



**UNIVERSITÄT PADERBORN**

*Die Universität der Informationsgesellschaft*

Faculty for Computer Science, Electrical Engineering and Mathematics

Department of Computer Science

Research Group Data Science

## Master's Thesis

Submitted to the Data Science Research Group  
in Partial Fulfilment of the Requirements for the Degree of

## Master of Science

# Grammar-based Compression of RDF Graphs

by  
PHILIP FRERK

Thesis Supervisors:

Prof. Dr. Axel-Cyrille Ngonga Ngomo

Prof. Dr. Stefan Böttcher

Paderborn, May 2019



# Erklärung

Ich versichere, dass ich die Arbeit ohne fremde Hilfe und ohne Benutzung anderer als der angegebenen Quellen angefertigt habe und dass die Arbeit in gleicher oder ähnlicher Form noch keiner anderen Prüfungsbehörde vorgelegen hat und von dieser als Teil einer Prüfungsleistung angenommen worden ist. Alle Ausführungen, die wörtlich oder sinngemäß übernommen worden sind, sind als solche gekennzeichnet.

---

Ort, Datum

---

Unterschrift



**Abstract.** We present a full documentation of the Paderborn University Computer Science thesis template (UPB-CS-TT) and how to use it. This document also serves as a demonstrator to show what documents UPB-CS-TT produces. Have fun!



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	General Idea . . . . .	1
1.3	Thesis Structure . . . . .	2
<b>2</b>	<b>Related Work</b>	<b>3</b>
2.1	RDF . . . . .	3
2.2	HDT . . . . .	3
2.3	Grammar-based Graph Compression . . . . .	5
<b>3</b>	<b>Approach</b>	<b>9</b>
3.1	Definitions . . . . .	9
3.2	GRP vs HDT . . . . .	9
3.3	Compression Improvements . . . . .	10
<b>4</b>	<b>Implementation</b>	<b>15</b>
4.1	GRP vs HDT . . . . .	15
4.2	Compression Improvements . . . . .	15
<b>5</b>	<b>Evaluation</b>	<b>19</b>
5.1	Experimental Setup . . . . .	19
5.2	GRP vs HDT . . . . .	19
5.3	Compression Improvements . . . . .	23
<b>6</b>	<b>Summary and Discussion</b>	<b>27</b>
6.1	Summary and Discussion . . . . .	27
6.2	Future Work . . . . .	27
	<b>Bibliography</b>	<b>30</b>





# Introduction

## 1.1 Motivation

The majority of data on the Internet is unstructured because it is mainly text. That makes it difficult for machines to extract knowledge from the data and answer specific queries. In the context of the Semantic Web, an attempt is made to build a knowledge base using structured data. Here the data is stored as graphs. A graph is an often used data structure which consists of nodes and edges. The edges connect the different nodes. This thesis will focus on knowledge graphs, which are about expressing knowledge or facts. A possible format for such graphs is the Resource Description Framework (RDF)<sup>1</sup>, in which a graph is represented by triples of the form (*subject*, *predicate*, *object*), where *subject* and *object* are nodes and *predicate* is an edge of the graph. A triple is typically used to express a certain fact.

In reality, such knowledge graphs can become very large with millions or even billions of nodes and edges. The following three use cases can then become very hard or even infeasible: storage, transmission and processing. One way to circumvent this problem is to use compression.

The currently most popular compressor for RDF data is HDT [FMPG<sup>+</sup>13]. Here the data is made smaller by a compact representation of the triples.

There is a completely different compression technique which is called grammar-based compression. This method has so far been tested very little for RDF graphs. As the name suggests, such a compression is based on the principle of a formal grammar, with productions that can be nested among each other. There are not yet many grammar-based compression algorithms for graphs. One of the best known is GraphRePair [MP18]. This is a compressor that works for any type of graph and is therefore also applicable for RDF. Why the algorithm is particularly suitable for RDF will be explained later on.

In this thesis it shall be investigated to what extent grammar-based compression is applicable for RDF and whether it delivers even better results than HDT.

## 1.2 General Idea

As already mentioned, there are essentially three use cases for RDF data:

1. Storage

---

<sup>1</sup><https://www.w3.org/RDF/>

2. Transmission

3. Processing/Consumption

The idea is that you can make all these use cases easier / possible by compression, especially if one is dealing with very large amounts of data. In the first two use cases, compression can be helpful by reducing the size of the data. Thus the data can be stored more compactly and above all can be transferred faster, which occurs frequently in reality and can become a problem due to possibly slow transfer speeds .

here possibly create numbers

The third use case (Processing/Consumption) is about reading or writing access to data. The more common case of read access is typically the execution of queries on RDF data. For example, you want to find out which other nodes a given node is connected to. But there are other examples which will be explained in the course of the thesis. A compression algorithm can help in this respect by making it possible to answer such queries directly on the compressed graphs. This may even be faster than on the original data due to the reduced size.

The aim of the thesis is to compare the two compressors (HDT and GRP) with regard to these use cases and finally determine which approach works better for which use cases.

## 1.3 Thesis Structure

In Ch. 2 the necessary fundamentals are presented. It will mainly be about the two compressors HDT and GRP.

Ch. 3 deals with the approaches for the different tasks of the thesis. There they will only be presented in a theoretical manner whereas in Ch. 4 the implementation details will be presented. In Ch. 5, the approaches to the different tasks will be evaluated using real RDF data in order to confirm the theoretical hypotheses stated before.

Finally, the results of the thesis are summarized in Ch. 6 and the future work on this topic is presented.

## Related Work

This chapter gives an introduction to the domain RDF and introduces the different compression algorithms discussed in this thesis.

### 2.1 RDF

RDF is a format for structured data on the web. It can be formally described as follows. Let the following sets be infinite and mutually disjoint:

- $U$  (URI references, typically represent unique entities or properties)
- $B$  (Blank nodes, represent an arbitrary value, necessary for displaying more complex logical statements)
- $L$  (Literals, represent fixed values, e.g. numbers or strings)

An RDF triple  $(s, p, o) \in (U \cup B) \times U \times (U \cup B \cup L)$  displays the statement that the subject  $s$  is related to the object  $o$  via the predicate  $p$ . So a subject can anything but a literal. A predicate can only be a URI and an object can be anything.

### 2.2 HDT

HDT ([FMPG<sup>+</sup>13]) is the best known approach to compressing RDF data because it is directly tailored to RDF. The idea behind it is to make RDF's extensive representation more compact, thus reducing memory size. In the normal RDF format (e.g. turtle), each triple is written down after another. This does not only require much space, but is also not good for query performance, since each single triple has to be traversed within a query execution.

It is easy to see that a more compact representation can be achieved by cleverly grouping the triples. That is exactly what HDT does. Imagine for example some triples that all have the same subject  $s$ . These would then be written in turtle format as follows:

$$(s, p_1, o_{11}), \dots, (s, p_1, o_{1n_1}), (s, p_2, o_{21}), \dots, \\ (s, p_2, o_{2n_2}), \dots, (s, p_k, o_{kn_k})$$

In HDT the triples are then grouped by  $s$ .

$$s \rightarrow [(p_1, (o_{1_1}, \dots, o_{1_{n_1}})), (p_2, (o_{2_1}, \dots, o_{2_{n_2}})), \dots, (p_k, (o_{k_{n_k}}))]$$

Since all triples have the same subject  $s$ , it does not have to be written down at the beginning of each triple anymore. On the second level the triples are then grouped according to the predicates. So all triples with the predicate  $p_1$  come first, then all with the predicate  $p_2$  and so on. Finally only the different objects have to be enumerated since the subject and predicate are already implicitly given. You will now see how HDT implements this mechanism.

First, all URIs and literals (which are usually quite long) are mapped to unique IDs (integers). From now on these IDs will be used. You can see the triples with the IDs in Fig. 2.1 on the left ('Plain Triples').

The above mentioned grouped representation is called 'Compact Triples' and can be seen in Fig. 2.1 in the middle. In the array 'Predicates' you can see the IDs of the predicates. A '0' means that from now on a new subject comes. For example, the first two entries in the 'Predicates' are linked to the first subject, since there is a '0' in the third position in the array. This works analogously for the objects, in the array 'Objects' the IDs of the objects are listed, and a '0' here means that from then on a new predicate comes. For instance, the first entry in 'Objects' is a '6', so the first derived triple would be (1, 2, 6). The second entry of 'Objects' is a '0' since there is a change from predicate '2' to '3' in 'Predicates'.

The last representation of triples is called 'Bitmap Triples'. This is a slight adaption of 'Compact Triples'. The array  $S_p$  has the same content as 'Predicates', but without the zeros. That job is done by the bit-array  $B_p$  in which a '1' at position  $i$  means that at position  $i$  in  $S_p$  there comes a change of the subject (this change was previously denoted by a '0' in 'Predicates').

That works analogously for the objects where  $S_o$  has the same content as 'Objects', but without zeros.

The advantage of 'Bitmap Triples' is that queries can be executed faster on this representation. Also, the compressed size is slightly smaller than with 'Compact Triples'. Therefore, 'Bitmap Triples' are always used in the current version of HDT.

HDT has also procedures for answering queries, but those will not be mentioned here, because they are not relevant for the thesis. [FMPG<sup>+</sup>13]

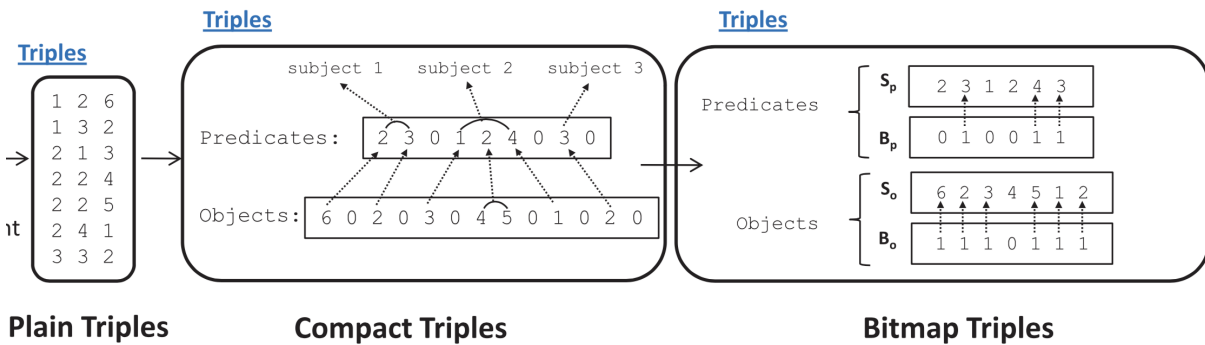


Figure 2.1: Three different representations of triples in HDT

## 2.3 Grammar-based Graph Compression

In this section we will discuss two different approaches to grammar-based graph compression. First, [Dü16] is introduced and briefly explained why the approach is less suitable for this thesis. Then [MP18], which the thesis will focus on, is introduced.

### 2.3.1 Dürksen’s Algorithm

An approach to grammar-based compression of graphs was developed in [Dü16]. Here the authors assume a hyper graph with node and edge labels. Such a graph can be seen in Fig. 2.2 on the left. To simplify the compression, a so-called transformation is now executed, in which a new node is inserted for each edge  $e$ , which then has the same label as  $e$ . The original structure of the graph is obtained by connecting two nodes, which were previously connected by an edge, indirectly by the new node. The result of the transformation is a graph, which only has node labels, but no edge labels. Moreover, hyper edges (edges with more than two incident nodes) are no longer present due to the transformation.

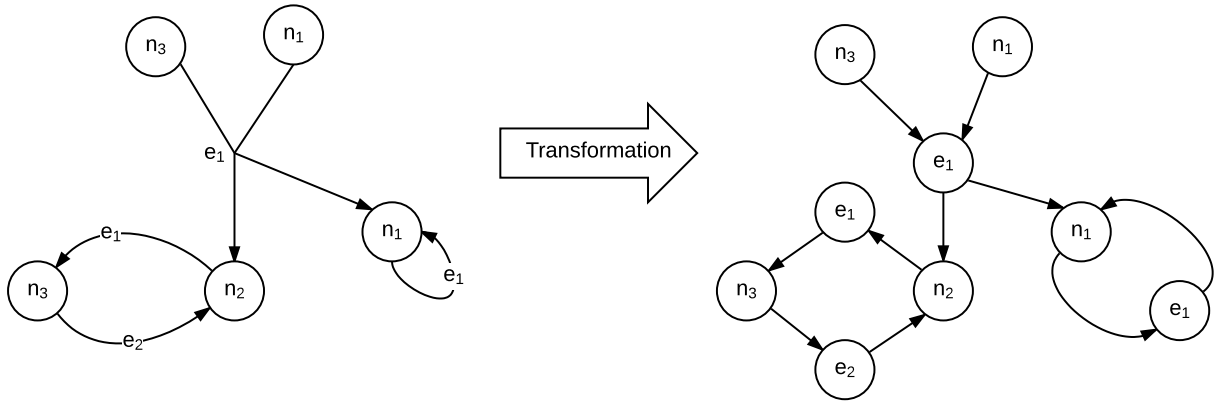


Figure 2.2: Transformation of a labeled graph (each edge is replaced by a new node having the label of the replaced edge).

Thus patterns like in Fig. 2.3 can be replaced. Here it happens twice that a node with the label  $X$  is connected with a node with the label  $Y$ . This pattern is then stored in a central location and can be referenced via the label  $X'$ . Thus only two nodes remain in the compressed graph (with  $X'$  as label) and the graph was reduced to a smaller size. There are some details which make this pattern replacement possible, but they are neglected at this point.

Dürksen’s algorithm does not seem to be suitable for RDF, because it is based on the fact that there are several nodes with the same label. However, since in RDF a node represents an entity that normally does not occur twice, Dürksen’s basic assumption is not fulfilled in RDF.

In addition, Dürksen’s algorithm is not yet in a mature state, i.e. there is only a rudimentary implementation that does not work reliably for large graphs. In addition, work is currently underway to find a compact representation of such a compressed graph in order to save the data with little memory. [Dü16]

For these reasons we will focus on the compressor GraphRePair, which is presented in chapter 2.3.2.

### 2.3.2 GraphRePair

This chapter deals with GraphRePair, the compressor from [MP18].

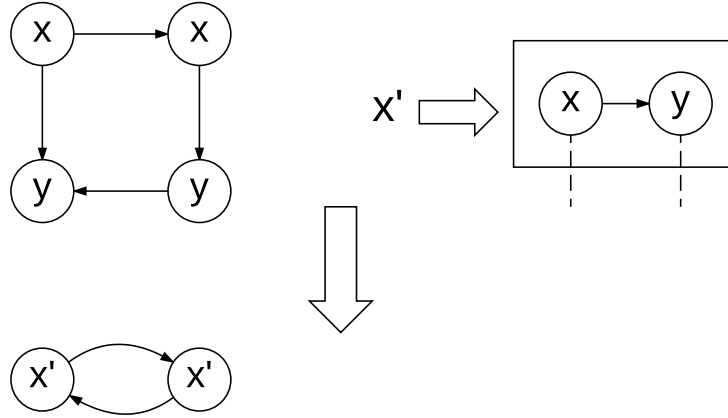


Figure 2.3: Replacement of two occurrences of the pattern  $X \rightarrow Y$  by nodes with the label  $X'$ .

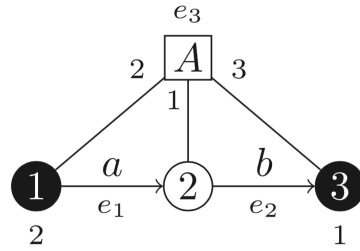


Figure 2.4: A hyper graph as defined in [MP18]. The black nodes 1 and 3 are external nodes whereas the white node 2 is internal.  $e_3$  (with the label A) is a hyper edge.

## Foundations

For a set  $M$ ,  $M^+ = \{x_1 \cdot x_2 \cdot \dots \cdot x_n \mid x_1, x_2, \dots, x_n \in M\}$  is defined as the set of all non-empty strings of  $M$  where  $\cdot$  stands for the concatenation of two symbols.  $M^* = M^+ \cup \{\epsilon\}$  is similar to  $M^+$ , but it also includes the empty string  $\epsilon$ .

A hypergraph over  $\Sigma$  is a tuple  $g = (V, E, att, lab, ext)$  where  $V = \{1, \dots, n\}$  is the set of nodes.  $E \subseteq \{(i, j) \mid i, j \in V\}$  is the set of edges.  $att : E \rightarrow V^+$  is the attachment mapping,  $lab : E \rightarrow \Sigma$  is the edge label mapping, and  $ext \in V^*$  is a series of external nodes.

A hypergraph does not contain multi-edges, which means for two edges  $e_1 \neq e_2$  it holds  $att(e_1) \neq att(e_2) \vee lab(e_1) \neq lab(e_2)$ .

For a hypergraph  $g = (V, E, att, lab, ext)$   $V_g, E_g, att_g, lab_g, ext_g$  are used to refer to its components.

An example for a hypergraph is illustrated in Fig. 2.4. Formally the graph can be described as  $V = \{1, 2, 3\}$ ,  $E = \{e_1, e_2, e_3\}$ ,  $att = \{e_1 \mapsto 1 \cdot 2, e_2 \mapsto 2 \cdot 3, e_3 \mapsto 2 \cdot 1 \cdot 3\}$ ,  $lab = \{e_1 \mapsto a, e_2 \mapsto b, e_3 \mapsto A\}$  and  $ext = 3 \cdot 1$ .

External nodes are black and the numbers below them indicate indices for their position in  $ext$ . Analogously, hyper edges ( $e_3$  in this case) have indices indicating the order of the attached nodes. [MP18]

## Digram

Digrams are like patterns, the same digram can occur multiple times in a graph. Based on that, GRP reduces a graph's size.

A digram  $d$  is a hyper graph with  $E_d = \{e_1, e_2\}$  so that the following conditions hold:

1.  $\forall v \in V_d : v \in att_d(e_1) \vee v \in att_d(e_2)$

2.  $\exists v \in V_d : v \in att_d(e_1) \wedge v \in att_d(e_2)$
3.  $ext_d \neq \epsilon$

Condition 1 ensures that all nodes in  $d$  are incident to one of the two edges of  $d$ . Conditions 2 is used to make sure that there is one 'middle node' incident to both edges of  $d$ . Finally, condition 3 ensures that there are external nodes. [MP18]

### Digram Occurrence

Let  $g$  be a hyper graph and  $d$  be a digram with the two edges  $e_1^d, e_2^d$ . Let  $o = \{e_1, e_2\} \subseteq E_g$  and let  $V_o$  be the set of nodes incident with edges in  $o$ . Then  $o$  is an occurrence of  $d$  in  $g$  if there exists a bijection  $b : V_o \rightarrow V_d$  so that for  $i \in \{1, 2\}$  and  $v \in V_o$  all following conditions hold:

1.  $b(v) \in att_d(e_i^d)$  iff  $v \in att_g(e_i)$
2.  $lab_d(e_i^d) = lab_g(e_i)$
3.  $b(v) \in ext_d$  iff  $v \in att_g(e)$  for some  $e \in E_g \setminus o$

Condition 1 and 2 ensure that the two edges of  $o$  form a graph isomorphic to  $d$ . Condition 3 makes sure that every external node of  $d$  is mapped to a node in  $g$  that is incident to at least one edge in  $g$  that is not contained in  $o$ . [MP18]

### Algorithm

Algorithm 1 is the main routine of GraphRePair. The algorithm take a graph as input and returns a grammar whereas  $N$  is the set of non terminals,  $P$  is the set of productions and  $S \in P$  is the start production. It maintains a list of digram occurrences in  $g$ . As long as this list contains multiple occurrences for one digram the loop will continue. In the loop, the most frequent digram is found and then replaced and the grammar is extended. After a digram replacement the occurrence list has to be updated, because the graph has now changed and former digram occurrences may not exist anymore. Moreover, new digram occurrences can be present now. This occurrence update is a complex process and will not be discussed here, since it is not relevant for the thesis. [MP18]

---

#### Algorithm 1 GraphRePair (Graph $g = (V, E, att, lab, ext)$ )

---

- 1:  $N, P \leftarrow \emptyset$
  - 2:  $S \leftarrow g$
  - 3:  $L(d) \leftarrow$  list of non-overlapping occurrences of every digram  $d$  appearing in  $g$
  - 4: **while**  $|L(d)| > 1$  for at least one digram  $d$  **do**
  - 5:    $mfd \leftarrow$  most frequent digram
  - 6:    $A \leftarrow$  new non terminal for  $mfd$
  - 7:   Replace every occurrence of  $mfd$  in  $g$
  - 8:    $N \leftarrow N \cup \{A\}$
  - 9:    $P \leftarrow P \cup \{A \rightarrow mfd\}$
  - 10:   Update the occurrence list  $L$
  - 11: **return** Grammar  $G = (N, P, S)$
-

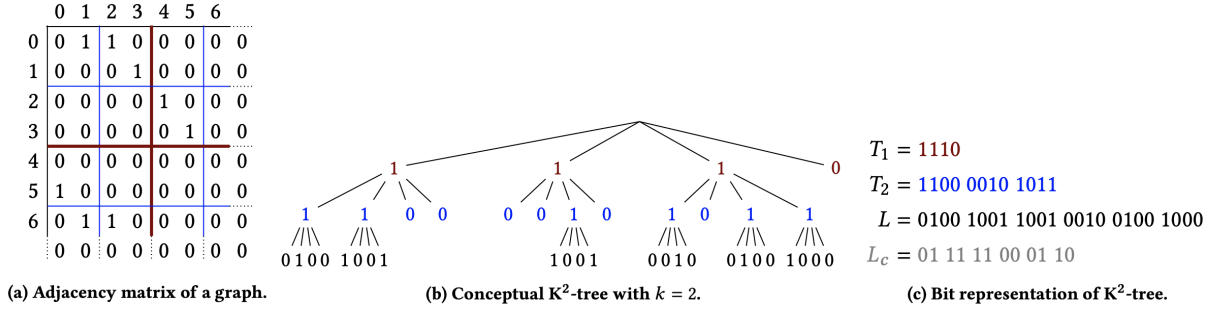


Figure 2.5: An adjacency matrix, its  $k^2$  tree representation and the tree's bit representation.

### Grammar Encoding

After GRP has constructed a grammar for some graph, this grammar has to be stored in an efficient way. The start production (the graph) is most often much bigger than the other productions and is therefore encoded differently.

The start production is encoded using  $k^2$  trees which is illustrated in Fig. 2.5.

First the matrix is extended with zeros to the next power of two. Then it is partitioned into  $k^2$  equally large partitions ( $k = 2$  here).

The tree's root now represents the whole matrix and its child order correspond to the order of the just created partitions. If a partition only contains zeros, a zero leaf is added to the tree (e.g. in Fig. 2.5 (a), partition 4, right bottom). Otherwise a one-node is added and the corresponding condition will be partitioned itself. The recursive procedures continues until there is zero-leaf for each path of the tree.

Afterwards, the tree is represented by bit-strings.

Since the graph has also edge labels, an adjacency matrix and its corresponding tree will be created for each of those labels.

The remaining productions of the grammar are encoded differently, because they are usually quite small compared to the start production. Here so-called  $\delta$ -codes with variable length are used. Essentially this is a way of displaying numbers as a bit-string in an efficient way (similar to a Huffman Code). To realize that, the right hand side of the components of the production have to be represented by numbers, e.g. the edges labels, the node IDs etc. [MP18]



## Approach

This chapter contains the solution approaches to the several tasks. It starts with some basic definitions, then continues with the comparison of HDT and GRP. Finally, approaches to improving the compression ratio will be presented.

### 3.1 Definitions

This chapter includes some basic definitions that will be used later on.

#### 3.1.1 Compression Ratio

One of the key metrics for a compressor is its compression ratio. The compression ratio depends on the input data and is defined as follows

$$\text{compression ratio} = \frac{\text{compressed size}}{\text{original size}} \in (0, 1]$$

The compression ratio is typically in the  $(0, 1]$  interval, since the compressed data will not be larger than the original data (there are cases where this happens, but it does not happen in the compressors considered here). Obviously, the compressed data cannot have a size of 0 or less. Normally the compression ratio is measured at the file size level (in byte). If a different measure is used, it will be mentioned at that point.

#### 3.1.2 (De-)Compression Time

Another key metric of a compressor is its (de-)compression time. This metric also depends on the input data and indicates the run time needed for compression and decompression of the data, respectively. The run time is typically measured in milliseconds.

### 3.2 GRP vs HDT

This chapter deals with the comparison between HDT and GRP. The main aim is to determine which of the two compressors achieves a lower compression ratio.

The obvious question is now whether there are certain properties/features that an RDF graph can have, and which have a positive or negative impact on the compression ratio of one or both algorithms.

### 3.2.1 Relation Between Structure of Data and Compression Ratio

First these features are considered for HDT. Fig. 3.1 is shown again. There you can see that the size of the data becomes smaller if there are only a few subjects. This is the case because the bit-array  $B_p$  contains a 1 every time a new subject is considered. For example, if there is only one subject, then  $B_p$  consists only of zeros.

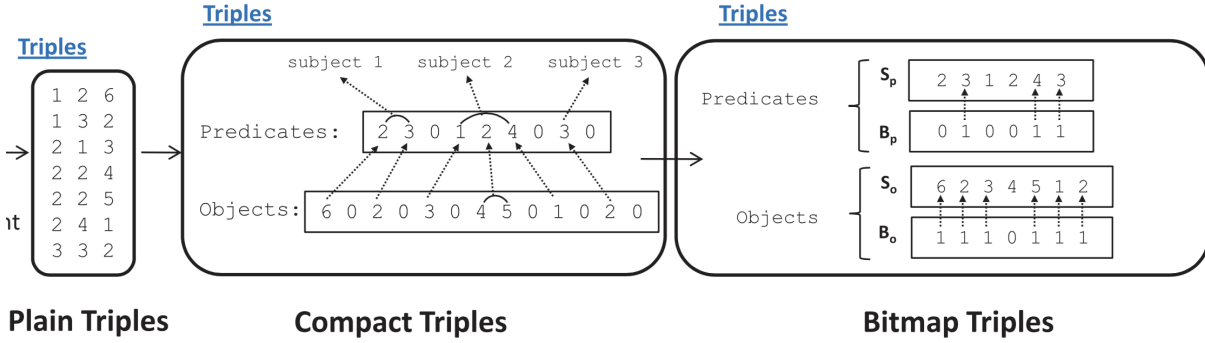


Figure 3.1: Three different representations of triples in HDT

For GRP that feature analysis is more complex. Since GRP constructs grammar rules by using the graph's structure, it is more dependent on the data's structure than HDT. There can be many features that can lead to different constructible grammar rules. Those will not be discussed here. But a simple insight is that GRP's compression ratio will be bigger when there are more different predicates in the graph. This is true, because GRP's grammar rules are based on repeating edge labels. That fact will be considered in the evaluation.

### 3.2.2 Dictionary Size

As already explained in Ch. 2, HDT divides an RDF file into its header, dictionary and triples component. This is partially also true for GRP except that GRP does not create a header. But it also assigns an ID to every URI or literal and then only works with these IDs. Unfortunately the authors of [MP18] did not work on efficient storing of the dictionary. In GRP the size for the dictionary component is just ignored. Therefore we have to add this size in order to compare GRP with HDT in a fair way. To achieve that we just add the same size to the compressed size of GRP that HDT would need to store the dictionary.

## 3.3 Compression Improvements

Ch. 5.2 will show that GRP always achieves a better compression ratio than HDT. Since according to [MP18] a graph compressed with GRP is less suitable for neighborhood queries (i.e. typical SPARQL queries) than an uncompressed graph, we will concentrate in the following on further improving the compression ratio. So we do not compare the query speeds between HDT and GRP.

### 3.3.1 Ontology Knowledge

RDF has meta data that contains knowledge about the actual data. This is also called ontology. An ontology is normally itself an RDF graph. There are two known languages for formulating

an ontology: RDFS <sup>1</sup> and OWL <sup>2</sup>. Of these, OWL is the more powerful and is therefore chosen here.

This chapter is about finding out whether one can change the structure of an RDF graph by applying knowledge from its ontology so that it is more compressible for GRP, but at the same time remains semantically equivalent to the original graph. In this way no data would be lost by compression.

In Ch. ?? it has already been mentioned that GRP depends more on the structure of the input graph than HDT does. It will therefore be interesting to see how applying ontology knowledge influences GRP's compression ratio.

This chapter is about elaborating the theoretical concepts of OWL and investigating how they can be used for grammar-based compression.

Let

$$elr = \frac{\text{number of different edge labels}}{\text{number of edges}}$$

(edge label ratio) be the ratio of the edge labels or properties to the total number of edges of the graph.

Generally it can be said that GRP can compress a graph better if *elr* is lower, because then there is more redundancy in the graph. However, if the graph structure becomes unfavorable for GRP, the compression ratio may still be worse at a lower value for *elr*.

### Symmetric Properties

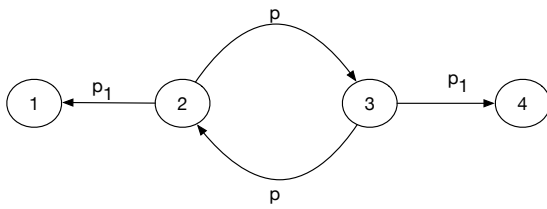
There is a predicate in [owl] called “owl:SymmetricProperty” which expresses that a certain other predicate *p* is symmetric. This means, if there is a triple (*s, p, o*) in the graph, then there can also be a triple (*o, p, s*) at the same time. In reality, however, it can happen that only one of the two triples really exists in the graph. The idea now is to always add the other triple to the graph in such a case. This makes the graph larger at first, but more grammar rules can be found. This is because you make *elr* smaller by adding it, which can lead to a better compression ratio. At the same time you should not get an unfavorable structure. The procedure is illustrated in Fig. 3.2a. That graph shall be seen as a sub graph of a much larger graph. Here the predicate *p* is symmetric, so the edge from node 3 to 2 was added, *p*<sub>1</sub> is not symmetric. Due to the addition, the digram of Fig. 3.2b can now be found twice, whereas it was previously found only once. These two occurrences overlap and therefore cannot be replaced both. However, it may be that one of the two occurrences cannot be replaced because the nodes involved are still connected to other nodes that are not shown in this figure. So the addition increases the probability that the digram can actually be replaced. At the same time the degree of nodes 1 and 2 is increased by one. But this should not really decrease the chance of finding other digrams in the graph, since 1 and 2 have already been connected before.

### Inverse Properties

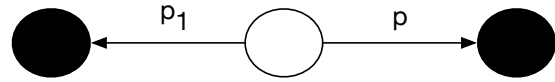
A further concept of [owl] is “owl:inverseOf”, which is defined for two properties *p*<sub>1</sub>, *p*<sub>2</sub>. If (*s, p*<sub>1</sub>, *o*) exists then (*o, p*<sub>2</sub>, *s*) should also exist and vice versa. However, it can happen that only one of the triples really exists. The approach is now similar to the symmetric properties. One can argue in a similar way for adding those edges to graph here. It can decrease *elr* if there are many occurrences of a few inverse properties. Also, adding those edges will not really make the graph's structure more complex since all the nodes have been connected before.

<sup>1</sup><https://www.w3.org/TR/rdf-schema/>

<sup>2</sup>[https://www.w3.org/TR/2012/REC-owl2-overview-20121211/#Documentation\\_Roadmap](https://www.w3.org/TR/2012/REC-owl2-overview-20121211/#Documentation_Roadmap)



(a) A sub graph to which the edge from node 3 to node 2 was added.



(b) The digram that can be found twice in the graph of Fig. 3.2a

Figure 3.2: Visualization of the benefits of adding symmetric edges to the graph.  $p$  is symmetric,  $p_1$  is not symmetric.

### Transitive Properties

### Equal Properties

In OWL there are the predicates “owl:equivalentProperty” and “owl:sameAs”. The first one denotes that two properties are equivalent, but it does not mean that they are equal. But the latter one is expressing equality. If there is  $p$  which is equal to other properties  $p_1, \dots, p_n$ , then one approach could be to replace each occurrence of  $p_1, \dots, p_n$  with  $p$ . This would reduce *elr* and, at the same time, not change the structure of the graph. But it is questionable if such cases exist in reality.

### 3.3.2 Dictionary Improvements

As already seen in Ch. ??, the dictionary makes up most of the memory of a compressed RDF graph. It is therefore worth investigating whether the dictionary can be compressed better. One can take advantage of certain features of the dictionary to achieve that.

As mentioned above, GRP does not have its own method for compressing the dictionary. We have therefore taken the compression method from HDT, and applied it in GRP to ensure a fair comparison.

HDT has a fairly mature mechanism for compressing the dictionary.

genauer erklären

### Literals

Objects in RDF can be literals. Literals typically contain constant values and usually have no common prefixes. Therefore the compression of HDT is not suitable for these. Since literals can often contain whole flow texts, a text compression would probably be well applicable. An example of such a text compression is a Huffman Code [Sha10]. Here the text is converted into a binary format. Every single character of a text is binary coded, whereby frequently occurring characters get short and rare characters get longer codes. These codes are expressed by a binary Huffman tree. An example can be seen in Fig. 4.3. Each leaf contains a symbol whereas the one and zeros on the path to the symbol define its code. The tree is constructed in such a way that paths to frequent characters are shorter than those to rare characters. The whole procedures can be seen in [Sha10].

This tree must then be stored in addition to the compressed data so that the original data can be recovered. It will be seen in Ch. 4 how that is done.

## Blank Nodes

Blank nodes are normally used when a node does not get a URI, but is still necessary to represent a statement. Such a node is often used to formulate more complex logical statements. The same blank node can occur in different triples. In order for it to be referenced uniquely, it gets an ID. These IDs are usually chosen arbitrarily and have no meaning beyond that. When reading an RDF graph with the Jena-API <sup>3</sup> (which is used by HDT) random UUID strings are assigned for the blank nodes, which are quite long. They also have no common prefixes, which makes the HDT dictionary compression ineffective again.

To improve compression, one can reassign the IDs of the blank nodes. For example you can use numbers from 1 to  $n$  ( $n$  = number of blank nodes) to have short IDs.

Another possibility is not to save the IDs of the blank nodes at all. In HDT all strings in the dictionary (including the blank node IDs) are mapped to short IDs. Thus the blank node IDs are in principle already stored. They can therefore be removed from the dictionary. HDT must then be changed so that it can handle the case in which it does not find a corresponding string in the dictionary for a certain short ID. At this point it would know directly that the considered node is a blank node and the longer blank node ID is unimportant. Such a situation will occur when a decompression is performed.

In Ch. ?? one can see the effects of shortening and becoming blank node IDs on the compression ratio.

---

<sup>3</sup><https://jena.apache.org/index.html>



# Implementation

## 4.1 GRP vs HDT

In this chapter, it will be explained how the comparison between HDT and GRP is implemented.

### 4.1.1 Synthetic Data Creation

As has just explained, HDT can compress a graph that is similar to a hub pattern (few subjects, many objects) very well. So it gets worse the further away the graph is from this pattern. This corresponds to the authority pattern, where there are few objects but many subjects.

The task now is to create a series of RDF graphs that first correspond to the hub pattern and then continue to change in the direction of the authority pattern. This is illustrated in Fig. 4.1. In this scenario all graphs ( $G_1$  to  $G_m$ ) have the same size, i.e. the same number of nodes and edges.  $G_1$  has only one subject connected to all objects.  $G_2$  then has two subjects more and correspondingly 2 objects less. This goes on and on until there is only one object that is connected to all subjects ( $G_m$ ). The edges are randomly distributed among the nodes, so that all nodes have a similar degree and each node has at least a degree of one.

It is also ensured that each of the generated files has exactly the same size. This is made possible by ensuring that each URI has the same length. Since the RDF graph also has the same number of triples, the files are of the same size. Since the evaluation compares the compression ratios for the different RDF files, it is important that all files are of the same size to ensure a fair comparison.

A section of such a file (for  $G_1$ ) is shown in Fig. 4.2. In that example, there is only one distinct predicate for all triples. The number of predicates is always one at first. The amount of predicates has a similar effect on both compressors and is therefore omitted at first. But in Ch. ?? that effect will be discussed in more detail.

Apart from that, blank nodes and literals are not used here. For both compressors, blank nodes and literals are being handled analogously to URI nodes. Therefore they are not needed at this point, in order to show the behavior of HDT and GRP in the above described scenario.

## 4.2 Compression Improvements

This chapter introduces the implementation details of Ch. 3.3. It starts with the applying the ontology in order to achieve a better compression ratio and then continues with dictionary compression improvements.

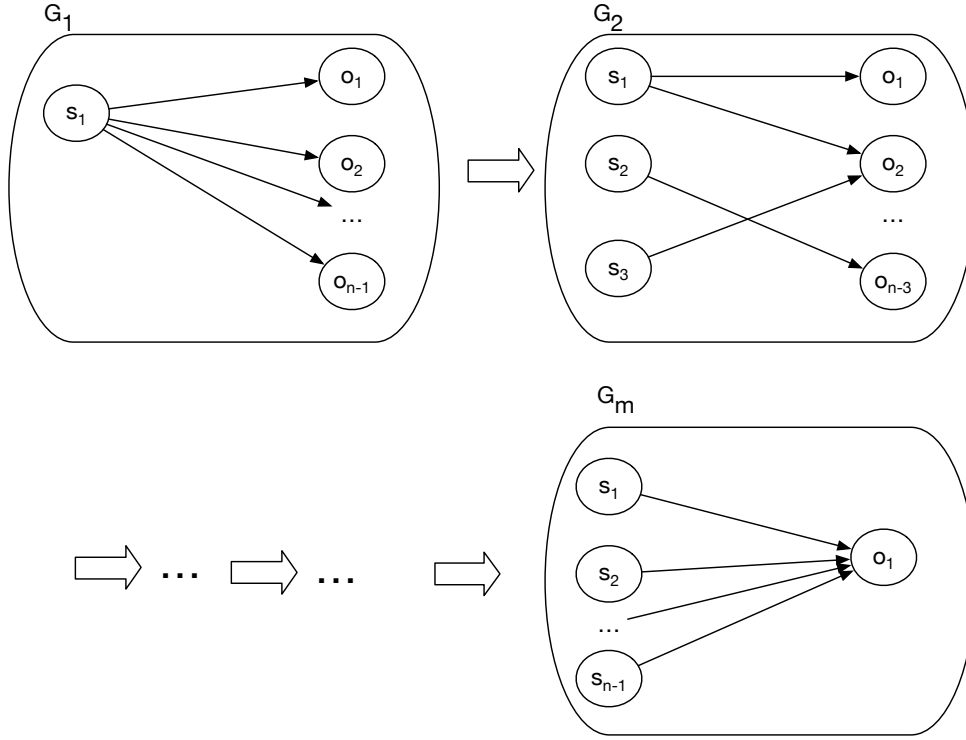


Figure 4.1: Step-by-step transition from hub pattern to authority pattern. The number of nodes  $n$  is the same for each graph. The number of edges is also the same for each graph.

```
<http://subject/0000000000> <http://predicate/0000000000> <http://object/0000001036> .
<http://subject/0000000000> <http://predicate/0000000000> <http://object/0000000854> .
<http://subject/0000000000> <http://predicate/0000000000> <http://object/0000000991> .
<http://subject/0000000000> <http://predicate/0000000000> <http://object/0000000450> .
<http://subject/0000000000> <http://predicate/0000000000> <http://object/0000000456> .
<http://subject/0000000000> <http://predicate/0000000000> <http://object/0000000689> .
<http://subject/0000000000> <http://predicate/0000000000> <http://object/0000001166> .
<http://subject/0000000000> <http://predicate/0000000000> <http://object/0000000960> .
```

Figure 4.2: Excerpt from the generated RDF file for  $G_1$  (see Fig. 4.1). Each triple has the same length.

### 4.2.1 Ontology Knowledge

#### Symmetric Properties

In order to remove or add symmetric properties SPARQL can be used. The code of listing 4.1 will be used to do that:

```
INSERT {?o ?p ?s}
WHERE{
    {?o ?p ?s}
    MINUS {?o ?p ?s}
}
```

Listing 4.1: SPARQL query for adding triples with the symmetric property  $p$

That update has to be executed for each symmetric property  $p$ . In the case where one wants to remove the second, a delete update would have to be executed.



### 4.2.2 Dictionary Improvements

For the dictionary improvements, HDT's dictionary compression is used as a foundation. Therefore, the HDT-Java code <sup>1</sup> has been extended.

#### Blank Nodes

HDT normally uses the arbitrary and long UUIDs generated by the Jena API and tries to compress them using prefix trees. As already explained in Ch. 3.3.2 it is possible to use shorter strings (e.g. numbers from 1 to  $n$ ). Then, during run time a mapping from old ID to new ID is maintained in order to make sure that the same blank node will get the same new ID if it occurs multiple times. That mapping does not have to be stored persistently and will therefore be omitted once the compression is finished.

The other possibility is to omit blank node ID completely. That can be achieved quite easily. The HDT code is changed in such a way that skips storing blank nodes.

Both approaches have been implemented and will be evaluated in Ch. 5.3.2.

#### Literals

There is no standard way for storing the binary tree. One approach can be seen in Algorithm 2 which has to be started with the root node. That method creates an unambiguous bit representation of the tree, since each node has either two or no children.

---

#### Algorithm 2 EncodeNode (TreeNode node)

---

```

1: if node is leaf then
2:   writeBit(1)
3:   writeCharacter(node.character)
4: else
5:   writeBit(0)
6:   EncodeNode(node.leftChild)
7:   EncodeNode(node.rightChild)

```

---

Alternatively, there are pre-computed Huffman trees for natural languages such as English. There they have already investigated which letter occurs how often in English texts and in this way a generally valid Huffman code has been established. The advantage is that one does not have to save the Huffman tree and does not have to calculate it oneself, which saves runtime. The disadvantage, however, is that the tree is not optimal for the text to be compressed, as it is more general. Another problem in our case is that the literals contain a lot of special characters that are not taken into account in prefabricated Huffman codes. Also, in Ch. 5.3.2 it will turn out that saving the Huffman tree does not require much additional memory. Therefore, we choose to calculate the Huffman code ourselves to compress the literals. All literals of the input graph have to be traversed in the beginning to compute the character frequencies.

---

<sup>1</sup><https://github.com/rdfhdt/hdt-java>

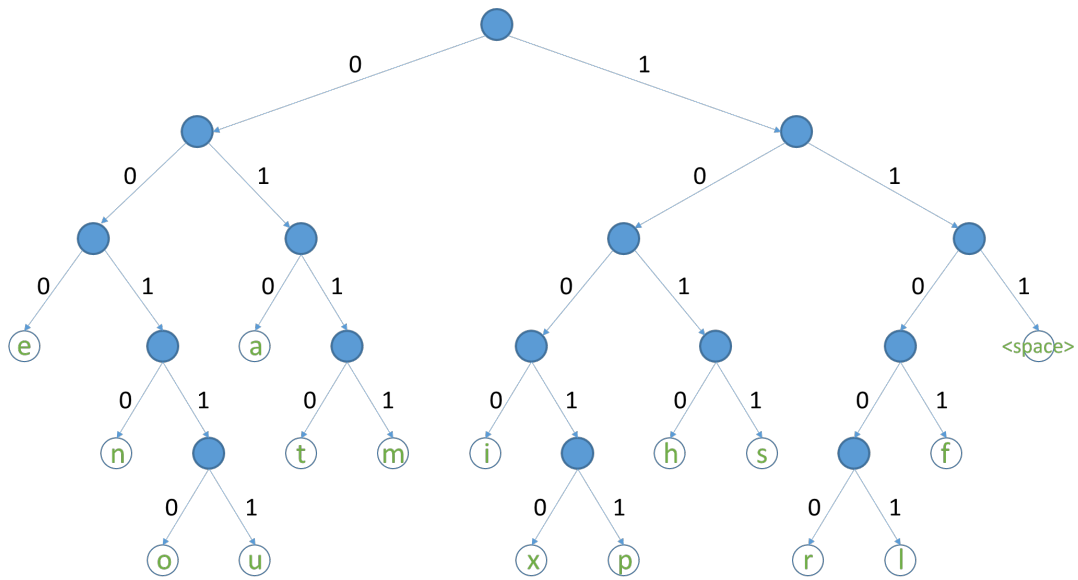


Figure 4.3: An example of a Huffman Tree.

This chapter is about the evaluation of the several approaches presented in Ch. 3 and 4.

## 5.1 Experimental Setup

For the following evaluations HDT-Java 2.0 <sup>1</sup> (the currently newest version) has been used. For GRP the implementation mentioned in [MP18] has been used (written in Scala). It is not open source and has been given to us by the authors of [MP18].

### 5.1.1 Datasets

In this chapter, an overview about the different datasets, that were evaluated, is given.

#### Semantic Web Dog Food

Semantic Web Dog Food <sup>2</sup> is a collection of RDF files from the RDF researchers community. It contains data about their conferences and workshops.

#### DBPedia

DBPedia <sup>3</sup> is an RDF version of the knowledge from Wikipedia. It contains many different data files of a quite big size, with some of them having hundreds of millions of triples. Also, DBPedia includes an ontology that is used for the evaluation.

## 5.2 GRP vs HDT

Fig. 5.1a shows the compression ratio for HDT (without dictionary size). As expected, the ratio gets higher the more similar the graph is to the authority pattern. In general it can be said that this effect is quite small. There is only a distance of 0.002 between the minimum and the maximum. This small effect can also be seen by looking at Fig. 5.1b. It can be seen that the size of the dictionary has a much bigger effect, since the compression ratio is now much larger and the curve behavior from Fig. 5.1a is no longer recognizable. It is noticeable that the

---

<sup>1</sup><https://github.com/rdfhdt/hdt-java/releases/tag/v2.0>

<sup>2</sup><http://www.scholarlydata.org/dumps/>

<sup>3</sup><https://wiki.dbpedia.org/Downloads2015-04>

dictionary size gets bigger when the graph is further away from the star pattern. The dictionary implementation of HDT seems to be more inefficient when there are about as many subject as objects.

dafür sorgen, da  
kurve rechts au  
den knick hat

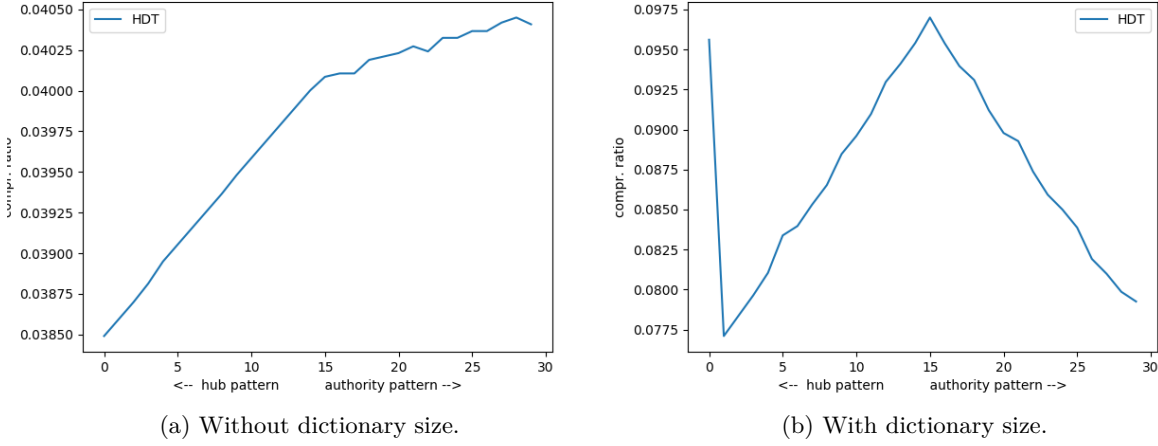


Figure 5.1: The compression ratios for HDT without and with dictionary sizes.

Next the compression ratio of GRP is considered, which is presented in Fig. 5.2a (without dictionary sizes). Here one can see that GRP has a better compression ratio if the graph is more similar to the star pattern (hub order authority pattern). This property of GRP has also been mentioned in [MP18]. It can also be seen that the effect on the compression ratio is bigger for GRP than for HDT (standard deviation is twice as high for GRP as for HDT). A grammar-based compression is therefore more dependent on the structure of the input data.

When Fig. 5.2b is considered, it can be seen that this curve behaves almost exactly like the one from Fig. 5.1b, since the size of the dictionary accounts for most of the compressed data size.

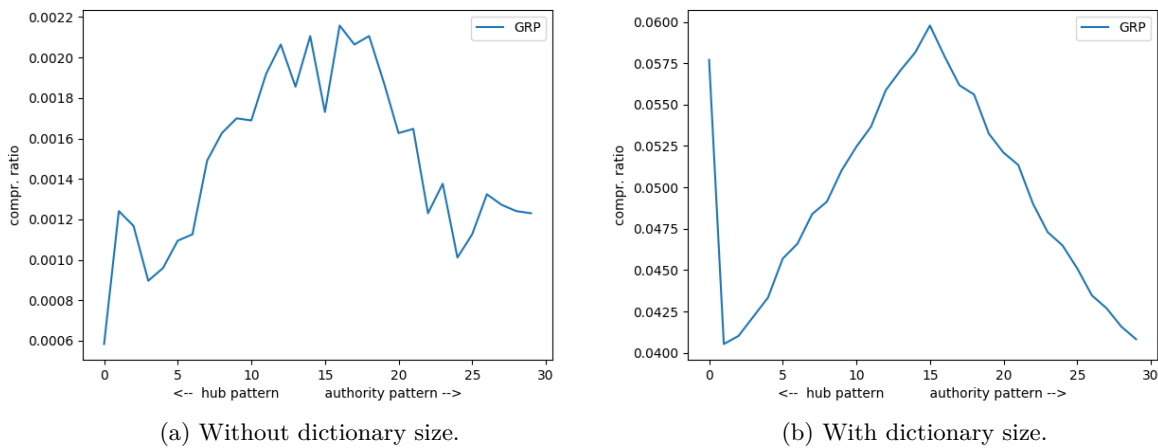


Figure 5.2: The compression ratios for GRP without and with dictionary sizes.

Finally, Fig. 5.3a and Fig. 5.3b show the compression ratios for both algorithms. Since both use the same method to compress the dictionary, the curves in Fig. 5.3b are very similar. However,

it becomes clear that GRP compresses better than HDT. In Fig. 5.3a, the ratio of HDT is 31 times higher on average. Of course, this factor becomes much smaller in Fig. 5.3b because the dictionary accounts for most of the memory size. Here the compression ratio of HDT is on average 1.8 times as high as that of GRP.

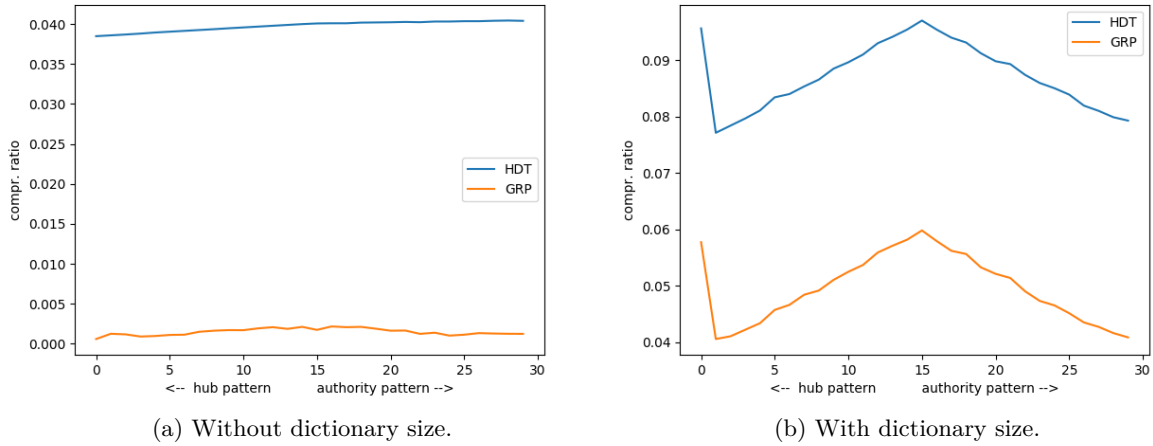


Figure 5.3: The compression ratios for GRP and HDT without and with dictionary sizes.

One could now argue that only one distinct predicate was used in that scenario and this is beneficial for GRP, as it gets worse as the number of predicates increases. Therefore a further evaluation is made in Fig. 5.4, where 1000 distinct predicates have been used. That is a quite high number considering the number of triples (1199) compared to real RDF data. One can see that the compression ratios are now higher for both compressors, but still GRP's ratio is always smaller than HDT's. HDT's ratio is still 1.7 times higher on average. So, the increasing number of predicates has a similar effect on both algorithms.

Apart from the compression ratio, the run time is also important for the overall performance. Fig. 5.5 shows the average run times of the two algorithms. For this the same scenario with the star pattern (and only one distinct predicate) was used. It has been executed 100 times to get a sophisticated run time measurement, because the run time also depends on the current CPU workload of the computer.

It can be seen that the runtime of HDT is significantly higher than that of GRP. It is on average ca. 48 times as high.

However, it should also be noted that the implementation of GRP is rather rudimentary (according to the authors of [MP18]), while that of HDT has been under development for some time. So they are not comparable in terms of quality. Unfortunately, one cannot say at this point whether a more professional implementation of GRP will also be slower than HDT.

In addition one can notice that GRP's run time fluctuates more than that of HDT. GRP has a standard deviation of about 134, while HDT only has a standard deviation of about 7. One reason for this is that GRP, in contrast to HDT, is non-deterministic because of the partly random search order of the graph. On the other hand, the high deviations are also a confirmation of the above mentioned hypothesis that the behavior of GRP depends more on the structure of the input data than HDT does.

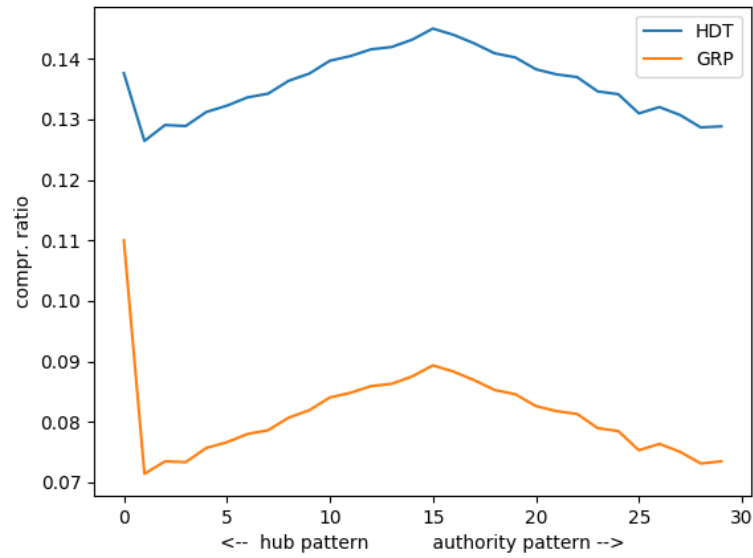


Figure 5.4: The compression ratios for GRP and HDT with dictionary sizes. Graphs have now 1000 distinct predicates.

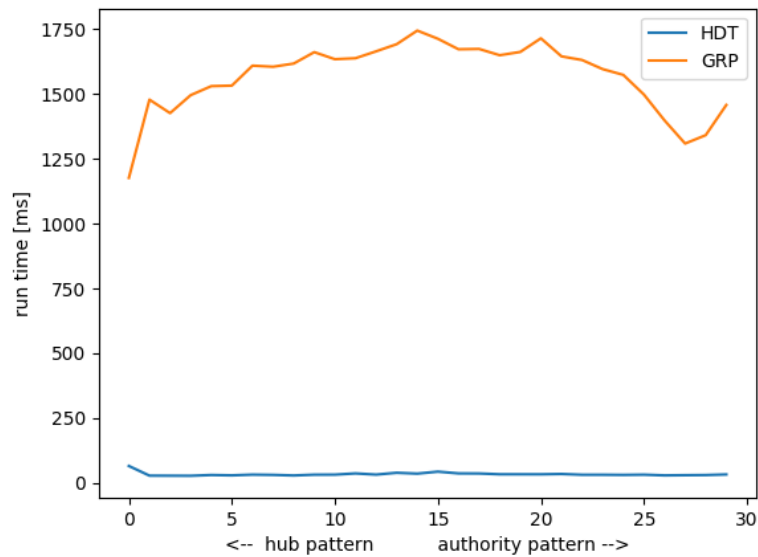


Figure 5.5: Run times of both algorithms (average run time of 100 consecutive executions).

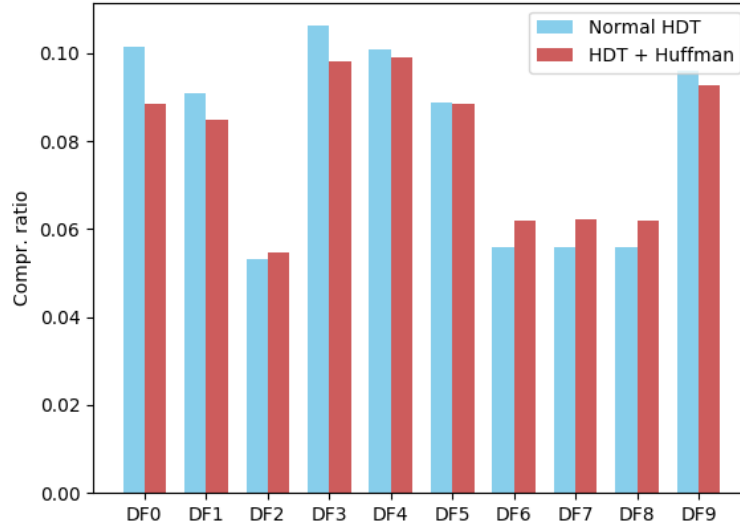


Figure 5.6: Compression ratios for Semantic Web Dog Food Files. Comparison between Normal HDT and HDT + Huffman.

## 5.3 Compression Improvements

### 5.3.1 Ontology Knowledge

### 5.3.2 Dictionary Improvements

#### Literals

As already mentioned, a self-generated Huffman code was used to compress the literals of an RDF graph, as one hopes for a better compression than the with the prefix-based compression of HDT.

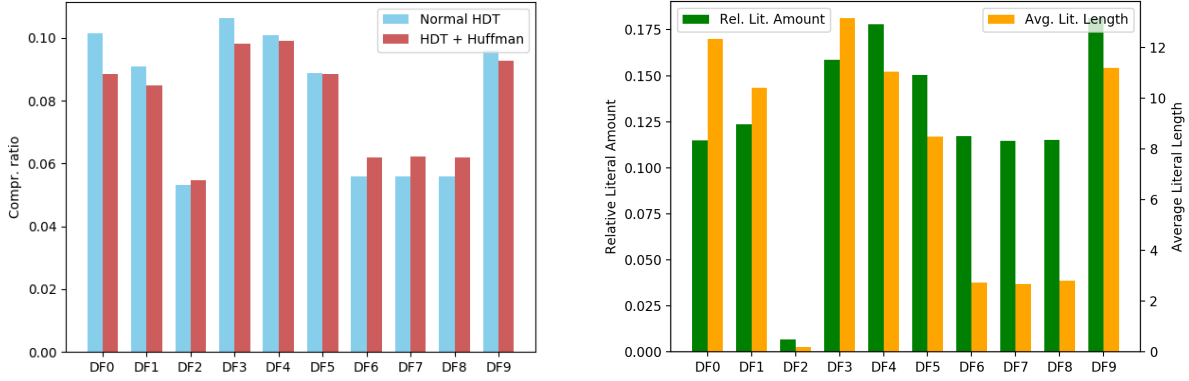
Now the results of the evaluation are presented. First, the Semantic Web Dog Food data set is used. This is well suited for an evaluation, since both literals and URIS are included as objects. Fig. 5.6 shows the compression rates for a subset of the data (DF0 - DF9).

One can see that the use of the Huffman code only brings a small improvement in some cases. In some cases Normal HDT even compresses better than HDT with Huffman.

This is because these data do not have a very high proportion of literals and literals are also rather short. This can be seen in Fig. 5.7, where Fig.5.7a lists compression rates again and Fig.5.7b shows the relative proportion of literals and the average length of a literal. The proportion of literals in the triples is never higher than 17.5% and the highest average literal length is about 12. So these are single words rather than whole texts. Only in cases where both the proportion and the literal length are relatively high, an improvement can be seen (e.g. with DF9).

Now the DBpedia data set is considered, because it has a different structure of data. Here, graphs are considered which contain the abstracts of Wikipedia, because these are longer texts. In addition there are abstracts in different languages, so you can see if some languages are better suited for Huffman than others.

Fig. 5.8a shows the compression rates for the abstracts. The abbreviations stand for the languages in which the abstracts are written. It can be seen that Huffman significantly improves the compression rates here. The reason for this can be seen in Fig. 5.8b, where the average

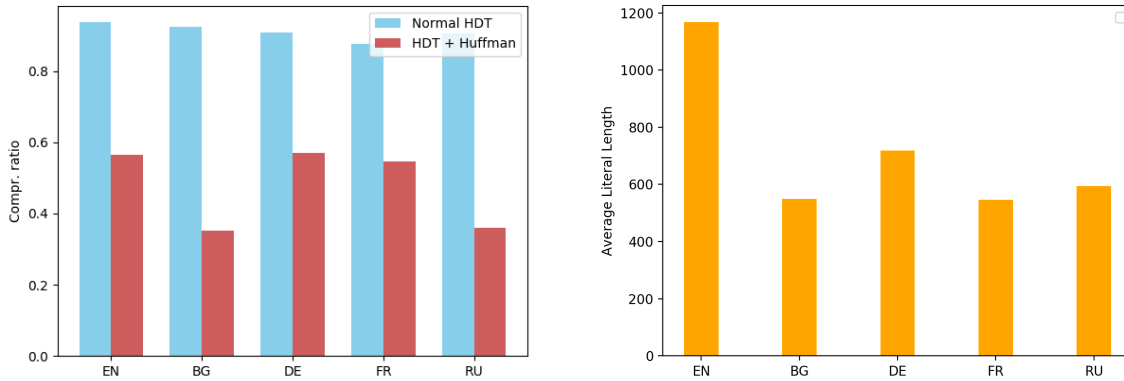


(a) Compression ratios for Normal HDT and HDT + Huffman (b) Relative amounts of literals and average literal length.

Figure 5.7: Relation between literal percentage and literal length (right) and compression ratios (left) for Semantic Web Dog Food Data.

literal length is illustrated. The relative proportion of literals is not shown this time, since it is 100% for all files. The average literal length varies between languages, but is generally much higher than in Semantic Web Dog Food. Therefore, the improvement is much greater here.

Another phenomenon can also be seen here: Although the English file contains by far the longest literals, the improvement by Huffman here is not as big as, for example, in the Bulgarian file. In the English case, the Huffman code is less effective because there are fewer different characters than in the other languages. Huffman is generally more efficient on large alphabets, according to [Sha10].



(a) Compression ratios for Normal HDT and HDT + Huffman

(b) Average literal length.

Figure 5.8: Relation between literal length (right) and compression ratios (left) for DBpedia Abstracts (Abbreviations denote the languages the abstracts are written in).

As mentioned above, saving the Huffman code means almost no additional storage effort. The average fraction of the Huffman code is about 0.1% and was therefore not displayed in the

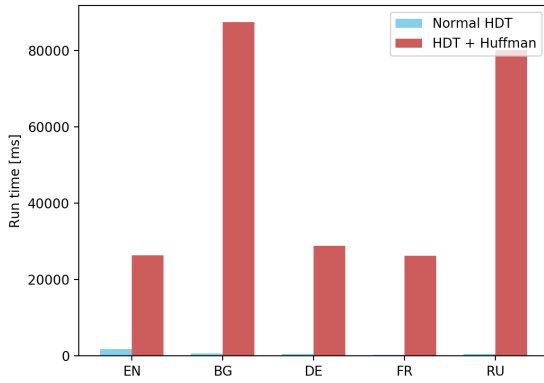


visualizations.

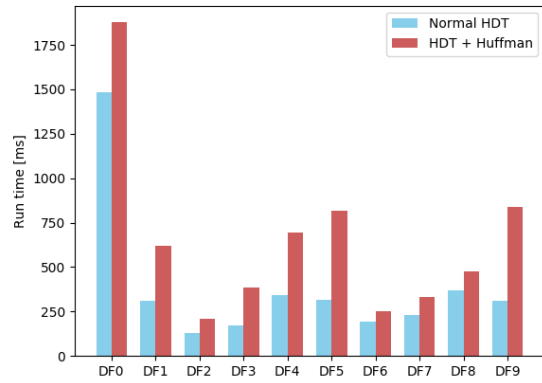
Since the calculation of the Huffman code increases the runtime of the whole compression, it is now considered how big this effect is. Fig. 5.9a shows the run times for the compression of the DBPedia abstract data. Here the runtime has been increased very much. This is because HDT can normally compress very little with this data because of the many long literals and is therefore finished quite quickly. With Huffman one can compress much better and it takes quite a long time.

In Fig. 5.9b you can see the run times for the Semantic Web Dog Food data, where the run times are only slightly longer, because Huffman hardly improves the compression.

So it can be said that the use of a standard Huffman code could be worthwhile because of the otherwise high runtime. At this point, however, the evaluation with such a standard code was omitted, since such existing codes do not contain all the special characters that occur in RDF data.



(a) Run times for DBPedia abstracts.



(b) Run times for Semantic Web Dog Food files.

Figure 5.9: Run times for Normal HDT and HDT + Huffman.

## Blank Nodes

### Literals and Blank Nodes Combined



## Summary and Discussion

### 6.1 Summary and Discussion

### 6.2 Future Work



## Todo list

■ hier evtl Beispiele von knowledge bases . . . . .	1
■ here possibly concrete numbers . . . . .	2
■ fix . . . . .	11
■ genauer erklären . . . . .	12
■ das passt vom Text her nicht . . . . .	17
■ dafür sorgen, dass kurve rechts auch den knick hat . . . . .	20
■ belegen . . . . .	21



# Bibliography

- [Dü16] Matthias Dürksen. Grammar-based Graph Compression. Bachelor’s thesis, University of Paderborn, 2016.
- [FMPG<sup>+</sup>13] Javier D. Fernández, Miguel A. Martínez-Prieto, Claudio Gutiérrez, Axel Polleres, and Mario Arias. Binary rdf representation for publication and exchange (hdt). *Web Semantics: Science, Services and Agents on the World Wide Web*, 19:22 – 41, 2013.
- [MP18] Sebastian Maneth and Fabian Peternek. Grammar-based graph compression. *Inf. Syst.*, 76:19–45, 2018.
- [owl] Owl reference. <https://www.w3.org/TR/owl-ref/>. Accessed: 2019-05-07.
- [Sha10] Mamta Sharma. Compression using huffman coding. 2010.