used interchangeably Ehrlinger and Wöß (2016). In this paper the term knowledge graph will be defined as follows: “A knowledge base is a data set with some formal semantics. This could include multiple axioms, definitions, rules, facts, statements, and other primitives.” Davies et al. (2006)

#### 2.1.2 Ontology

Like knowledge base, ontology does not have one clear definition. This is because the term have many different interpretations, and is often confused with other terms Feilmayr and Wöß (2016). Most definitions of ontologies say that ontologies represent a schema,

## 2 Background Theory

basic theory, or conceptualization of a domain, which knowledge bases do not. Davies et al. (2006) To differentiate ontologies and knowledge bases in this paper the definition of an ontology will be “An ontology is an extended knowledge base that allows for semantic modeling of knowledge.” With this definition, ontologies can be considered a specialized form of knowledge bases.

### 2.1.3 Knowledge Graphs

A broad definition of knowledge graph is “A knowledge graph acquires and integrates information into an ontology and applies a reasoner to derive new knowledge.” Ehrlinger and Wöß (2016) This definition encompasses multiple technologies. In this paper we will use the more specific definition that fits the data used “We define a Knowledge Graph as an RDF graph. An RDF graph consists of a set of RDF triples where each RDF triple  $(s, p, o)$  is an ordered set of the following RDF terms: a subject  $s \in U \cup B$ , a predicate  $p \in U$ , and an object  $U \cup B \cup L$ . An RDF term is either a URI  $u \in U$ , a blank node  $b \in B$ , or a literal  $l \in L$ ” Färber et al. (2017)

#### Subject

In a RDF triple the subject is a URI or a blank node. The URI when used in the subject identifies an entity, or is an alias, different language or other variation of an entity. This subject can have relations to objects describing the same entity.

#### Predicate

Predicates are always a URI. URIs used for predicates differ from the ones used for subjects and objects in that predicates are of a type. A type is an identifier used to describe the relation between the subject and object.

#### Object

Objects have the widest range of possible entries. Like subjects and predicates, objects can be URIs. Object URIs can be entity identifiers like subjects, they can be class identifiers, or they can contain some data and a data type. Blank nodes are just that, blank, and literals are an atomic value.

### 2.1.4 Facts and entities

A fact is a term often used when describing knowledge bases, ontologies and knowledge graphs. Usually a fact is the smallest piece of information in such a system. In a knowledge graph this is a single RFD triple. Another term often used is entity. An entity is a collection of facts, usually from the same article. Entities can be linked together through facts. In such a fact, one entity is used as the subject, the predicate describes the relation, and the object is the other entity. In such a relation both entities will be URIs.



## 2.2 Existing knowledge graphs and ontologies

Currently two of the largest open technologies for knowledge graphs are Yago and DBPedia. Both these projects use automatic extraction from Wikipedia to create the graph, but differ in the ontology used to build the graphs. Yago also includes data from WordNet and GeoNames to accurately assign entities to classes. Both projects use RDF triples to create a knowledge graph.

### 2.2.1 Uses for Knowledge graphs

One of the uses of knowledge graphs today is to find and display a info box in search engines. This information is a compact set of facts that tries to fit the search query. Because of the graph structure of knowledge graphs the information in the info box can be adapted to the query by choosing the predicates and related facts closest related to the query. This makes it possible to create a set of information that can give the user a quick overview of the information retrieved by the query.

### 2.2.2 Yago

Yago is an acronym for Yet Another Great Ontology, and is main data source in this paper. The project describes it self as a knowledge base Suchanek et al. (2007) and an ontology Mahdisoltani et al. (2013), but is often described as a knowledge graph by others. In this paper Yago is described as a knowledge graph.

### Temporal data

Yago have many entities with facts describing date and spatial data. Suchanek et al. (2007) Date facts follows the ISO 8601 format, YYYY-MM-DD, and introduces # as a wildcard symbol. A fact can only hold information on a single point in time, and uses yagoDate as a data type in addition to information on the date. Suchanek et al. (2007) In a date fact, the object holds the date information, and the predicate describes a connection between subject and date. “Nidaros\_Cathedral wasCreatedOnDate 1300-##-##.” In this example the predicate wasCreatedOn is used to describe a relation between the subject and a date.

To describe a time span two facts are required. One of the facts describes a start date, and the second an end date. Since an entity can have multiple date facts connected, all date predicates are also assigned to a class. Start dates are assigned to a predicate with a type that has a “creation” class, such as “StartedOn-Date”. End dates are assigned to “destruction” type predicates. This makes it possible to deduce a time span for a given subject and predicate combination. Suchanek et al. (2007)

## 2 Background Theory

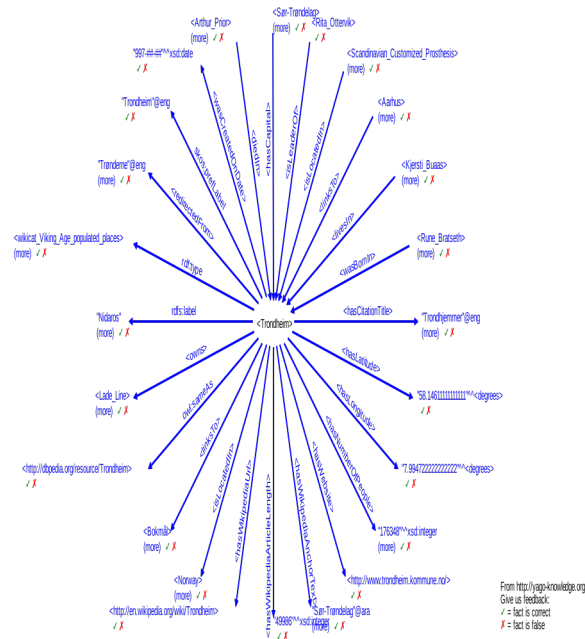


Figure 2.1: Trondheim as a subject in YAGO

## Spatial data

Yago only contains permanent spatial data for entities on earth. This means that entities like cities, buildings, rivers and mountains are given a spatial dimension. In addition, events, people, groups and artifacts can be given a spatial dimension by relating the entity to a specific place. All spatial facts must have a predicate that fall under the `yagoGeoEntity` class, and all objects used in a fact with a `yagoGeoEntity` must have a relation containing both “hasLatitude” and “hasLongitude”.

### 2.2.3 DBPedia

### 2.2.4 Common tools and technology

### 2.2.5 Jena

Apache Jena is a Java framework created for linked data stores. This is used to create a store that can be queried, and traversed using SPARQL as a query language.

## 2.3 Introducing Figures

## 2.4 Introducing Tables in the Report

### 3 Related Work

There has been done some work keyword search on RDF graphs. This includes Tran et al. (2009) and Elbassuoni and Blanco (2011). None of these takes spatial or temporal data into consideration. The methods for keyword queries outlined in the papers are based on a combination of indexing the graph and traversing subgraphs. The indexing is similar in both papers, but retrieval, scoring and traversal methods varies. One method of retrieval and scoring that is used as a baseline is a breadth first search for traversal and a minimal subgraph for scoring.

In paper Shi et al. (2016) some methods for indexing, searching and ranking keyword searches on spatial RDF graphs is proposed. These methods builds upon the methods outlined in Tran et al. (2009); Elbassuoni and Blanco (2011). The most substantial difference is the use of R-trees to index the spatial dimension of the graph. This is done so that a subgraph containing the keywords can be retrieved given a location. This method starts by finding the closest spatial vertex and uses that as the root for the subgraph. Subgraphs rooted in a spatial node requires different ranking, and there is proposed a set of rules in the paper.

R-trees are designed around the spatial dimension, but some work on modifications to include a temporal dimension has been done. Some research done on spatiotemporal tree structures include Tao et al. (2003); Šaltenis et al. (2000). Most of the spatiotemporal trees are however created to be able to query and index moving or frequently updated objects. To include a temporal dimension some differences in the tree structures are made. Most approaches builds on the R-tree, but modifies the tree in different aspects. For the purpose of this thesis a different tree structure is not needed.



# 4 Architecture

## 4.1 Search

### 4.1.1 BFS search

The most basic method of finding a match for keywords is using a breadth first search (BFS) He et al. (2007). For each keyword in the query the algorithm will find all vertices that contain the keyword. From that set of vertices the BFS search finds the first vertex that can connect all the vertices in the set. The vertex that connects the the rest of the set is the one that best fits the keywords in the search. There is however no guarantee that this set of vertices contains any spatial or temporal data. If this is the case, the search will continue until the a vertex containing spatial and temporal data is found. The first set of vertices containing all the keywords, a place, and time will be the best result using this search. The BFS search can however contain multiple times, or multiple places. To determine what time or what place best fits the query, a separate ranking algorithm can be used. BFS is also a time consuming algorithm and is used as a proof of concept and baseline in this thesis.

---

**Algorithm 1:** GetFullResultTree( $p, t, Q_t$ )

---

**Result:** Set containing all nodes with at least one keyword

Queue  $Q$  ADD( $p$ );

Threshold  $t$ ;

Set  $n$  ;

**while**  $Q \neq \emptyset$  **do**

    Clear( $n$ )  $e = \text{POP}(Q)$  **if**  $\text{GetDistance}(e) > t$  **then**  
        | continue;

**end**

**foreach**  $\text{Term } t \in Q_t$  **do**

**if**  $t \in e$  **then**

            |  $p.\text{AddMatchChild}(e)$ ;

**end**

**end**

**if**  $Q_t \subseteq e$  **then**

        |  $t = \text{GetDistance}(e)$ ;

**end**

$n = \text{GetNeighbors}(e)$  ;

$Q$  ADD( $n$ );

**end**

---

## 4 Architecture

The first part of the search is to find a set of root nodes for sub-graphs. This is done by collecting all neighbors from the place queried. The neighbors can be connected by any predicate in the first developed algorithm, but as explain in 4.2 this selecting specific predicates can greatly increase speed by limiting the amount of nodes traversed. The BFS algorithm described above is used for each of the possible roots of subgraphs, stopping when all keywords are matched, or the depth of the graph exceeds the three or exceeds the currently shallowest subgraph that has matched all keywords. The reason for limiting the depth of subgraphs is to ensure at least a partial hit within a reasonable time. Limiting the depth of subgraphs to the current shallowest subgraph is done because no a deeper graph cannot be a more accurate hit.

- Tokens
- Objects

When traversing the graph and finding a match, a node object is created. This object is used to keep track of the pseudo hierarchy in the graph. All objects contain information on the depth of the node, matched query terms, and relation to parent and children, if any. In addition root objects contains a list of all query terms hit in the sub graph. If a child node is found within multiple subgraphs, a new object will be created representing the node for each subgraph. This creates some objects that are nearly identical, but with different relations and possible different depth.

---

**Algorithm 2:** minimum spanning graph

---

1: **procedure** MINIMUM For each node in hits: for each minNode in min list: if node has same hits and higher depth than minNode: Break if node and minNode is equal in depth and hits: continue if node contains other words than minNode: add node remove potential common words from minNode

---

- Start by finding all trees with the same minimum depth for each node (Roots for different trees) related to the place in the query.

After traversing the main graph we have found all subgraphs containing at least one query term. These subgraphs contain many nodes which hit the same terms. Before ranking the subgraphs the minimum spanning graph needs to be found. The minimum spanning tree is found using a greedy algorithm that iterates through all nodes in the subgraphs, then keeping the nodes containing the most terms, and lowest depth. The minimum tree will contain as many terms as possible, a term will only be found in one node, and the graph will be as shallow as possible.

- Find the minimum spanning tree for each of the root nodes containing the maximum amount of the query words.
  - Rank based on 1. Query words hit 2. Nodes in the tree 3. depth of the tree
- The final piece is to rank the minimum subgraphs. This is done using the nr. 2 formula found in Shi et al. (2016).

## 4.2 Pruning

When using the BFS search method, all possible spatial vertices close to the queried place will be explored. This is expensive and many of the vertices will be irrelevant. Pruning the potential place vertices will reduce the amount of sub-graphs traversed, and will in turn reduce the overall time used to find results for the query.

### 4.2.1 Predicate pruning

When traversing the graph a lot of unnecessary predicates are followed, resulting in many extra nodes added to the search. Specifying a set of predicates that contain the relevant information can greatly increase the speed of the algorithm, and also keep the memory requirements a lot lower. When selecting predicates for traversal the information expected from the search should be the top priority. Because of this, all predicates that may contain spatial or temporal data should be kept.

When using the entirety of the Yago data set, most nodes are highly connected. Many of the links in the graph are from predicates such as “linksTo” or “redirectedFrom”. These predicates creates a highly connected graph, and ensures a hit within a few nodes of the start. The same predicates will also often add the same nodes multiple times, create circular graphs, and take up unnecessary CPU power and memory.

When pruning predicates there are two methods that are possible to implement. The first will remove the predicates that contain little or no new information, such as “linksTo” or “redirectedFrom” mentioned above. This will still keep the graph connected, and keeps the predicates containing more useful information.

An other method of pruning is to create a list of predicates to be followed. This can drastically reduce the connections in the graph, but the results will only contain information relevant to the query. When preselecting predicates there is a much greater chance of not finding a match for a query. In addition a lot of metadata could be lost, if the metadata predicates are not added to the list of predicate to be explored.

### 4.2.2 Unqualified place pruning

For any spatial search, a lot of places will be found. The graph also contains information on the relationship between places. To find the best possible matches for a spatial query, only places of a higher resolution should be chosen. This means that places of similar or smaller expanse should be allows to be queried, e.g. a query of Boston can return the entirety of Boston, close to Boston, or within Boston. Massachusetts has multiple predicates linking it to Boston, but should not be queried because the information

## *4 Architecture*

returned would be less detailed than what is queried.

### **4.2.3 Bounding**

When selecting places or times for a query, the boundaries should be set so that they encompass the entirety of the queried time and or place.



# 5 Experiments and Results

## 5.1 Experimental Plan

The goal of the experiments is to gather data for comparison of search methods. when comparing search methods, speed and accuracy will be the metrics used. Three methods for searching is used, one following all gathering all objects connected to an input subject, regardless of the predicate. One method following Yago predicates leading directly to a Yago URI. And the final method will follow the same predicates as the second, and will also add a condition where all nodes must match at least one query word.

### 5.1.1 Time

When timing the search, two different times will be measured, time for each query, and time for each root node. In addition to these, avg. nodes visited for each root will also be used. Taking the average of these over a large input set should generate an appropriate result. Both time metrics measure how fast results are found, so the results should be similar as long as the input data does not contain a high number of highly connected nodes.

The timing of the algorithm should be the same as any breath first search,  $\Theta(V + E)$  so that the best case is visiting only the first vertex, and the worst case is traversing the entire graph.

### 5.1.2 Scoring and ranking

Each query is given a score. This score will be used to see how pruning, and difference in predicates can lead to differences in the result tree. Two metrics for accuracy will be used, avg. accuracy for each result and avg. highest accuracy for each query.

### BFS

### BFS with pruning

## 5.2 Experimental Setup

### 5.3 Data set

When selecting a data set an important feature was accurate spatial and temporal data. For this Yago was selected. Yago allows for download of subsets of data, and contains

## 5 Experiments and Results

all the necessary data. The high accuracy in Yago was also a benefit of this data set. From Yago, the data selected was the entire Yago taxonomy, the entire Yago core, and from Genomes the sets GeonamesOnlyData, GeonamesClassIds, GeonamesGlosses, GeonamesTypes, and GeonamesClasses were selected. These sets makes it possible to build a complete graph of the entities in Yago, but keeps the disk usage as low as possible.

- Short description of each of the parts
- preprocessing: replace non-unicode chars, replacing space with underscore in URIs, terminate triplet where missing termination, remove double quotes in URIs, remove backslash unless legal escape sequence
- tdbloader for persistent queryable storage
- Structure of graph? Here or in Yago description?

### 5.4 Queries

When selecting the keywords used in queries, a semi random selection method was used at first. This method created three pools of words, rare words (less than 1000 occurrences) uncommon words (less than 100 000 occurrences) and common words (more than 100 000 occurrences). The three pools where created by stemming all words, then counting occurrences. From each pool of words, a set of 100 where randomly selected, and from those 100, 10 where selected manually to ensure a range of occurrences, and to ensure no stop words, misspelled words or other errors were in the final set of keywords.

A second set of words were also chosen to ensure a larger hit rate. This set was a random selection of 20 words from the top 150 words. From the 20 random words 10 where manually selected, the 150 word number was selected based on Zipfs law Piantadosi (2014).

When selecting places for spatial queries, a set of 20 random places were selected from the YagoGeonamesOnlyData set, and from that set, 5 were chosen manually for the final set. In the final set, three places where added manually, Oslo, London, and New York City were added to ensure variation in placement and node connections.

Generating random queries can create results that will not have any hits. Find a good method of creating queries.

### 5.5 Pruning

A simple method to improve performance when searching in large data sets is to remove data that is unnecessary to search.

### 5.5.1 Predicate pruning

Predicate pruning is the process of eliminating predicates from the search. The taxonomy of Yago creates assigns predicates to one or more group or class. By selecting specific groups or classes it is possible to follow the types of data searched for, while eliminating a lot of other data. When pruning based on classes or groups the type of data returned must be known beforehand, to ensure that no possible hits is overlooked.

Another possible method for pruning is to specify the exact predicates to follow. This will remove a lot of data, but is possibility when searching for something specific. This will also greatly increase the speed of the search, and reduce the memory usage for the search.

## 5.6 Experimental Results

### 5.7 BFS search

BFS search without any pruning and search depth of 3 were unable to finish due to memory constraints and taking too long. Search with a depth of 2 may finish, depending on on the connectedness of the nodes discovered, and the amount of root nodes found in the search.

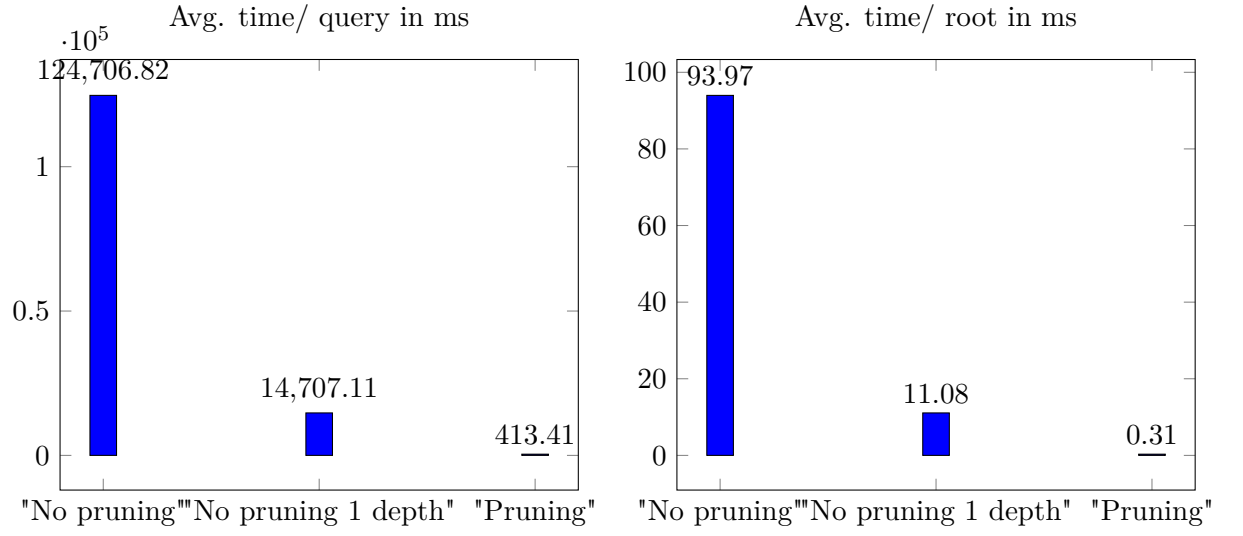
#### 5.7.1 Time

When running the query, two sets are displayed, one with mostly random input data, and one with more closely selected data.

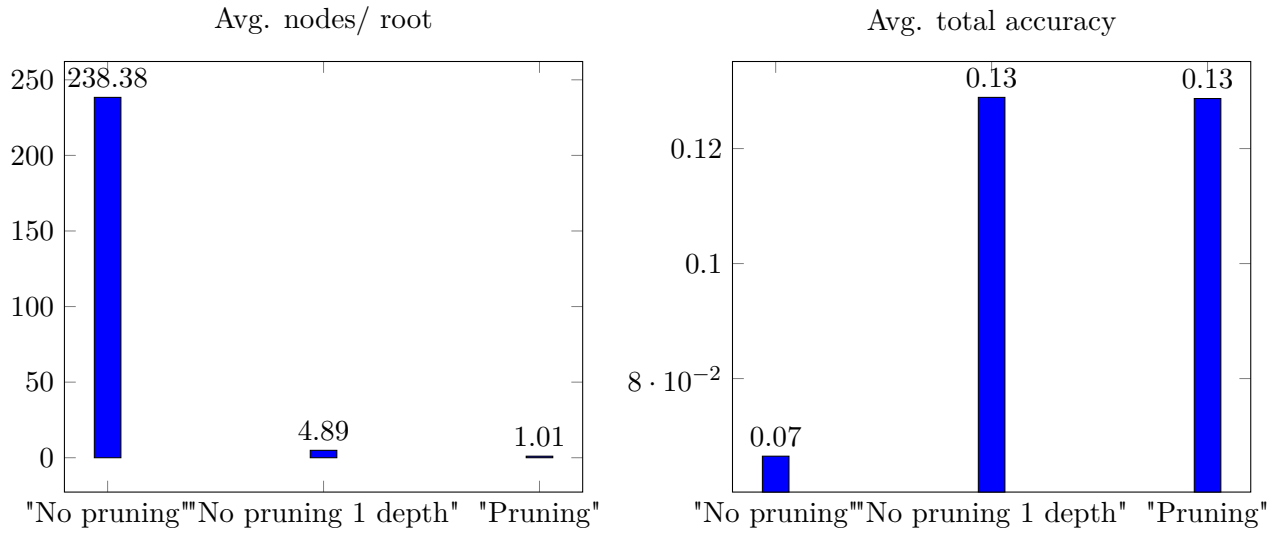
Because most of the randomly selected places had few possible root nodes, the table contains only the manually selected.

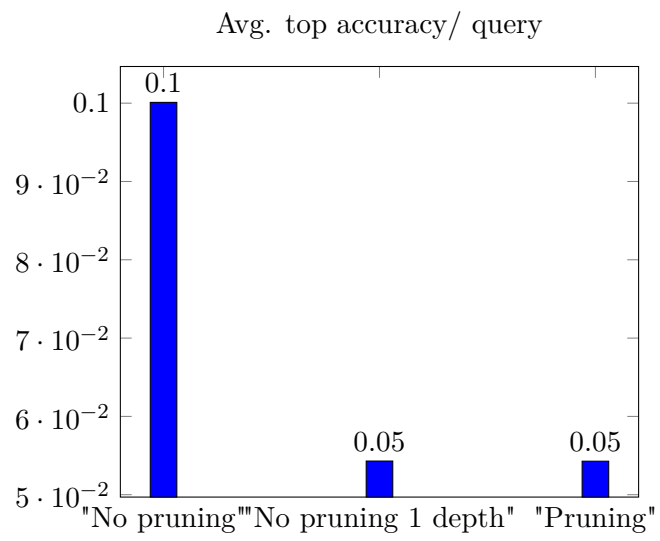
- time/query
- time/root
- roots found/ query
  
- accuracy/ result
- avg highest accuracy/ query

## 5 Experiments and Results



### 5.7.2 Accuracy







# 6 Evaluation and Discussion

## 6.1 Evaluation

## 6.2 Discussion

When traversing an rdf graph, the predicates chosen will have a great effect on the time used. Knowing the data searching through will help when choosing the predicate to follow. When choosing what predicates to use in the search, the goal should be to minimize the amount of unnecessary nodes found and followed. When removing predicates, the obvious downside is the possibility of missing some results, and decreasing the accuracy.

- limitations
- Effectiveness of pruning
- choice of query terms
- Root nodes





# 7 Conclusion and Future Work

## 7.1 Contributions

## 7.2 Future Work



# Bibliography

- John Davies, Rudi Studer, and Paul Warren. *Semantic Web technologies: trends and research in ontology-based systems*. John Wiley & Sons, 2006.
- Lisa Ehrlinger and Wolfram Wöß. Towards a definition of knowledge graphs. 09 2016.
- Shady Elbassuoni and Roi Blanco. Keyword search over rdf graphs. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management, CIKM '11*, pages 237–242, New York, NY, USA, 2011. ACM. ISBN 978-1-4503-0717-8. doi: 10.1145/2063576.2063615. URL <http://doi.acm.org/10.1145/2063576.2063615>.
- Christina Feilmayr and Wolfram Wöß. An analysis of ontologies and their success factors for application to business. *Data & Knowledge Engineering*, 101:1 – 23, 2016. ISSN 0169-023X. doi: <https://doi.org/10.1016/j.datak.2015.11.003>. URL <http://www.sciencedirect.com/science/article/pii/S0169023X1500110X>.
- Michael Färber, Frederic Bartscherer, Carsten Menne, and Achim Rettinger. Linked data quality of dbpedia, freebase, opencyc, wikidata, and yago. *Semantic Web*, 9:1–53, 03 2017. doi: 10.3233/SW-170275.
- Hao He, Haixun Wang, Jun Yang, and Philip S. Yu. Blinks: Ranked keyword searches on graphs. In *Proceedings of the 2007 ACM SIGMOD International Conference on Management of Data, SIGMOD '07*, pages 305–316, New York, NY, USA, 2007. ACM. ISBN 978-1-59593-686-8. doi: 10.1145/1247480.1247516. URL <http://doi.acm.org/10.1145/1247480.1247516>.
- Farzaneh Mahdisoltani, Joanna Biega, and Fabian M. Suchanek. YAGO3: A Knowledge Base from Multilingual Wikipedias. In *CIDR*, Asilomar, United States, January 2013. URL <https://hal-imt.archives-ouvertes.fr/hal-01699874>.
- Steven T Piantadosi. Zipf’s word frequency law in natural language: A critical review and future directions. *Psychon Bull Rev*, 21:1112–1130, 2014.
- Jieming Shi, Dingming Wu, and Nikos Mamoulis. Top-k relevant semantic place retrieval on spatial rdf data. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD '16*, pages 1977–1990, New York, NY, USA, 2016. ACM. ISBN 978-1-4503-3531-7. doi: 10.1145/2882903.2882941.

## Bibliography

- Fabian M. Suchanek, Gjergji Kasneci, and Gerhard Weikum. Yago: A Core of Semantic Knowledge. In *16th International Conference on the World Wide Web*, pages 697–706, 2007.
- Yufei Tao, Dimitris Papadias, and Jimeng Sun. The tpr\*-tree: An optimized spatio-temporal access method for predictive queries. In *Proceedings of the 29th International Conference on Very Large Data Bases - Volume 29*, VLDB '03, pages 790–801. VLDB Endowment, 2003. ISBN 0-12-722442-4. URL <http://dl.acm.org/citation.cfm?id=1315451.1315519>.
- T. Tran, H. Wang, S. Rudolph, and P. Cimiano. Top-k exploration of query candidates for efficient keyword search on graph-shaped (rdf) data. In *2009 IEEE 25th International Conference on Data Engineering*, pages 405–416, March 2009. doi: 10.1109/ICDE.2009.119.
- Kevin Wilkinson, Craig Sayers, Harumi Kuno, Dave Reynolds, and Jena Database. Efficient rdf storage and retrieval in jena2. *Proceedings of SWDB*, 02 2004.
- Simonas Šaltenis, Christian Jensen, Christian (codirector, Michael Böhlen, Curtis Dyreson, Heidi Gregersen, Dieter Simonas, Janne Skyt, Giedrius Slivinskas, Kristian Torp, Richard (codirector, Bongki Moon, Michael Soo, Amazon Com, and Andreas Timeconsult. R-tree based indexing of general spatio-temporal data. 01 2000.

# Appendices