

	B.Triples				Uniprot				US Census			
	Dictionary		Triples		Dictionary		Triples		Dictionary		Triples	
	Orig.	Compr.	Orig.	Compr.	Orig.	Compr.	Orig.	Compr.	Orig.	Compr.	Orig.	Compr.
Original	153.2	28.0 (18.3%)	65.5	9.5 (4.2%)	26.3	1.4 (5.3%)	58.6	8.1 (13.8%)	6.8	0.9 (13.2%)	49.6	10.7 (21.6%)
Delta	105.3	28.4 (18.5%)			4.4	1.0 (3.8%)	59.6		5.7	1.3 (19.1%)	52.8	7.7 (15.5%)
Tree	114.0	29.9 (19.5%)			15.9	1.6 (6.1%)						

Table 2: Original and compressed size (in MB) of Dictionary+Triples decomposition.

	B.Triples	Uniprot	US Census
Original Size	541.5	239.4	148.2
bzip2	46.9 (8.7%)	10.8 (4.5%)	9.7 (6.6%)
gzip	55.2 (10.2%)	18.3 (7.7%)	13.4 (9.0%)
ppmdi-6	37.9 (7.0%)	7.5 (3.1%)	6.4 (4.4%)

Table 3: Direct Compression (in MB).

	B.Triples	Uniprot	US Census
Original Size	541.5	239.4	148.2
Subject A.L.	35.5 (6.6%)	5.3 (2.2%)	4.8 (3.3%)
Predicate A.L.	50.7 (9.4%)	6.60 (2.8%)	7.9 (5.3%)
Object A.L.	42.5 (7.9%)	6.6 (2.8%)	7.5 (5.0%)

Table 4: Adjacency Lists Representations (in MB).

lion Triples data set results in weaker compression, while the numeric data nature of U.S. Census is also punished.

b) Adjacency Lists. A second approximation focus data repeatability by using Adjacency Lists (Figure 1b). For example, the set of triples $\{(s, p_1, o_{11}), \dots, (s, p_1, o_{1n_1}), \dots, (s, p_2, o_{21}), \dots, (s, p_2, o_{2n_2}), \dots, (s, p_k, o_{kn_k})\}$ would be written as the adjacency list $s \rightarrow [(p_1, \text{ObjList}_1), \dots, (p_k, \text{ObjList}_k)]$.

We considered three types of adjacency lists: subject (present in Turtle and N3), predicate and object headed. Predicates and subjects precede objects in sublists, and both lists and sublists are ordered lexicographically. Table 4 shows that subject adjacency lists compression (using ppmdi to compare best results) improves direct compression in every case, while predicate and object have better results only in Uniprot. Adjacency lists levels of compression depend on both their compact power and their ability to produce repeated elements in sublists.

c,d) Dictionary+Triples. The next two tests are focused on the graph nature of RDF. Note that standard Web graph compression is hardly applicable to RDF because they exploit locality and similarity features of their link structure[1]. Locality implies that the source and the target of a link tend to share the same domain, therefore lexicographical order benefits the compression. By means of similarity, some successors are shared by pages in the same domain, facilitating compression too. In RDF, these properties are rarely present. Nevertheless, power law assumption and high compression levels previously presented, show a regularity in triples that might be exploited.

A simple graph representation is proposed to test RDF compressibility: we split the data into the dictionary of elements and the triples substituting for each element, the corresponding number assignation in the dictionary. Triples, in turn, can be represented as one type of adjacency lists. Literals in the dictionary were sorted lexicographically and managed independently, as we focused on URIs, exploiting the redundancy of their long shared prefixes. Figure 1c stands for a commonly used delta coding, in which URIs are also sorted lexicographically and each one is coded by two integers and a string. Integers delimit the number of characters

shared with the previous symbol and the number of different characters, while the string represents the portion of the URI that differs from the previous one.

Delta encoding is a compression-oriented representation, so that operations with the dictionary (mainly to find the identifier of an element and vice versa, commonly in RDF triples stores) become more complex. In order to facilitate these operations, we considered a compact tree representation (referred to as DFUDS) built on a sequence of balanced parentheses, as shown in Figure 1d. We stored the preorder traversal sequence of the tree and two bits per node to represent the structure. When using this succinct representation in-memory, an auxiliary structure is needed, with cost $o(n)$, where n is the number of nodes. This structure facilitates common operations without uncompression.

Table 2 details different Dictionary+Triples decomposition. In compression, triples and delta integers are coded by canonical bit-oriented Huffman, while the preorder sequence and the literals are compressed with ppmdi. Compressed values of triples correspond to subject adjacency lists. As a general observation, both delta and tree coding dictionaries reduce original size. In contrast, identifiers re-assignment can slightly increase the size of triples representation. B.T and U.S datasets are composed of a great variety of literals, so that the improvement of both delta and tree coding do not compensate for the uncompressibility of numeric literals in U.S or the variability in B.T. In these cases, they do not reach original compressing levels.

Uniprot is highly compressible due to the massive presence of URIs, which benefits the compression of the dictionary. Most of these URIs are named sequentially, so that tree coding size is bigger than delta before compression, as it only stores the difference and tree repeats the whole identifier.

3. CONCLUSIONS

Table 1 summarizes the results of the different approaches tested. From this study we can conclude:

1. RDF data at big scale is highly compressible.
2. Dedicated data structures, *e.g.* adjacency lists, code triples efficiently and facilitate compression (both string with ppmdi and integer with Huffman).
3. RDF URIs are prone to efficient compression with standard techniques, but compression of literals deserve finer approaches.
4. The structure of RDF graphs differs from XML or Web data, hence, classical approaches such as [1] are not directly applicable.

4. REFERENCES

- [1] P. Boldi and S. Vigna. The webgraph framework i: compression techniques. In *WWW*, pp. 595–602, 2004.
- [2] L. Ding and T. Finin. Characterizing the semantic web on the web. In *ISWC*, pp. 242–257, 2006.
- [3] LOD data sets. <http://linkeddata.org/data-sets>.