

# Compressed Indexes for Fast Search of Semantic Data

(Extended Abstract)

Raffaele Perego  
ISTI-CNR, Pisa, Italy

Giulio Ermanno Pibiri  
ISTI-CNR, Pisa, Italy

Rossano Venturini  
University of Pisa, Pisa, Italy

**Abstract**—The sheer increase in volume of RDF data demands efficient solutions for the *triple indexing problem*, that is devising a compressed data structure to compactly represent RDF triples by guaranteeing, at the same time, fast pattern matching operations. This problem lies at the heart of delivering good practical performance for the resolution of complex SPARQL queries on large RDF datasets. We propose a trie-based index layout to solve the problem and introduce two novel techniques to reduce its space of representation for improved effectiveness. The extensive experimental analysis reveals that our best space/time trade-off configuration substantially outperforms existing solutions at the state-of-the-art, by taking 30–60% less space and speeding up query execution by a factor of 2–81 times.

## I. THE PERMUTED TRIE INDEX

The Resource Description Framework (RDF) is a W3C standard offering a general graph-based model for describing information as a set of (*subject, predicate, object*) relations, known as triples. Each one of these components is a URI string (or even a literal in the case of an object). Since URI strings can be very long and the same URI generally appears in many RDF statements, the components of triples are commonly mapped to an integer ID to save space, so that each triple in the dataset can be represented with three integers. In this work, we focus on reducing the space of representation of the integer triples while attaining to efficient resolution of all possible selection patterns, noting that this is crucial for guaranteeing practical SPARQL query evaluation [Perego et al., 2020].

As a high-level overview, our solution maintains three different permutations of the triples, with each permutation sorted to allow efficient searches and effective compression as we are going to detail next. The permutations chosen are SPO, POS and OSP in order to (symmetrically) support all the six different triple selection patterns with one or two wildcard symbols: SP? and S?? over SPO; ?PO and ?P? over POS; S?O and ??O over OSP. The two additional patterns with, respectively, all symbols specified or none, can be resolved over any permutation, e.g., over the canonical SPO in order to avoid permuting back each returned triple.

Each permutation of the triples is represented as a trie with 3 levels, where nodes at the same level concatenated together to form an integer sequence. We keep track of where groups of siblings begin and end in the concatenated sequence of nodes by storing such pointers as absolute positions in the sequence of nodes. Therefore, the pointers are integer sequences as

well. The advantage of this layout is *two-fold*: first, we can effectively compress the integer sequences that constitute the levels of the tries to achieve small storage requirements; second, the triple selection patterns are made *cache-friendly* and, hence, efficient by requiring to simply *scan* ranges of consecutive nodes in the trie levels. In the following, we refer to this solution as the 3T index.

**Cross Compression.** It is possible to employ levels of the tries to compress levels of other tries, thus holistically *cross-compressing* the different permutations. Cross-compression works by noting this crucial property: *the nodes belonging to a sub-tree rooted in the second level of trie  $j$  are a subset of the nodes belonging to a sub-tree rooted in the first level of trie  $i$ , with  $j = (i + 2) \bmod 3$ , for  $i = 0, 1, 2$* . The correctness of this property follows automatically by taking into account that the triples indexed by each permutation are the same. Therefore, the children of  $x$  in the second level of trie  $j$  can be re-written as the *positions* they take in the (larger, enclosing) set of children of  $x$  in the first level of trie  $i$ . Re-writing the node IDs as positions relative to the set of children of a sub-tree yields a clear space optimization because the number of children of a given node is much smaller (on average) than the number of distinct subjects or objects (see also Table II for the precise statistics).

**Eliminating a Permutation.** The low number of predicates exhibited by RDF data leads us to consider a different **select** algorithm for the resolution of the query pattern S?O, able to take advantage of such skewed distribution. The idea is to pattern match S?O directly over the SPO permutation. For a given subject  $s$  and object  $o$ , in short, we operate as follows. We consider the set of all the predicates that are children of  $s$ . For each predicate  $p$  in the set, we determine if the object  $o$  is a child of  $p$  with a search operation: if it is, then  $(s, p, o)$  is a triple to return. In the light of this algorithm, we consider another index layout. In fact, now five of out the eight different selection patterns can be solved efficiently by the trie SPO, i.e.: SPO, SP?, S??, S?O and ????. In order to support other two selection patterns, we can either chose to: (1) materialize the permutation POS for *predicate-based* retrieval (query patterns ?PO and ?P?); (2) materialize the permutation OPS for *object-based* retrieval (query patterns ?PO and ??O). The choice of which permutation to maintain depends on the statistics of the selection patterns that have to be supported.

Index	DBLP	Geonames	DBpedia	Freebase	
	bits/triple	bits/triple	bits/triple	bits/triple	
3T	75.24(+31%)	71.59 (+32%)	80.64(+33%)	74.20(+30%)	
CC	63.54(+18%)	67.04 (+27%)	66.91(+19%)	70.46(+26%)	
2To	56.46 (+8%)	53.23 (+8%)	57.51 (+6%)	55.72 (+6%)	
2Tp	<b>51.99</b>	<b>48.98</b>	<b>54.14</b>	<b>52.17</b>	
	ns/triple	ns/triple	ns/triple	ns/triple	
SPO <i>all</i>	203	221	353	521	
SP? <i>all</i>	197	347	11	3	
S?? <i>all</i>	28	40	10	3	
??? <i>all</i>	11	13	9	9	
S?O	3T,CC 2To,2Tp	2490 (5.6×) <b>445</b>	3767 (7.7×) <b>490</b>	1833 (2.6×) <b>692</b>	6547 (1.8×) <b>3736</b>
?PO	3T,2To,2Tp CC	<b>5</b> 12 (2.4×)	<b>5</b> 15 (3.0×)	<b>5</b> 16 (3.2×)	<b>5</b> 14 (2.8×)
??O	3T,CC 2To 2Tp	12 (2.4×) <b>5</b> 5 (1.0×)	12 (2.4×) <b>5</b> 5 (1.0×)	12 (2.4×) <b>5</b> 6 (1.2×)	10 (2.0×) <b>5</b> 10 (2.0×)
?P?	3T,2Tp CC 2To	<b>9</b> 21 (2.3×) 81 (9.0×)	<b>8</b> 36 (4.5×) 138(17.2×)	<b>6</b> 30 (5.0×) 22 (3.7×)	<b>6</b> 29 (4.8×) 18 (3.0×)

TABLE I: Comparison between the performance of 3T, CC and 2T indexes, expressed as the total space in bits/triple and in average ns/triple for all the different selection patterns.

Dataset	Triples	Subjects (S)	Predicates (P)	Objects (O)
DBLP	88,150,324	5,125,936	27	36,413,780
Geonames	123,020,821	8,345,450	26	42,728,317
DBpedia	351,592,624	27,318,781	1,480	115,872,941
Freebase	2,067,068,154	102,001,451	770,415	438,832,462

TABLE II: Datasets statistics.

We stress that the introduced algorithm allows us to actually *save the space* for a third permutation that costs roughly 1/3 of the whole space of the index. We call this solution the 2T index; with two concrete instantiations 2Tp (predicate-based) and 2To (object-based).

## II. EXPERIMENTS

We perform our experimental analysis on large and publicly available standard datasets, whose statistics are summarized in Table II. All the experiments are performed on a server machine with 4 Intel i7-7700 cores (@3.6 GHz), 64 GB of RAM DDR3 (@2.133 GHz) and running Linux 4.4.0, 64 bits. The C++14 implementation of our indexes is freely available at [https://github.com/jermp/rdf\\_indexes](https://github.com/jermp/rdf_indexes). We compiled the code with gcc 7.3.0 using the highest optimization setting. To measure the query processing speed, we use a set of 5000 triples drawn at random from the datasets and set 0, 1 or 2 wildcard symbols. In all tables, percentages and speed up factors are taken with respect to the values in bold font.

**The Permuted Trie Index.** By looking at the results reported in Table I, we can conclude that: (1) the 3T index is the one delivering best *worst-case* performance guarantee for all triple selection patterns; (2) the 2T variants reduce its space of representation by 25–33% *without* affecting or even improving the retrieval efficiency on most triple selection patterns (*only one* out of the eight possible has a lower query throughput in the worst-case); (3) the cross-compression technique is outperformed by the 2T index layouts for space usage but

Index	DBLP	Geonames	DBpedia	Freebase	
	bits/triple	bits/triple	bits/triple	bits/triple	
2Tp	<b>51.99</b>	<b>48.98</b>	<b>54.14</b>	<b>52.17</b>	
HDT-FoQ	76.89 (+32%)	88.73 (+45%)	76.66 (+29%)	83.11 (+37%)	
TripleBit	125.10 (+58%)	120.03 (+59%)	130.07 (+58%)	—	
	ns/triple	ns/triple	ns/triple	ns/triple	
2Tp	5	5	5	5	
?PO HDT-FoQ	12 (2.4×)	13 (2.6×)	14 (2.8×)	13 (2.6×)	
TripleBit	15 (3.0×)	13 (2.6×)	14 (2.8×)	—	
S?O	2Tp HDT-FoQ TripleBit	<b>445</b> 1789 (4.0×) 11872(26.7×)	<b>490</b> 2097 (4.3×) 13008(26.5×)	<b>692</b> 3010 (4.3×) 18023(26.0×)	<b>3736</b> 0.7×10 <sup>7</sup> (2057×) —
SP?	2Tp HDT-FoQ TripleBit	<b>197</b> 640 (3.2×) 1222 (6.2×)	<b>347</b> 897 (2.6×) 927 (2.7×)	<b>11</b> 30 (2.7×) 42 (3.8×)	<b>3</b> 9 (3.0×) —
S??	2Tp HDT-FoQ TripleBit	<b>28</b> 110 (3.9×) 2275(81.2×)	<b>40</b> 154 (3.9×) 3261(81.5×)	<b>10</b> 29 (2.9×) 490(49.0×)	<b>3</b> 9 (3.0×) —
?P?	2Tp HDT-FoQ TripleBit	<b>9</b> 108(12.0×) 28 (3.1×)	<b>8</b> 173(21.6×) 28 (3.5×)	<b>6</b> 32 (5.3×) 40 (6.7×)	<b>4</b> 41 (6.8×) —
??O	2Tp HDT-FoQ TripleBit	<b>5</b> 17 (3.4×) 24 (4.8×)	<b>5</b> 17 (3.4×) 60(12.0×)	<b>6</b> 18 (3.0×) 24 (4.0×)	<b>10</b> 18 (1.8×) —

TABLE III: Comparison between the performance of different indexes, expressed as the total space in bits/triple and in average ns/triple.

offers a better worst-case performance guarantee than 2Tp for the pattern ?P?. Therefore, as a reasonable trade-off between space and time, we elect 2Tp as the solution to compare against the state-of-the-art alternatives.

**Overall Comparison.** We now compare the performance of our selected solution 2Tp against the competitive approaches HDT-FoQ [Martínez-Prieto et al., 2012] and TripleBit [Yuan et al., 2013]. We use the C++ libraries provided by the corresponding authors: <https://github.com/rdfhdt/hdt-cpp>, and <https://github.com/nitingupta910/TripleBit>, respectively. Table III reports the space of the indexes and the timings for the different selection patterns, but excluding (due to page limit) the ones for SPO and ???: our approach is anyway faster for both by at least a factor of 3× (TripleBit does not support the query pattern SPO). Concerning the space, we see that the 2Tp index is significantly more compact, specifically by 30% and almost 60% compared to HDT-FoQ and TripleBit respectively, on average across all different datasets (TripleBit fails in building the index on Freebase). Concerning the speed of triple selection patterns, most factors of improved efficiency range in the interval 2–5× and, depending on the pattern examined, we report peaks of 26×, 49×, 81×, or even 2057×.

## REFERENCES

- M. A. Martínez-Prieto, M. A. Gallego, and J. D. Fernández. Exchange and consumption of huge rdf data. In *Extended Semantic Web Conference*, pages 437–452. Springer, 2012.
- R. Perego, G. E. Pibiri, and R. Venturini. Compressed indexes for fast search of semantic data. *IEEE Transactions on Knowledge and Data Engineering. To appear*, 2020.
- P. Yuan, P. Liu, B. Wu, H. Jin, W. Zhang, and L. Liu. Triplebit: a fast and compact system for large scale rdf data. *PVLDB*, 6(7):517–528, 2013.