

DBpedia SPARQL Benchmark – Performance Assessment with Real Queries on Real Data

Mohamed Morsey, Jens Lehmann, Sören Auer, and Axel-Cyrille Ngonga Ngomo

Department of Computer Science, University of Leipzig
Johannisgasse 26, 04103 Leipzig, Germany.
{lastname}@informatik.uni-leipzig.de

Abstract. Triple stores are the backbone of increasingly many Data Web applications. It is thus evident that the performance of those stores is mission critical for individual projects as well as for data integration on the Data Web in general. Consequently, it is of central importance during the implementation of any of these applications to have a clear picture of the weaknesses and strengths of current triple store implementations. In this paper, we propose a generic SPARQL benchmark creation procedure, which we apply to the DBpedia knowledge base. Previous approaches often compared relational and triple stores and, thus, settled on measuring performance against a relational database which had been converted to RDF by using SQL-like queries. In contrast to those approaches, our benchmark is based on queries that were actually issued by humans and applications against existing RDF data not resembling a relational schema. Our generic procedure for benchmark creation is based on query-log mining, clustering and SPARQL feature analysis. We argue that a pure SPARQL benchmark is more useful to compare existing triple stores and provide results for the popular triple store implementations Virtuoso, Sesame, Jena-TDB, and BigOWLIM. The subsequent comparison of our results with other benchmark results indicates that the performance of triple stores is by far less homogeneous than suggested by previous benchmarks.¹

1 Introduction

Triple stores, which use IRIs for entity identification and store information adhering to the RDF data model [9] are the backbone of increasingly many Data Web applications. The RDF data model resembles directed labeled graphs, in which each labeled edge (called *predicate*) connects a *subject* to an *object*. The intended semantics is that the *object* denotes the value of the *subject's* property *predicate*. With the W3C SPARQL standard [17] a vendor-independent query language for the RDF triple data model exists. SPARQL is based on powerful graph matching allowing to bind variables to fragments in the input RDF graph. In addition, operators akin to the relational joins, unions, left outer joins, selections and projections can be used to build more expressive queries [18]. It is evident that the performance of triple stores offering a SPARQL query interface is mission critical for individual projects as well as for data integration on the Web in

¹ This work was supported by a grant from the European Union's 7th Framework Programme provided for the project LOD2 (GA no. 257943).

general. It is consequently of central importance during the implementation of any Data Web application to have a clear picture of the weaknesses and strengths of current triple store implementations.

Existing SPARQL benchmark efforts such as LUBM [16], BSBM [4] and SP² [18] resemble relational database benchmarks. Especially the data structures underlying these benchmarks are basically relational data structures, with relatively few and homogeneously structured classes. However, RDF knowledge bases are increasingly heterogeneous. Thus, they do not resemble relational structures and are not easily representable as such. Examples of such knowledge bases are curated bio-medical ontologies such as those contained in Bio2RDF [2] as well as knowledge bases extracted from unstructured or semi-structured sources such as DBpedia [10] or LinkedGeoData [1]. DBpedia (version 3.6) for example contains 289,016 classes of which 275 classes belong to the DBpedia ontology. Moreover, it contains 42,016 properties, of which 1335 are DBpedia-specific. Also, various datatypes and object references of different types are used in property values. Such knowledge bases can *not* be easily represented according to the relational data model and hence performance characteristics for loading, querying and updating these knowledge bases might potentially be fundamentally different from knowledge bases resembling relational data structures.

In this article, we propose a generic SPARQL benchmark creation methodology. This methodology is based on a flexible data generation mimicking an input data source, query-log mining, clustering and SPARQL feature analysis. We apply the proposed methodology to datasets of various sizes derived from the DBpedia knowledge base. In contrast to previous benchmarks, we perform measurements on *real* queries that were issued by humans or Data Web applications against existing RDF data. We evaluate two different methods *data generation* approaches and show how a representative set of resources that preserves important dataset characteristics such as indegree and outdegree can be obtained by sampling across classes in the dataset. In order to obtain a representative set of *prototypical queries* reflecting the typical workload of a SPARQL endpoint, we perform a query analysis and clustering on queries that were sent to the official DBpedia SPARQL endpoint. From the highest-ranked query clusters (in terms of aggregated query frequency), we derive a set of 25 SPARQL query templates, which cover most commonly used SPARQL features and are used to generate the actual benchmark queries by parametrization. We call the benchmark resulting from this dataset and query generation methodology *DBPSB* (i.e. DBpedia SPARQL Benchmark). The benchmark methodology and results are also available online². Although we apply this methodology to the DBpedia dataset and its SPARQL query log in this case, the same methodology can be used to obtain application-specific benchmarks for other knowledge bases and query workloads. Since the DBPSB can change with the data and queries in DBpedia, we envision to update it in yearly increments and publish results on the above website. In general, our methodology follows the four key requirements for domain specific benchmarks are postulated in the Benchmark Handbook [8], i.e. it is (1) relevant, thus testing typical operations within the specific domain, (2) portable, i.e. executable on different platforms, (3) scalable, e.g. it is possible to run the benchmark on both small and very large data sets, and (4) it is understandable.

² <http://aksw.org/Projects/DBPSB>

We apply the DBPSB to assess the performance and scalability of the popular triple stores *Virtuoso* [7], *Sesame* [5], *Jena-TDB* [15], and *BigOWLIM* [3] and compare our results with those obtained with previous benchmarks. Our experiments reveal that the performance and scalability is by far less homogeneous than other benchmarks indicate. As we explain in more detail later, we believe this is due to the different nature of DBPSB compared to the previous approaches resembling relational databases benchmarks. For example, we observed query performance differences of several orders of magnitude much more often than with other RDF benchmarks when looking at the runtimes of individual queries. The main observation in our benchmark is that previously observed differences in performance between different triple stores amplify when they are confronted with actually asked SPARQL queries, i.e. there is now a wider gap in performance compared to essentially relational benchmarks.

The remainder of the paper is organized as follows. In Section 2, we describe the dataset generation process in detail. We show the process of query analysis and clustering in detail in Section 3. In Section 4, we present our approach to selecting SPARQL features and to query variability. The assessment of four triple stores via the DBPSB is described in Sections 5 and 6. The results of the experiment are discussed in Section 7. We present related work in Section 8 and conclude our paper in Section 9.

2 Dataset Generation

A crucial step in each benchmark is the generation of suitable datasets. Although we describe the dataset generation here with the example of DBpedia, the methodology we pursue is dataset-agnostic.

The data generation for DBPSB is guided by the following requirements:

- The DBPSB data should resemble the original data (i.e., DBpedia data in our case) as much as possible, in particular the large number of classes, properties, the heterogeneous property value spaces as well as the large taxonomic structures of the category system should be preserved.
- The data generation process should allow to generate knowledge bases of various sizes ranging from a few million to several hundred million or even billion triples.
- Basic network characteristics of different sizes of the network should be similar, in particular the in- and outdegree.
- The data generation process should be easily repeatable with new versions of the considered dataset.

The proposed dataset creation process starts with an input dataset. For the case of DBpedia, it consists of the datasets loaded into the official SPARQL endpoint³. Datasets of multiple size of the original data are created by duplicating all triples and changing their namespaces. This procedure can be applied for any scale factors. While simple, this procedure is efficient to execute and fulfills the above requirements.

For generating smaller datasets, we investigated two different methods. The first method (called “rand”) consists of selecting an appropriate fraction of all triples of

³ Endpoint: <http://dbpedia.org/sparql>, Loaded datasets: <http://wiki.dbpedia.org/DatasetsLoaded>

Dataset	Indegree w/ literals	Outdegree w/ literals	Indegree w/o literals	Outdegree w/o literals	No. of nodes	No. of triples
Full DBpedia	5.45	30.52	3.09	15.57	27,665,352	153,737,776
10% dataset (seed)	6.54	45.53	3.98	23.05	2,090,714	15,267,418
10% dataset (rand)	3.82	6.76	2.04	3.41	5,260,753	16,739,055
50% dataset (seed)	6.79	38.08	3.82	18.64	11,317,362	74,889,154
50% dataset (rand)	7.09	26.79	3.33	10.73	9,581,470	78,336,781

Table 1: Statistical analysis of DBPSB datasets.

the original dataset randomly. If RDF graphs are considered as small world graphs, removing edges in such graphs should preserve the properties of the original graph. The second method (called “seed”) is based on the assumption that a representative set of resources can be obtained by sampling across classes in the dataset. Let x be the desired scale factor in percent, e.g. $x = 10$. The method first selects $x\%$ of the classes in the dataset. For each selected class, 10% of its instances are retrieved and added to a queue. For each element of the queue, its concise bound description (CBD) [19] is retrieved. This can lead to new resources, which are appended at the end of the queue. This process is iterated until the target dataset size, measured in number of triples, is reached.

Since the selection of the appropriate method for generating small datasets is an important issue, we performed a statistical analysis on the generated datasets for DBpedia. The statistical parameters used to judge the datasets are the average indegree, the average outdegree, and the number of nodes, i.e. number of distinct IRIs in the graph. We calculated both the in- and the outdegree for datasets once with literals ignored, and another time with literals taken into consideration, as it gives more insight on the degree of similarity between the dataset of interest and the full DBpedia dataset. The statistics of those datasets are given in Table 1. According to this analysis, the seed method fits our purpose of maintaining basic network characteristics better, as the average in- and outdegree of nodes are closer to the original dataset. For this reason, we selected this method for generating the DBPSB.

3 Query Analysis and Clustering

The goal of the query analysis and clustering is to detect prototypical queries that were sent to the official DBpedia SPARQL endpoint based on a query-similarity graph. Note that two types of similarity measures can be used on queries, i.e. string similarities and graph similarities. Yet, since graph similarities are very time-consuming and do not bear the specific mathematical characteristics necessary to compute similarity scores efficiently, we picked string similarities for our experiments. In the query analysis and clustering step, we follow a four-step approach. First, we select queries that were executed frequently on the input data source. Second, we strip common syntactic constructs (e.g., namespace prefix definitions) from these query strings in order to increase the conciseness of the query strings. Then, we compute a query similarity graph

from the stripped queries. Finally, we use a soft graph clustering algorithm for computing clusters on this graph. These clusters are subsequently used to devise the query generation patterns used in the benchmark. In the following, we describe each of the four steps in more detail.

Query Selection For the DBPSB, we use the DBpedia SPARQL query-log which contains all queries posed to the official DBpedia SPARQL endpoint for a three-month period in 2010⁴. For the generation of the current benchmark, we used the log for the period from April to July 2010. Overall, 31.5 million queries were posed to the endpoint within this period. In order to obtain a small number of distinctive queries for benchmarking triple stores, we reduce those queries in the following two ways:

- *Query variations.* Often, the same or slight variations of the same query are posed to the endpoint frequently. A particular cause of this is the renaming of query variables. We solve this issue by renaming all query variables in a consecutive sequence as they appear in the query, i.e., *var0*, *var1*, *var2*, and so on. As a result, distinguishing query constructs such as REGEX or DISTINCT are a higher influence on the clustering.
- *Query frequency.* We discard queries with a low frequency (below 10) because they do not contribute much to the overall query performance.

The application of both methods to the query log data set at hand reduced the number of queries from 31.5 million to just 35,965. This reduction allows our benchmark to capture the essence of the queries posed to DBpedia within the timespan covered by the query log and reduces the runtime of the subsequent steps substantially.

String Stripping Every SPARQL query contains substrings that segment it into different clauses. Although these strings are essential during the evaluation of the query, they are a major source of noise when computing query similarity, as they boost the similarity score without the query patterns being similar per se. Therefore, we remove all SPARQL syntax keywords such as PREFIX, SELECT, FROM and WHERE. In addition, common prefixes (such as <http://www.w3.org/2000/01/rdf-schema#> for RDF-Schema) are removed as they appear in most queries.

Similarity Computation The goal of the third step is to compute the similarity of the stripped queries. Computing the Cartesian product of the queries would lead to a quadratic runtime, i.e., almost 1.3 billion similarity computations. To reduce the runtime of the benchmark compilation, we use the LIMES framework [13]⁵. The LIMES approach makes use of the interchangeability of similarities and distances. It presupposes a metric space in which the queries are expressed as single points. Instead of aiming to find all pairs of queries such that $sim(q, p) \geq \theta$, LIMES aims to find all pairs of queries such that $d(q, p) \leq \tau$, where sim is a similarity measure and d is the corresponding metric. To achieve this goal, when given a set of n queries, it first computes

⁴ The DBpedia SPARQL endpoint is available at: <http://dbpedia.org/sparql/> and the query log excerpt at: <ftp://download.openlinksw.com/support/dbpedia/>.

⁵ Available online at: <http://limes.sf.net>

\sqrt{n} so-called *exemplars*, which are prototypical points in the affine space that subdivide it into regions of high heterogeneity. Then, each query is mapped to the exemplar it is least distant to. The characteristics of metrics spaces (especially the triangle inequality) ensures that the distances from each query q to any other query p obeys the following inequality

$$d(q, e) - d(e, p) \leq d(q, p) \leq d(q, e) + d(e, p), \quad (1)$$

where e is an exemplar and d is a metric. Consequently,

$$d(q, e) - d(e, p) > \tau \Rightarrow d(q, p) > \tau. \quad (2)$$

Given that $d(q, e)$ is constant, q must only be compared to the elements of the list of queries mapped to e that fulfill the inequality above. By these means, the number of similarity computation can be reduced significantly. In this particular use case, we cut down the number of computations to only 16.6% of the Cartesian product without any loss in recall. For the current version of the benchmark, we used the *Levenshtein* string similarity measure and a threshold of 0.9.

Clustering The final step of our approach is to apply graph clustering to the query similarity graph computed above. The goal of this step is to discover very similar groups queries out of which prototypical queries can be generated. As a given query can obey the patterns of more than one prototypical query, we opt for using the soft clustering approach implemented by the BorderFlow algorithm⁶.

BorderFlow [12] implements a seed-based approach to graph clustering. The default setting for the seeds consists of taking all nodes in the input graph as seeds. For each seed v , the algorithm begins with an initial cluster X containing only v . Then, it expands X iteratively by adding nodes from the direct neighborhood of X to X until X is node-maximal with respect to a function called the border flow ratio. The same procedure is repeated over all seeds. As different seeds can lead to the same cluster, identical clusters (i.e., clusters containing exactly the same nodes) that resulted from different seeds are subsequently collapsed to one cluster. The set of collapsed clusters and the mapping between each cluster and its seeds are returned as result. Applying BorderFlow to the input queries led to 12272 clusters, of which 24% contained only one node, hinting towards a long-tail distribution of query types. To generate the patterns used in the benchmark, we only considered clusters of size 5 and above.

4 SPARQL Feature Selection and Query Variability

After the completion of the detection of similar queries and their clustering, our aim is now to select a number of frequently executed queries that cover most SPARQL features and allow us to assess the performance of queries with single as well as combinations of features. The SPARQL features we consider are:

- the overall number of triple patterns contained in the query (*GP*),
- the graph pattern constructors UNION (*UON*), OPTIONAL (*OPT*),

⁶ An implementation of the algorithm can be found at <http://borderflow.sf.net>

```

1 SELECT * WHERE {
2   { ?v2 a dbp-owl:Settlement ;
3     rdfs:label %%v%% .
4     ?v6 a dbp-owl:Airport . }
5   { ?v6 dbp-owl:city ?v2 . }
6   UNION
7   { ?v6 dbp-owl:location ?v2 . }
8   { ?v6 dbp-prop:iata ?v5 . }
9   UNION
10  { ?v6 dbp-owl:iataLocationIdentifier ?v5 . }
11  OPTIONAL { ?v6 foaf:homepage ?v7 . }
12  OPTIONAL { ?v6 dbp-prop:nativename ?v8 . }
13 }

```

Fig. 1: Sample query with placeholder.

- the solution sequences and modifiers *DISTINCT (DST)*,
- as well as the filter conditions and operators *FILTER (FLT)*, *LANG (LNG)*, *REGEX (REG)* and *STR (STR)*.

We pick different numbers of triple patterns in order to include the efficiency of JOIN operations in triple stores. The other features were selected because they frequently occurred in the query log. We rank the clusters by the sum of the frequency of all queries they contain. Thereafter, we select 25 queries as follows: For each of the features, we choose the highest ranked cluster containing queries having this feature. From that particular cluster we select the query with the highest frequency.

In order to convert the selected queries into query templates, we manually select a part of the query to be varied. This is usually an IRI, a literal or a filter condition. In Figure 1 those varying parts are indicated by `%%v%%` or in the case of multiple varying parts `%%vn%%`. We exemplify our approach to replacing varying parts of queries by using Query 9, which results in the query shown in Figure 1. This query selects a specific settlement along with the airport belonging to that settlement as indicated in Figure 1. The variability of this query template was determined by getting a list of all settlements using the query shown in Figure 2. By selecting suitable placeholders, we ensured that the variability is sufficiently high (≥ 1000 per query template). Note that the triple store used for computing the variability was different from the triple store that we later benchmarked in order to avoid potential caching effects.

For the benchmarking we then used the list of thus retrieved concrete values to replace the `%%v%%` placeholders within the query template. This method ensures, that (a) the actually executed queries during the benchmarking differ, but (b) always return results. This change imposed on the original query avoids the effect of simple caching.

5 Experimental Setup

This section presents the setup we used when applying the DBPSB on four triple stores commonly used in Data Web applications. We first describe the triple stores and their configuration, followed by our experimental strategy and finally the obtained results. All experiments were conducted on a typical server machine with an AMD Opteron

```

1 SELECT DISTINCT ?v WHERE {
2   { ?v2 a dbp-owl:Settlement ;
3     rdfs:label ?v .
4     ?v6 a dbp-owl:Airport . }
5   { ?v6 dbp-owl:city ?v2 . }
6   UNION
7   { ?v6 dbp-owl:location ?v2 . }
8   { ?v6 dbp-prop:iata ?v5 . }
9   UNION
10  { ?v6 dbp-owl:iataLocationIdentifier ?v5 . }
11  OPTIONAL { ?v6 foaf:homepage ?v7 . }
12  OPTIONAL { ?v6 dbp-prop:nativename ?v8 . }
13 } LIMIT 1000

```

Fig. 2: Sample auxiliary query returning potential values a placeholder can assume.

6 Core CPU with 2.8 GHz, 32 GB RAM, 3 TB RAID-5 HDD running Linux Kernel 2.6.35-23-server and Java 1.6 installed. The benchmark program and the triple store were run on the same machine to avoid network latency.

Triple Stores Setup We carried out our experiments by using the triple stores *Virtuoso* [7], *Sesame* [5], *Jena-TDB* [15], and *BigOWLIM* [3]. The configuration and the version of each triple store were as follows:

1. **Virtuoso** Open-Source Edition version 6.1.2: We set the following memory-related parameters: `NumberOfBuffers = 1048576`, `MaxDirtyBuffers = 786432`.
2. **Sesame** Version 2.3.2 with Tomcat 6.0 as HTTP interface: We used the native storage layout and set the `spoc`, `posc`, `opsc` indices in the native storage configuration. We set the Java heap size to 8GB.
3. **Jena-TDB** Version 0.8.7 with Joseki 3.4.3 as HTTP interface: We configured the TDB optimizer to use statistics. This mode is most commonly employed for the TDB optimizer, whereas the other modes are mainly used for investigating the optimizer strategy. We also set the Java heap size to 8GB.
4. **BigOWLIM** Version 3.4, with Tomcat 6.0 as HTTP interface: We set the entity index size to 45,000,000 and enabled the predicate list. The rule set was empty. We set the Java heap size to 8GB.

In summary, we configured all triple stores to use 8GB of memory and used default values otherwise. This strategy aims on the one hand at benchmarking each triple store in a real context, as in real environment a triple store cannot dispose of the whole memory up. On the other hand it ensures that the whole dataset cannot fit into memory, in order to avoid caching.

Benchmark Execution Once the triple stores loaded the DBpedia datasets with different scale factors, i.e. 10%, 50%, 100%, and 200%, the benchmark execution phase began. It comprised the following stages:

1. **System Restart:** Before running the experiment, the triple store and its associated programs were restarted in order to clear memory caches.

2. **Warm-up Phase:** In order to measure the performance of a triple store under normal operational conditions, a warm-up phase was used. In the warm-up phase, query mixes were posed to the triple store. The queries posed during the warm-up phase were disjoint with the queries posed in the hot-run phase. For DBPSB, we used a warm-up period of 20 minutes.
3. **Hot-run Phase:** During this phase, the benchmark query mixes were sent to the tested store. We kept track of the average execution time of each query as well as the number of query mixes per hour (QMpH). The duration of the hot-run phase in DBPSB was 60 minutes.

Since some benchmark queries did not respond within reasonable time, we specified a 180 second timeout after which a query was aborted and the 180 second maximum query time was used as the runtime for the given query even though no results were returned. The benchmarking code along with the DBPSB queries is freely available⁷.

6 Results

We evaluated the performance of the triple stores with respect to two main metrics: their overall performance on the benchmark and their query-based performance.

The overall performance of any triple store was measured by computing its query mixes per hour (QMpH) as shown in Figure 4. Please note that we used a logarithmic scale in this figure due to the high performance differences we observed. In general, Virtuoso was clearly the fastest triple store, followed by BigOWLIM, Sesame and Jena-TDB. The highest observed ratio in QMpH between the fastest and slowest triple store was 63.5 and it reached more than 10 000 for single queries. The scalability of stores did not vary as much as the overall performance. There was on average a linear decline in query performance with increasing dataset size. Details will be discussed in Section 7.

We tested the queries that each triple store failed to executed withing the 180s timeout and noticed that even much larger timeouts would not have been sufficient most of those queries. We did not exclude the queries completely from the overall assessment, since this would have affected a large number of the queries and adversely penalized stores, which complete queries within the time frame. We penalized failure queries with 180s, similar to what was done in the SP2-Benchmark [18]. Virtuoso was the only store, which completed all queries in time. For Sesame and OWLIM only rarely a few particular queries timed out. Jena-TDB had always severe problems with queries 7, 10 and 20 as well as 3, 9, 12 for the larger two datasets.

The metric used for query-based performance evaluation is Queries per Second (QpS). QpS is computed by summing up the runtime of each query in each iteration, dividing it by the QMpH value and scaling it to seconds. The QpS results for all triple stores and for the 10%, 50%, 100%, and 200% datasets are depicted in Figure 3.

The outliers, i.e. queries with very low QpS, will significantly affect the mean value of QpS for each store. So, we additionally calculated the geometric mean of all the QpS timings of queries for each store. The geometric mean for all triple stores is also depicted in Figure 4. By reducing the effect of outliers, we obtained additional information from this figure as we will describe in the subsequent section.

⁷ <https://akswbenchmark.svn.sourceforge.net/svnroot/akswbenchmark/>

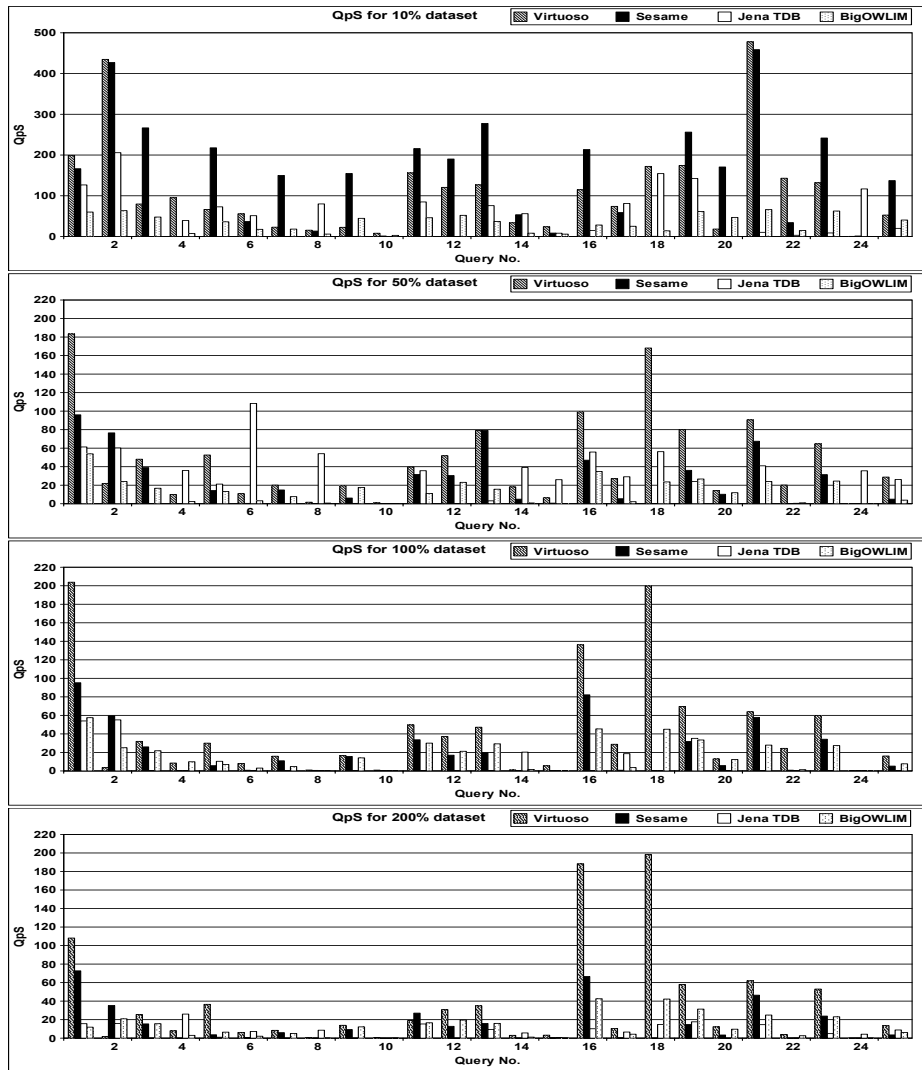


Fig. 3: Queries per Second (QpS) for all triple stores for 10%, 50%, 100%, and 200% datasets.

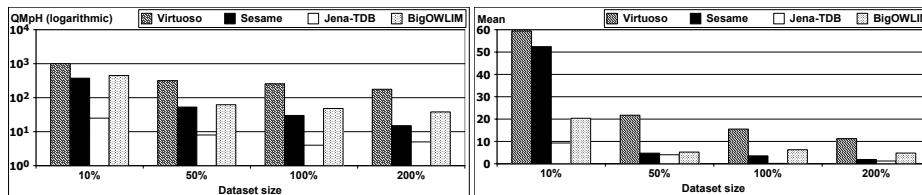


Fig. 4: QMpH for all triple stores (left). Geometric mean of QpS (right).

7 Discussion

This section consists of three parts: First, we compare the general performance of the systems under test. Then we look individual queries and the SPARQL features used within those queries in more detail to observe particular strengths and weaknesses of stores. Thereafter, we compare our results with those obtained with previous benchmarks and elucidate some of the main differences between them.

General Performance Figure 4 depicts the benchmark results for query mixes per hour for the four systems and dataset sizes. Virtuoso leads the field with a substantial head start of double the performance for the 10% dataset (and even quadruple for other dataset sizes) compared to the second best system (BigOWLIM). While Sesame is able to keep up with BigOWLIM for the smaller two datasets it considerably loses ground for the larger datasets. Jena-TDB can in general not deliver competitive performance with being by a factor 30-50 slower than the fastest system.

If we look at the geometric mean of all QpS results in Figure 4, we observe similar insights. The spreading effect is weakened, since the geometric mean reduces the effect of outliers. Still Virtuoso is the fastest system, although Sesame manages to get pretty close for the 10% dataset. This shows that most, but not all, queries are fast in Sesame for low dataset sizes. For the larger datasets, BigOWLIM is the second best system and shows promising scalability, but it is still by a factor of two slower than Virtuoso.

Scalability, Individual Queries and SPARQL Features Our first observation with respect to individual performance of the triple stores is that Virtuoso demonstrates a good scaling factor on the DBPSB. When dataset size changes by factor 5 (from 10% to 50%), the performance of the triple store only degrades by factor 3.12. Further dataset increases (i.e. the doubling to the 100% and 200% datasets) result in only relatively small performance decreases by 20% and respectively 30%.

Virtuoso outperforms Sesame for all datasets. In addition, Sesame does not scale as well as Virtuoso for small dataset sizes, as its performance degrades sevenfold when the dataset size changes from 10% to 50%. However, when the dataset size doubles from the 50% to the 100% dataset and from 100% to 200% the performance degrades by just half.

The performance of Jena-TDB is the lowest of all triple stores and for all dataset sizes. The performance degradation factor of Jena-TDB is not as pronounced as that of Sesame and almost equal to that of Virtuoso when changing from the 10% to the 50% dataset. However, the performance of Jena-TDB only degrades by a factor of 2 for the transition between the 50% and 100% dataset, and reaches 0.8 between the 100% and 200% dataset, leading to a slight increase of its QMPH.

BigOWLIM is the second fastest triple store for all dataset sizes, after Virtuoso. BigOWLIM degrades with a factor of 7.2 in transition from 10% to 50% datasets, but it decreases dramatically to 1.29 with dataset size 100%, and eventually reaches 1.26 with dataset size 200%.

Due to the high diversity in the performance of different SPARQL queries, we also computed the geometric mean of the QpS values of all queries as described in the previous section and illustrated in Figure 4. By using the geometric mean, the resulting

values are less prone to be dominated by a few outliers (slow queries) compared to standard QMPH values. This allows for some interesting observations in DBPSB by comparing Figure 4 and 4. For instance, it is evident that Virtuoso has the best QpS values for all dataset sizes.

With respect to Virtuoso, query 10 performs quite poorly. This query involves the features FILTER, DISTINCT, as well as OPTIONAL. Also, the well performing query 1 involves the DISTINCT feature. Query 3 involves a OPTIONAL resulting in worse performance. Query 2 involving a FILTER condition results in the worst performance of all of them. This indicates that using complex FILTER in conjunction with additional OPTIONAL, and DISTINCT adversely affects the overall runtime of the query.

Regarding Sesame, queries 4 and 18 are the slowest queries. Query 4 includes UNION along with several free variables, which indicates that using UNION with several free variables causes problems for Sesame. Query 18 involves the features UNION, FILTER, STR and LANG. Query 15 involves the features UNION, FILTER, and LANG, and its performance is also pretty slow, which leads to the conclusion that introducing this combination of features is difficult for Sesame. Adding the STR feature to that feature combination affects the performance dramatically and prevents the query from being successfully executed.

For Jena-TDB, there are several queries that timeout with large dataset sizes, but queries 10 and 20 always timeout. The problem with query 10 is already discussed with Virtuoso. Query 20 contains FILTER, OPTIONAL, UNION, and LANG. Query 2 contains FILTER only, query 3 contains OPTIONAL, and query 4 contains UNION only. All of those queries run smoothly with Jena-TDB, which indicates that using the LANG feature, along with those features affects the runtime dramatically.

For BigOWLIM, queries 10, and 15 are slow queries. Query 10 was already problematic for Virtuoso, as was query 15 for Sesame.

Query 24 is slow on Virtuoso, Sesame, and BigOWLIM, whereas it is faster on Jena-TDB. This is due to the fact that most of the time this query returns many results. Virtuoso, and BigOWLIM return a bulk of results at once, which takes long time. Jena-TDB just returns the first result as a starting point, and iteratively returns the remaining results via a buffer.

It is interesting to note that BigOWLIM shows in general good performance, but almost never manages to outperform any of the other stores. Queries 11, 13, 19, 21 and 25 were performed with relatively similar results across triple stores thus indicating that the features of these queries (i.e. UON, REG, FLT) are already relatively well supported. With queries 3, 4, 7, 9, 12, 18, 20 we observed dramatic differences between the different implementations with factors between slowest and fastest store being higher than 1000. It seems that a reason for this could be the poor support for OPT (in queries 3, 7, 9, 20) as well as certain filter conditions such as LNG in some implementations, which demonstrates the need for further optimizations.

Comparison with Previous Benchmarks In order to visualize the performance improvement or degradation of a certain triple store compared to its competitors, we calculated the relative performance for each store compared to the average and depicted it for each dataset size in Figure 5. We also performed this calculation for BSBM version 2 and version 3. Overall, the benchmarking results with DBPSB were less homogeneous than the

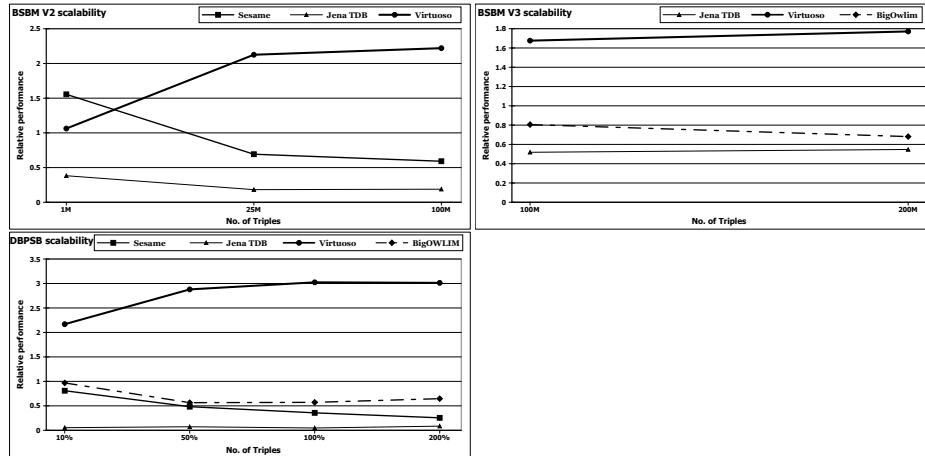


Fig. 5: Comparison of triple store scalability between BSBM V2, BSBM V3, DBPSB.

results of previous benchmarks. While with other benchmarks the ratio between fastest and slowest query rarely exceeds a factor of 50, the factor for the DBPSB queries (derived from real DBpedia SPARQL endpoint queries) reaches more than 1 000 in some cases.

As with the other benchmarks, Virtuoso was also fastest in our measurements. However, the performance difference is even higher than reported previously: Virtuoso reaches a factor of 3 in our benchmark compared to 1.8 in BSBM V3. BSBM V2 and our benchmark both show that Sesame is more suited to smaller datasets and does not scale as well as other stores. Jena-TDB is the slowest store in BSBM V3 and DBPSB, but in our case they fall much further behind to the point that Jena-TDB can hardly be used for some of the queries, which are asked to DBpedia. The main observation in our benchmark is that previously observed differences in performance between different triple stores amplify when they are confronted with actually asked SPARQL queries, i.e. there is now a wider gap in performance compared to essentially relational benchmarks.

8 Related work

Several RDF benchmarks were previously developed. The *Lehigh University Benchmark* (LUBM) [16] was one of the first RDF benchmarks. LUBM uses an artificial data generator, which generates synthetic data for universities, their departments, their professors, employees, courses and publications. This small number of classes limits the variability of data and makes LUBM inherent structure more repetitive. Moreover, the SPARQL queries used for benchmarking in LUBM are all plain queries, i.e. they contain only triple patterns with no other SPARQL features (e.g. FILTER, or REGEX). LUBM performs each query 10 consecutive times, and then it calculates the average re-

	LUBM	SP²Bench	BSBM V2	BSBM V3	DBPSB
RDF stores tested	DLDB-OWL, Sesame, OWL-JessKB	ARQ, Redland, SDB, Sesame, Virtuoso	Virtuoso, Sesame, Jena-TDB, Jena-SDB	Virtuoso, 4store, BigData, Jena-TDB, BigOwlim	Virtuoso, Jena-TDB, BigOWLIM, Sesame
Test data	Synthetic	Synthetic	Synthetic	Synthetic	Real
Test queries	Synthetic	Synthetic	Synthetic	Synthetic	Real
Size of tested datasets	0.1M, 0.6M, 1.3M, 2.8M, 6.9M	10k, 50k, 250k, 1M,	1M, 25M, 100M, 5M, 25M	100M, 200M	14M, 75M, 150M, 300M
Dist. queries	14	12	12	12	25
Multi-client	–	–	x	x	–
Use case	Universities	DBLP	E-commerce	E-commerce	DBpedia
Classes	43	8	8	8	239 (internal) +300K(YAGO)
Properties	32	22	51	51	1200

Table 2: Comparison of different RDF benchmarks.

sponse time of that query. Executing the same query several times without introducing any variation enables query caching, which affects the overall average query times.

SP²Bench [18] is another more recent benchmark for RDF stores. Its RDF data is based on the Digital Bibliography & Library Project (DBLP) and includes information about publications and their authors. It uses the SP²Bench Generator to generate its synthetic test data, which is in its schema heterogeneity even more limited than LUBM. The main advantage of SP²Bench over LUBM is that its test queries include a variety of SPARQL features (such as FILTER, and OPTIONAL). The main difference between the DBpedia benchmark and SP²Bench is that both test data and queries are synthetic in SP²Bench. In addition, SP²Bench only published results for up to 25M triples, which is relatively small with regard to datasets such as DBpedia and LinkedGeoData.

Another benchmark described in [14] compares the performance of BigOWLIM and AllegroGraph. The size of its underlying synthetic dataset is 235 million triples, which is sufficiently large. The benchmark measures the performance of a variety of SPARQL constructs for both stores when running in single and in multi-threaded modes. It also measures the performance of adding data, both using bulk-adding and partitioned-adding. The downside of that benchmark is that it compares the performance of only two triple stores. Also the performance of each triple store is not assessed for different dataset sizes, which prevents scalability comparisons.

The Berlin SPARQL Benchmark (BSBM) [4] is a benchmark for RDF stores, which is applied to various triple stores, such as Sesame, Virtuoso, and Jena-TDB. It is based on an e-commerce use case in which a set of products is provided by a set of vendors and consumers post reviews regarding those products. It tests various SPARQL features on those triple stores. It tries to mimic a real user operation, i.e. it orders the query in a manner that resembles a real sequence of operations performed by a human user. This is an effective testing strategy. However, BSBM data and queries are artificial and the data schema is very homogeneous and resembles a relational database. This is reasonable

for comparing the performance of triple stores with RDBMS, but does not give many insights regarding the specifics of RDF data management.

A comparison between benchmarks is shown in Table 2. In addition to general purpose RDF benchmarks it is reasonable to develop benchmarks for specific RDF data management aspects. One particular important feature in practical RDF triple store usage scenarios (as was also confirmed by DBPSB) is full-text search on RDF literals. In [11] the LUBM benchmark is extended with synthetic scalable fulltext data and corresponding queries for fulltext-related query performance evaluation. RDF stores are benchmarked for basic fulltext queries (classic IR queries) as well as hybrid queries (structured and fulltext queries).

9 Conclusions and Future Work

We proposed the DBPSB benchmark for evaluating the performance of triple stores based on non-artificial data and queries. Our solution was implemented for the DBpedia dataset and tested with 4 different triple stores, namely Virtuoso, Sesame, Jena-TDB, and BigOWLIM. The main advantage of our benchmark over previous work is that it uses real RDF data with typical graph characteristics including a large and heterogeneous schema part. Furthermore, by basing the benchmark on queries asked to DBpedia, we intend to spur innovation in triple store performance optimisation towards scenarios, which are actually important for end users and applications. We applied query analysis and clustering techniques to obtain a diverse set of queries corresponding to feature combinations of SPARQL queries. Query variability was introduced to render simple caching techniques of triple stores ineffective.

The benchmarking results we obtained reveal that real-world usage scenarios can have substantially different characteristics than the scenarios assumed by prior RDF benchmarks. Our results are more diverse and indicate less homogeneity than what is suggested by other benchmarks. The creativity and inaptness of real users while constructing SPARQL queries is reflected by DBPSB and unveils for a certain triple store and dataset size the most costly SPARQL feature combinations.

Several improvements can be envisioned in future work to cover a wider spectrum of features in DBPSB: Several improvements can be envisioned in future work to cover a wider spectrum of features in DBPSB:

- Coverage of more SPARQL 1.1 features, e.g. reasoning and subqueries.
- Inclusion of further triple stores and continuous usage of the most recent DBpedia query logs.
- Testing of SPARQL update performance via DBpedia Live, which is modified several thousand times each day. In particular, an analysis of the dependency of query performance on the dataset update rate could be performed.
- We may also further investigate the data generation process in future work, in particular based on recent work in [6].

References

1. Sören Auer, Jens Lehmann, and Sebastian Hellmann. LinkedGeoData - adding a spatial dimension to the web of data. In *Proc. of 8th International Semantic Web Conference (ISWC)*, 2009.
2. François Belleau, Marc-Alexandre Nolin, Nicole Tourigny, Philippe Rigault, and Jean Morissette. Bio2rdf: Towards a mashup to build bioinformatics knowledge systems. *Journal of Biomedical Informatics*, 41(5):706–716, 2008.
3. Barry Bishop, Atanas Kiryakov, Damyan Ognyanoff, Ivan Peikov, Zdravko Tashev, and Ruslan Velkov. Owlrim: A family of scalable semantic repositories. *Semantic Web*, 2(1):1–10, 2011.
4. Christian Bizer and Andreas Schultz. The Berlin SPARQL Benchmark. *Int. J. Semantic Web Inf. Syst.*, 5(2):1–24, 2009.
5. Jeen Broekstra, Arjohn Kampman, and Frank van Harmelen. Sesame: A generic architecture for storing and querying RDF and RDF schema. In *ISWC*, number 2342 in LNCS, pages 54–68. Springer, July 2002.
6. Songyun Duan, Anastasios Kementsietsidis, Kavitha Srinivas, and Octavian Udrea. Apples and oranges: a comparison of RDF benchmarks and real RDF datasets. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 145–156. ACM, 2011.
7. Orri Erling and Ivan Mikhailov. RDF support in the virtuoso DBMS. In Sören Auer, Christian Bizer, Claudia Müller, and Anna V. Zhdanova, editors, *CSSW*, volume 113 of *LNI*, pages 59–68. GI, 2007.
8. Jim Gray, editor. *The Benchmark Handbook for Database and Transaction Systems (1st Edition)*. Morgan Kaufmann, 1991.
9. Graham Klyne and Jeremy J. Carroll. Resource description framework (RDF): Concepts and abstract syntax. W3C Recommendation, February 2004.
10. Jens Lehmann, Chris Bizer, Georgi Kobilarov, Sören Auer, Christian Becker, Richard Cyganiak, and Sebastian Hellmann. DBpedia - a crystallization point for the web of data. *Journal of Web Semantics*, 7(3):154–165, 2009.
11. Enrico Minack, Wolf Siberski, and Wolfgang Nejdl. Benchmarking fulltext search performance of RDF stores. In *ESWC2009*, pages 81–95, June 2009.
12. Axel-Cyrille Ngonga Ngomo and Frank Schumacher. Borderflow: A local graph clustering algorithm for natural language processing. In *CICLing*, pages 547–558, 2009.
13. Axel-Cyrille Ngonga Ngomo and Sören Auer. Limes - a time-efficient approach for large-scale link discovery on the web of data. In *Proceedings of IJCAI*, 2011.
14. Alisdair Owens, Nick Gibbins, and mc schraefel. Effective benchmarking for rdf stores using synthetic data, May 2008.
15. Alisdair Owens, Andy Seaborne, Nick Gibbins, and mc schraefel. Clustered TDB: A clustered triple store for jena. Technical report, Electronics and Computer Science, University of Southampton, 2008.
16. Zhengxiang Pan, Yuanbo Guo, , and Jeff Heflin. LUBM: A benchmark for OWL knowledge base systems. In *Journal of Web Semantics*, volume 3, pages 158–182, 2005.
17. Eric Prud'hommeaux and Andy Seaborne. SPARQL Query Language for RDF. W3C Recommendation, 2008.
18. Michael Schmidt, Thomas Hornung, Georg Lausen, and Christoph Pinkel. SP2Bench: A SPARQL performance benchmark. In *ICDE*, pages 222–233. IEEE, 2009.
19. Patrick Stickler. CBD - concise bounded description, 2005. Retrieved February 15, 2011, from <http://www.w3.org/Submission/CBD/>.