

# HDT

## Queryable compression format for Linked Data

Wouter Beek, VU University Amsterdam

Javier D. Fernández, Vienna University of  
Economics and Business, Austria

Ruben Verborgh, Ghent University – imec,  
Belgium

ISWC 2017

21<sup>TH</sup> SEPTEMBER 2017



# Practical tips

- WIFI
  - SSID: wu-conference
  - Username: wu0041
  - Password: SWSA2017
- Get slides: <https://rdfhdt.github.io/ISWC2017/>
- Please write/follow the minutes at: <https://board.net/p/HDT-ISWC17>
- Twitter: @rdfhdt

# General agenda

- **Session I** (09:00 - 9:40) Welcome and HDT foundations
- **Session II** (9:40 - 10:10) Practical uses I: Linked Data Fragments
- **Session III** (10:10 - 10:30) Practical uses II: LOD Laundromat
-  10:30 – 11:00
- **Session IV** (11:00 - 11:45) Presentations - Different perspectives on the topic
  - "Practical use of HDT within **data.world**". Byron Jacob and Jonathan Ortiz, data.world
  - "HDT as the backend engine for **Query Answering** over the Web of Data". Dennis Diefenbach, Laboratoire Hubert Curien, Saint Etienne, France.
  - "Counting to K: **Top k Shortest Path** Queries with HDT". Vadim Savenkov, Vienna University of Economics and Business, Austria.
- **Session V** (11:45 - 12:30) Discussion: concrete next steps and closing remarks



# Knowledge of HDT

0. Zero knowledge
1. I have just heard of HDT, I know some of the potential application scenarios
2. I know the basic foundations of HDT and I gave it a try
3. I often manage HDT datasets and tools

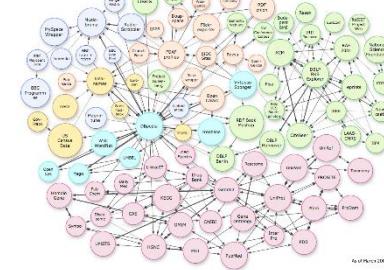


# Motivation. Origins

- 'Simple' task in 2009 (by Claudio Gutiérrez  )
- Let's inspect what people are publishing in RDF
  - Find RDF datasets
  - Download them
  - Do some (simple) queries to inspect the content
- Problems?
  - Discover datasets
  - Huge resources to download (large) datasets
  - Messiness of the data
  - Huge resources to index (large) datasets locally
  - Huge resources to query (large) datasets locally and to serve them online



2009



Is it much better now in 2017 ??

img: Beth Scupham

# The problem is in the roots (Big tree = big roots)

- Publication, Exchange and Consumption of **large RDF datasets**
    - Most RDF formats (N3, XML, Turtle) are text serializations, designed for human readability (not for machines)
      - Verbose = High costs to write/exchange/parse
      - A basic offline search = (decompress)+ index the file + search



“the published RDF dumps are actually bulks with no structure, no design, no final user in mind. They resemble unwanted creatures whose owners are keen to be rid of them”



# Header Dictionary Triples (HDT)

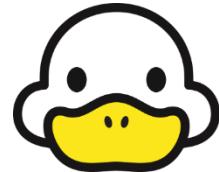


“

HDT compresses big RDF datasets while maintaining search operations

“HDT might not be the most compressed format or the fastest RDF store, but it should be the most efficient way to publish, exchange and consume RDF datasets.

Like a duck, it may not walk, swim and fly as fast as other animals ... but he does everything!”



# Motivation

- Lightweight Binary RDF (HDT)
  - Highly compact serialization of RDF
  - Allows fast RDF retrieval in compressed space (without prior decompression)
    - Includes internal indexes to solve basic queries with small (3%) memory footprint.
    - Very fast on basic queries (triple patterns), x 1.5 faster than Virtuoso, RDF3x.
    - Complex queries (joins) on the same scale of current solutions (Virtuoso, RDF3x).

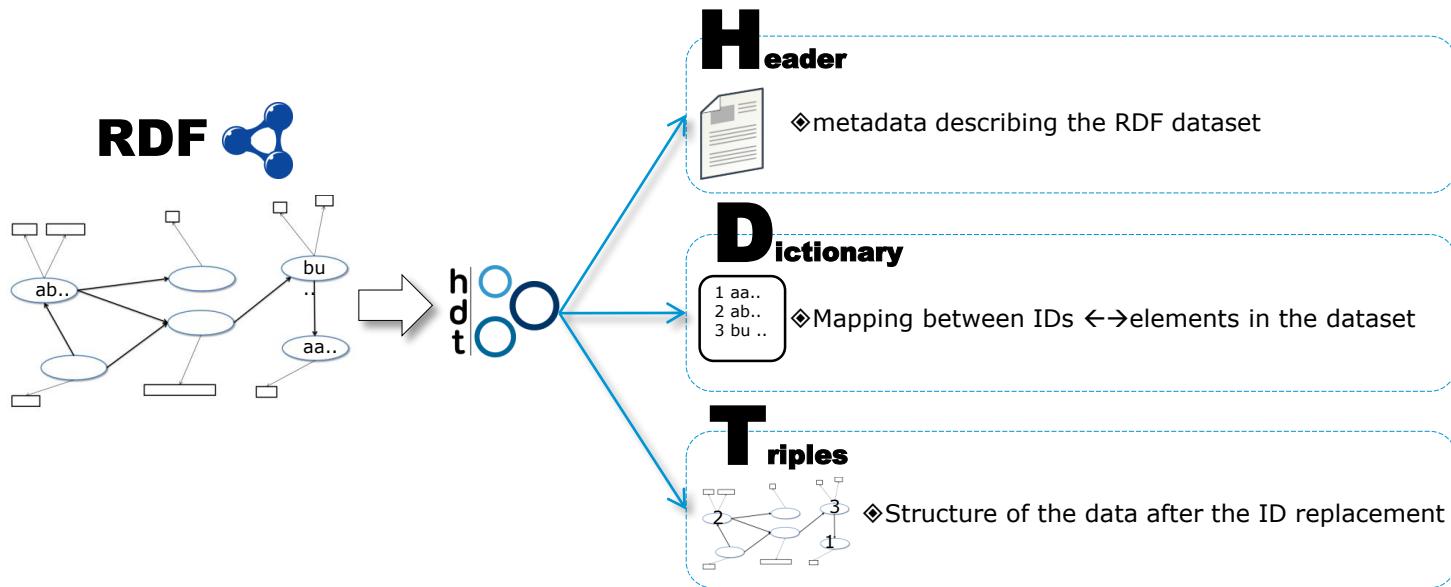


[rdfhdt.org](http://rdfhdt.org)

# Application scenarios

- Scalable storage
  - Store and serve thousands of (large) datasets (e.g. LOD Laundromat)
- Better consumer-centric publication
  - Compress and share ready-to-consume RDF datasets
  - Consumption with limited resources
    - smartphones, standard laptops
- Fast –low cost- SPARQL Query Engine
  - Via Linked Data Fragments (shown in next session)
  - Via HDT-Jena

# HDT (Header-Dictionary-Triples) Overview



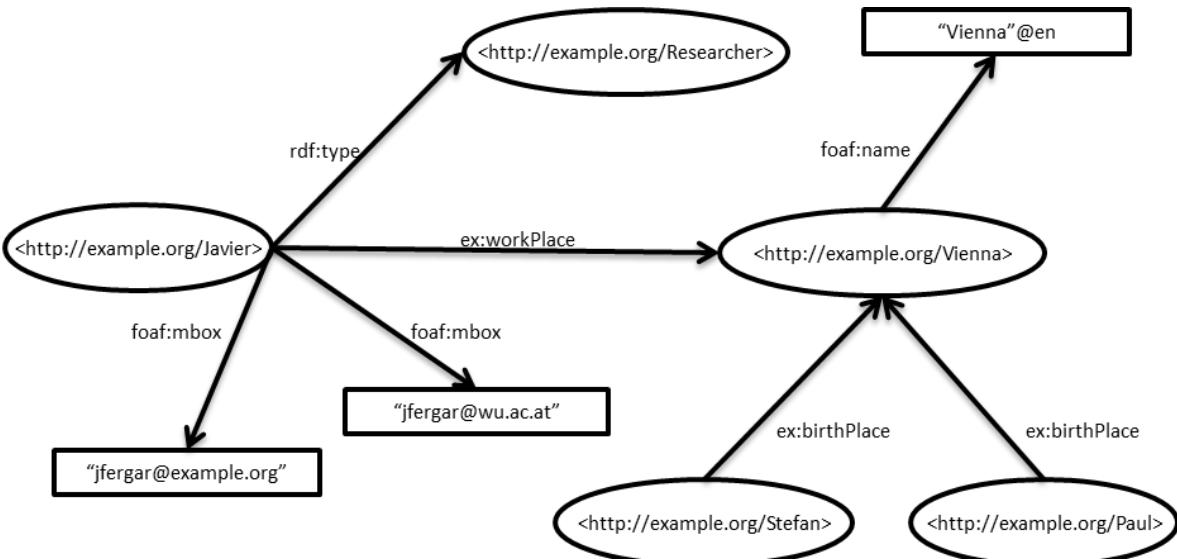
# Header

- Publication Metadata:
  - Publisher, Public Endpoint, ...
- Statistical Metadata:
  - Number of triples, subjects, entities, histograms...
- Format Metadata:
  - Description of the data structure, e.g. Triples Order.
- Additional Metadata:
  - Domain-specific.
  - ... In RDF

# Header

	Subject	Predicate	Object
1	<http://example.org>	<http://www.w3.org/1999/02/22-rdf-syntax-ns#type>	<http://purl.org/HDT/hdt#Dataset>
2	<http://example.org>	<http://www.w3.org/1999/02/22-rdf-syntax-ns#type>	<http://rdfs.org/ns/void#Dataset>
3	<http://example.org>	<http://rdfs.org/ns/void#triples>	"4789981"
4	<http://example.org>	<http://rdfs.org/ns/void#properties>	"212"
5	<http://example.org>	<http://rdfs.org/ns/void#distinctSubjects>	"505761"
6	<http://example.org>	<http://rdfs.org/ns/void#distinctObjects>	"1717095"
7	<http://example.org>	<http://purl.org/HDT/hdt#statisticalInformation>	_:statistics
8	<http://example.org>	<http://purl.org/HDT/hdt#publicationInformation>	_:publicationInformation
9	<http://example.org>	<http://purl.org/HDT/hdt#formatInformation>	_:format
10	_:format	<http://purl.org/HDT/hdt#dictionary>	_:dictionary
11	_:format	<http://purl.org/HDT/hdt#triples>	_:triples
12	_:dictionary	<http://purl.org/dc/terms/format>	<http://purl.org/HDT/hdt#dictionaryFour>
13	_:dictionary	<http://purl.org/HDT/hdt#dictionarynumSharedSubjectObject>	"383152"
14	_:dictionary	<http://purl.org/HDT/hdt#dictionarymapping>	"1"
15	_:dictionary	<http://purl.org/HDT/hdt#dictionarysizeStrings>	"22120359"
16	_:dictionary	<http://purl.org/HDT/hdt#dictionaryblockSize>	"16"
17	_:triples	<http://purl.org/dc/terms/format>	<http://purl.org/HDT/hdt#triplesBitmap>
18	_:triples	<http://purl.org/HDT/hdt#triplesnumTriples>	"4789981"
19	_:triples	<http://purl.org/HDT/hdt#triplesOrder>	"SPO"
20	_:statistics	<http://purl.org/HDT/hdt#originalSize>	"688768256"
21	_:statistics	<http://purl.org/HDT/hdt#hdtSize>	"41646252"
22	_:publicationInformation	<http://purl.org/dc/terms/issued>	"2015-09-03T03:23:05+0200"

# Dictionary+Triples partition

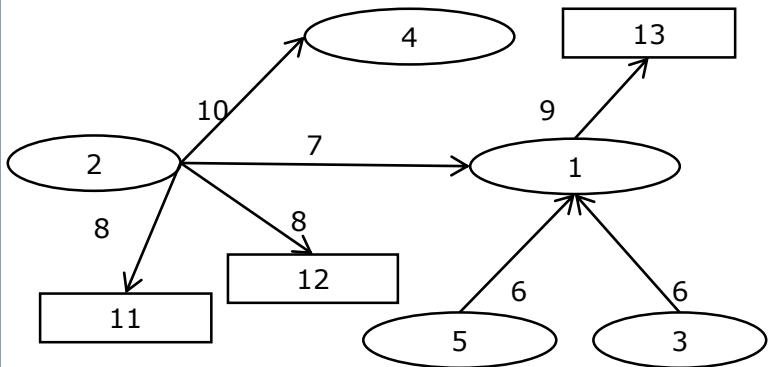


# Dictionary+Triples partition

```

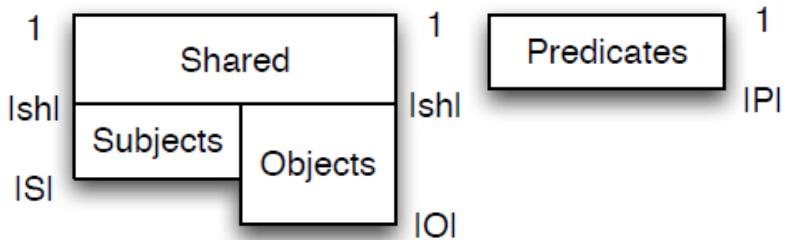
1 <http://example.org/Vienna>
2 <http://example.org/Javier>
3 <http://example.org/Paul>
4 <http://example.org/Researcher>
5 <http://example.org/Stefan>
6 ex:birthPlace
7 ex:workPlace
8 foaf:mbox
9 foaf:name
10 rdf:type
11 "jfergar@example.org"
12 "jfergar@wu.ac.at"
13 "Vienna"@en

```



# Dictionary

- Mapping of strings to correlative IDs. {1..n}
- Lexicographically sorted, no duplicates.
- Prefix-Based compression in each section.
  - Efficient ID $\leftarrow\rightarrow$ String operations



ID	Dictionary
1	<http://example.org/Vienna>
2	<http://example.org/Javier>
3	<http://example.org/Paul>
4	<http://example.org/Stefan>
2	<http://example.org/Researcher>
3	"jfergar@example.org"
4	"jfergar@wu.ac.at"
5	"Vienna"@en
1	ex:birthPlace
2	ex:workPlace
3	foaf:mbox
4	foaf:name
5	rdf:type

# Dictionary. Plain Front Coding (PFC)

- Prefix-Based compression used in each section
  1. Each string is encoded with two values
    - An integer representing the number of characters shared with the previous string
    - A sequence of characters representing the suffix that is not shared with the previous string

A  
An  
Ant  
Antivirus  
Antivirus Software  
Best

→ (0,a) (1,n) (2,t) (3,ivirus) (9, Software) (0,Best)

# Dictionary. Plain Front Coding (PFC)

- Prefix-Based compression used in each section
- 2. The vocabulary is split in buckets, each of them storing “ $b$ ” strings
  - The first string of each bucket (*header*) is coded explicitly (i.e. full string)
  - The subsequent  $b-1$  strings (*internal strings*) are coded differentially

A  
An  
Ant  
Antivirus  
Antivirus Software  
Best

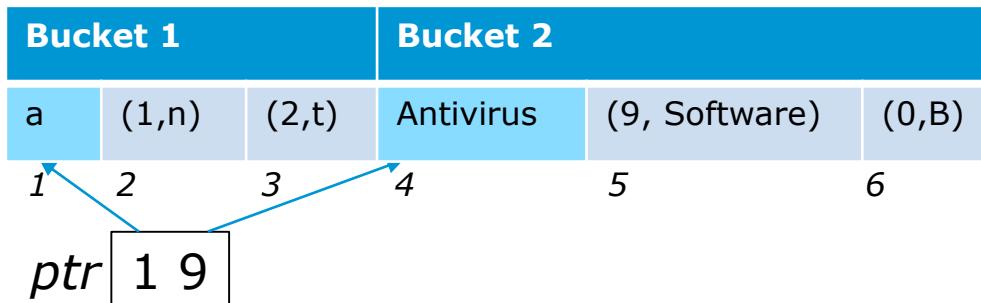


Bucket 1			Bucket 2		
1	2	3	4	5	6
a	(1,n)	(2,t)	Antivirus	(9, Software)	(0,Best)

# Dictionary. Plain Front Coding (PFC)

- Prefix-Based compression used in each section
- PFC is encoded with a byte sequence and an array of pointers (*ptr*) to denote the first byte of each bucket

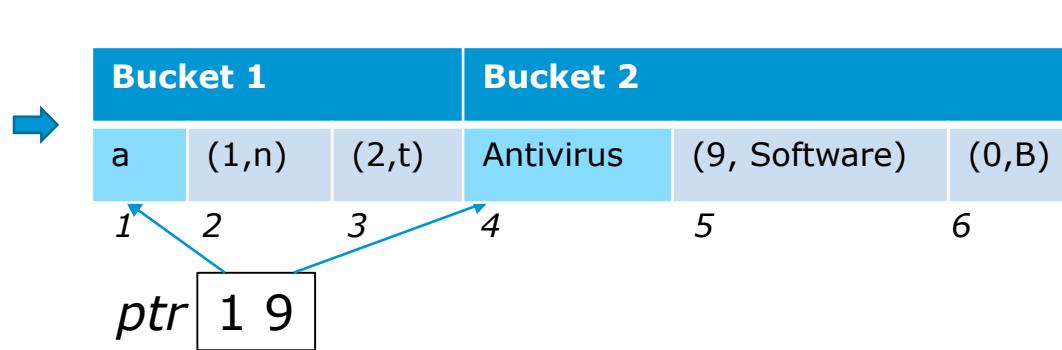
A  
An  
Ant  
Antivirus  
Antivirus Software  
Best



# Dictionary. Plain Front Coding (PFC)

- Prefix-Based compression used in each section
- Locate (string) performs a binary search in the headers + sequential decoding of internal strings
  - e.g. locate (Antivirus Software)=5
- Extract (id), finds the bucket id/b, and decodes until the given position
  - E.g. extract (5) = Antivirus Software

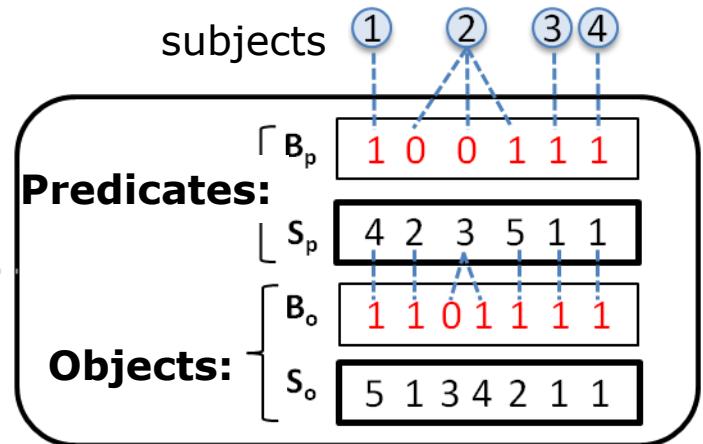
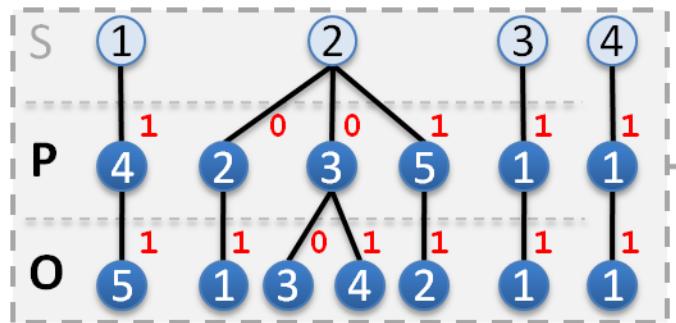
A  
An  
Ant  
Antivirus  
Antivirus Software  
Best



More on Compressed Dictionaries: Martínez-Prieto, M. A., Brisaboa, N., Cánovas, R., Claude, F., & Navarro, G. (2016). **Practical compressed string dictionaries**. *Information Systems*, 56, 73-108.

# Bitmap Triples Encoding

- Bitmap Triples:



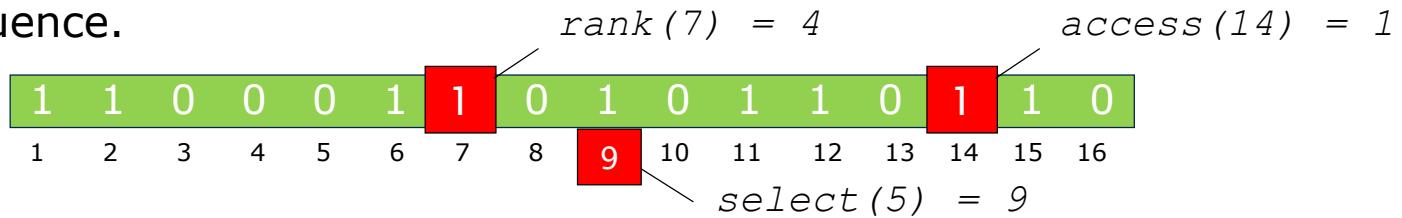
- We index the bitsequences to provide a SPO index

# Remember.... Succinct Data Structures

- Represent and index large volumes of data
- ~ Theoretical minimum space while serving efficient operations:
  - Mostly based on 3 operations:
    - Access
    - Rank
    - Select

# Bit Sequence Coding

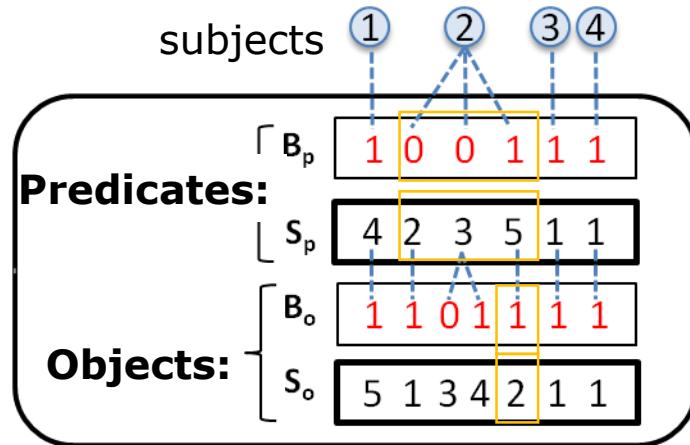
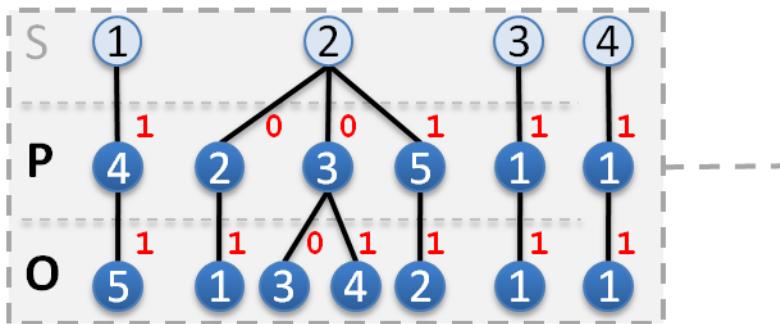
- Bitmap Sequence.



- Operations in constant time
  - $access(position) = Value$ .
  - $rank(position) = \text{"Number of ones, up to position"}$ .
  - $select(i) = \text{"Position where the one has } i \text{ occurrences"}$ .
- Implementation:
  - $n + o(n)$  bits
  - Adjustable space overhead: In practice, 37,5 % overhead

# Bitmap Triples Encoding

- **Bitmap Triples:**



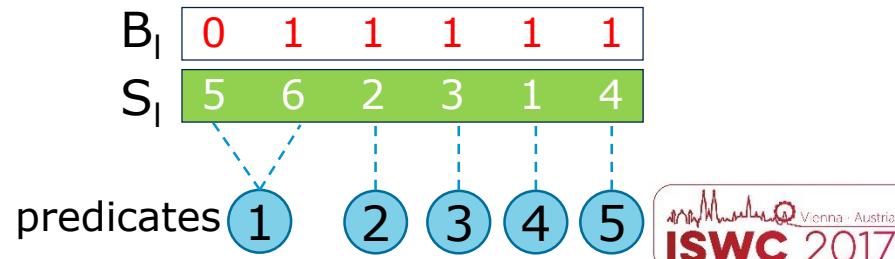
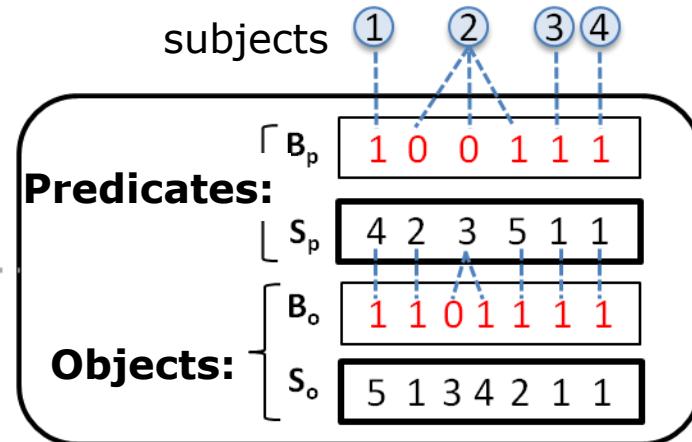
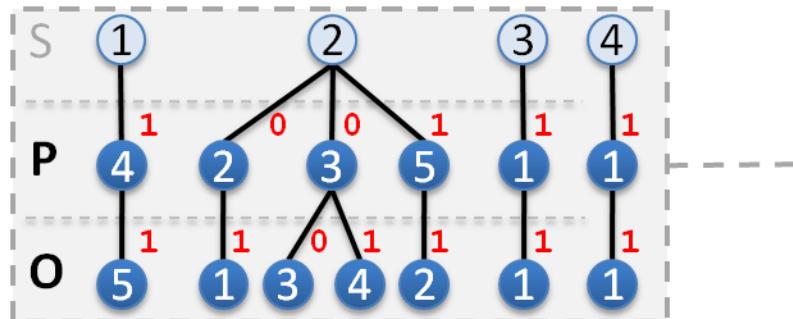
- E.g. retrieve (2,5,?)
  - Find the position of the second '1'-bit in Bp (select)
  - Binary search on the list of predicates looking for 5
  - Note that such predicate 5 is in position 4 of Sp
  - Find the position of the four '1'-bit in Bo (select)



# Additional Predicate Index

? P ?

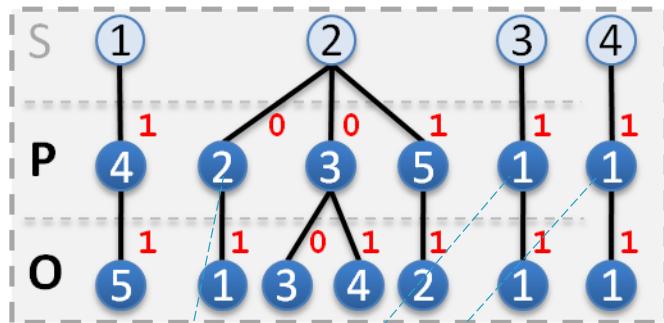
- (Historically) Wavelet Tree
    - Compact Sequence of Integers  $\{1, |P|\}$  with Operations in  $O(\log |P|)$  time
  - (Currently) Position list



# Additional Object Index

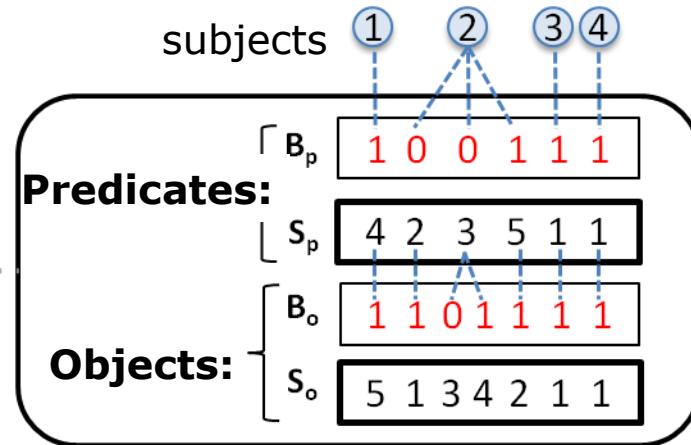
??O  
?PO

- Bitmap Triples:



B <sub>I</sub>	0	0	1	1	1	1	1
S <sub>I</sub>	2	5	6	4	3	3	1

objects



B <sub>I</sub>	0	1	1	1	1	1
S <sub>I</sub>	5	6	2	3	1	4



# On-the-fly indexes: HDT-FoQ

- From the exchanged HDT to the functional HDT-FoQ:
  - Publish and Exchange HDT
  - At the consumer:



Process	Type of Index	Patterns
1 index the bitsequences	Subject SPO	SPO, SP?, S??, S?O, ???
2 We index the position of each predicate (with a wavelet*)	Predicate PSO	?P?
3 We index the position of each object (just a position list)	Object OPS	?PO, ??O

HDT

# Achievements and Challenges



# HDT Acknowledged as W3C member submission:

- <http://www.w3.org/Submission/2011/03/>

Binary RDF Representation for Publication and Exchange (HDT)

<http://www.w3.org/Submission/2011/SUBM-HDT-20110330/>

**Member Submission**

Binary RDF Representation for Publication and Exchange (HDT)

W3C Member Submission 30 March 2011

This version:  
<http://www.w3.org/Submission/2011/SUBM-HDT-20110330/>

Latest version:  
<http://www.w3.org/Submission/HDT/>

Editor:  
Javier D. Fernández

Authors:  
Javier D. Fernández  
Miguel A. Martínez-Prieto  
Claudio Gutierrez  
Axel Polleres

Alcatel-Lucent



profi um



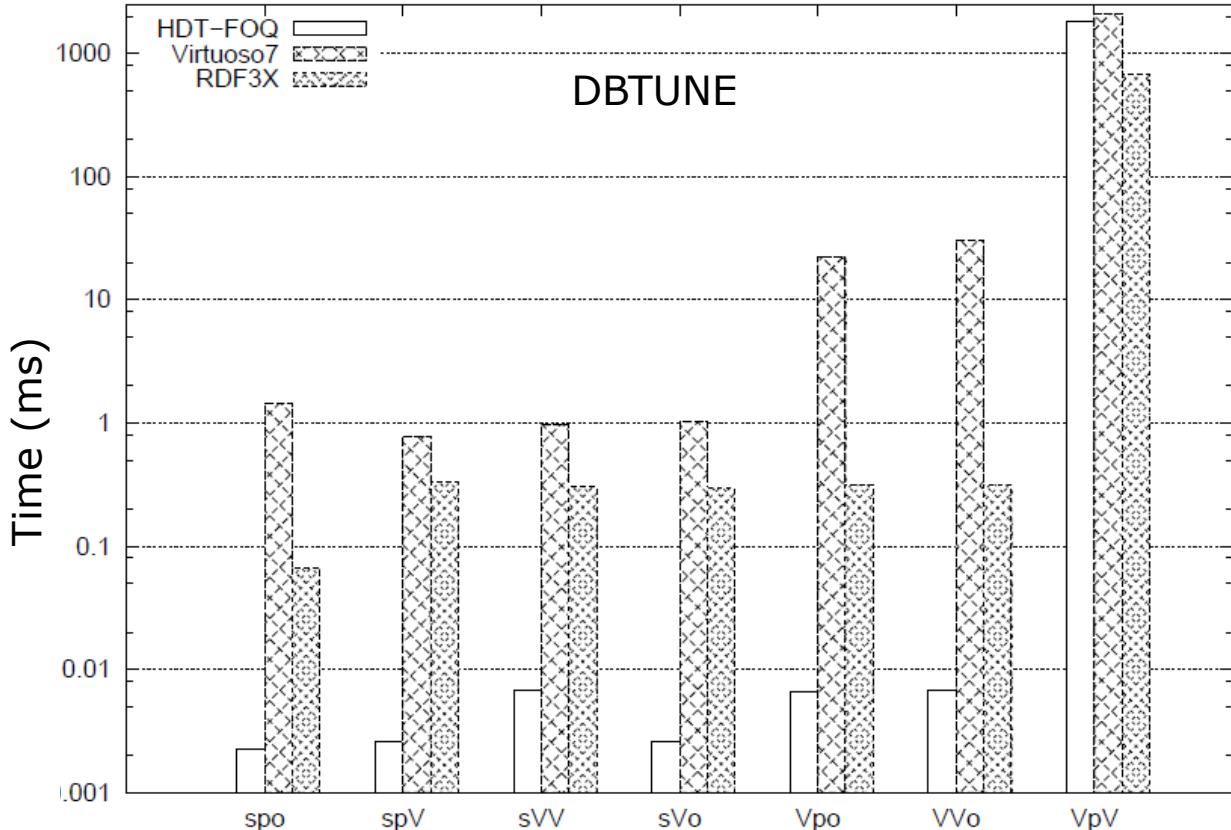
# Some numbers on size

<http://dataweb.infor.uva.es/projects/hdt-mr/>

*José M. Giménez-García, Javier D. Fernández, and Miguel A. Martínez-Prieto.* **HDT-MR:** A Scalable Solution for RDF Compression with **HDT** and **MapReduce**. In Proc. of International Semantic Web Conference (ISWC),. 2015

DATASET	TRIPLES	SO	S	O	P	NT	Size (GB)		
							NT+1zo	HDT	HDT+gz
LinkedGeoData	0.27BN	41.5M	10.4M	80.3M	18.3K	38.5	4.4	6.4	1.9
DBpedia	0.43BN	22.0M	2.8M	86.9M	58.3K	61.6	8.6	6.4	2.7
Ike	0.51BN	114.5M	0	145.1K	10	100.3	4.9	4.8	0.6
Mashup	1.22BN	178.0M	13.2M	167.2M	76.6K	200.3	18.0	17.1	4.6
LUBM-1000	0.13BN	5.0M	16.7M	11.2M	18	18.0	1.3	0.7	0.2
LUBM-2000	0.27BN	10.0M	33.5M	22.3M	18	36.2	2.7	1.5	0.5
LUBM-3000	0.40BN	14.9M	50.2M	33.5M	18	54.4	4.0	2.3	0.8
LUBM-4000	0.53BN	19.9M	67.0M	44.7M	18	72.7	5.3	3.1	1.0
LUBM-5000	0.67BN	24.9M	83.7M	55.8M	18	90.9	6.6	3.9	1.3
LUBM-6000	0.80BN	29.9M	100.5M	67.0M	18	109.1	8.0	4.7	1.6
LUBM-7000	0.93BN	34.9M	117.2M	78.2M	18	127.3	9.3	5.5	1.9
LUBM-8000	1.07BN	39.8M	134.0M	89.3M	18	145.5	10.6	6.3	2.2
LUBM-12000	1.60BN	59.8M	200.9M	133.9M	18	218.8	15.9	9.6	2.9
LUBM-16000	2.14BN	79.7M	267.8M	178.6M	18	292.4	21.2	12.8	3.8
LUBM-20000	2.67BN	99.6M	334.8M	223.2M	18	366.0	26.6	16.3	5.5
LUBM-24000	3.74BN	119.5M	401.7M	267.8M	18	439.6	31.9	19.6	6.6
LUBM-28000	3.74BN	139.5M	468.7M	312.4M	18	513.2	37.2	22.9	7.7
LUBM-32000	4.27BN	159.4M	535.7M	357.1M	18	586.8	42.5	26.1	8.8
LUBM-36000	4.81BN	179.3M	602.7M	401.8M	18	660.5	47.8	30.0	9.4
LUBM-40000	5.32BN	198.4M	666.7M	444.5M	18	730.9	52.9	33.2	10.4

# Evaluation. Querying (TP)



# Results

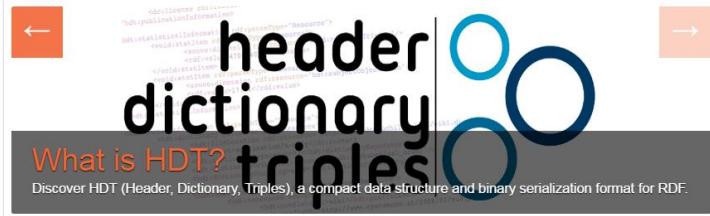
- Data is ready to be consumed 10-15x faster.
  - HDT << any other RDF format || RDF engine
- Competitive query performance.
  - Very fast on triple patterns, x 1.5 faster (Virtuoso, RDF3x).
- Integration with Jena
  - Joins on the same scale of existing solutions (Virtuoso, RDF3x).

# rdfhdt.org community

HDT - Your binary format for RDF



"HDT compresses big RDF datasets while maintaining search and browse operations" ([Read more](#))



DEVELOPERS:

- DataWebResearch**  
[http://datasets.interesník.es](#)
- DERI**
- 

**What is HDT?** Discover HDT (Header, Dictionary, Triples), a compact data structure and binary serialization format for RDF.

**GUI tool**

 Get the HDT-it GUI tool for Windows, Mac and Linux. Enjoy downloading some datasets in HDT and browse them within seconds!

[Downloads](#)

**Development and Integration**

 Keep your application data in HDT. Use our C++/Java libraries to get advantage of HDT's search performance.

[Know how](#)

**Research**

 Interested in the details? Review research publications about RDF/HDT in international venues.

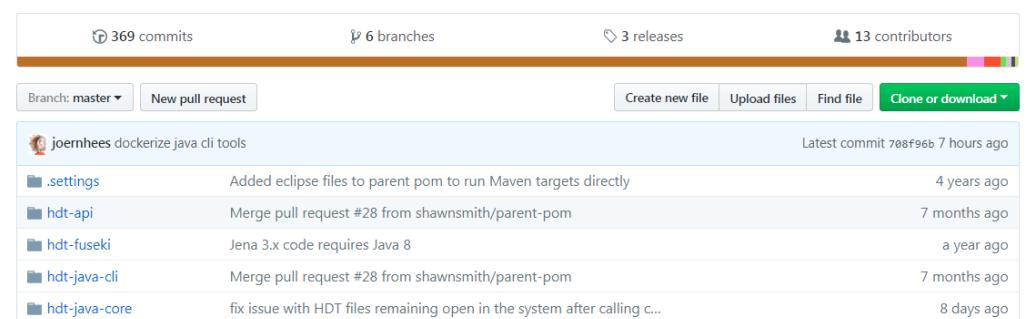
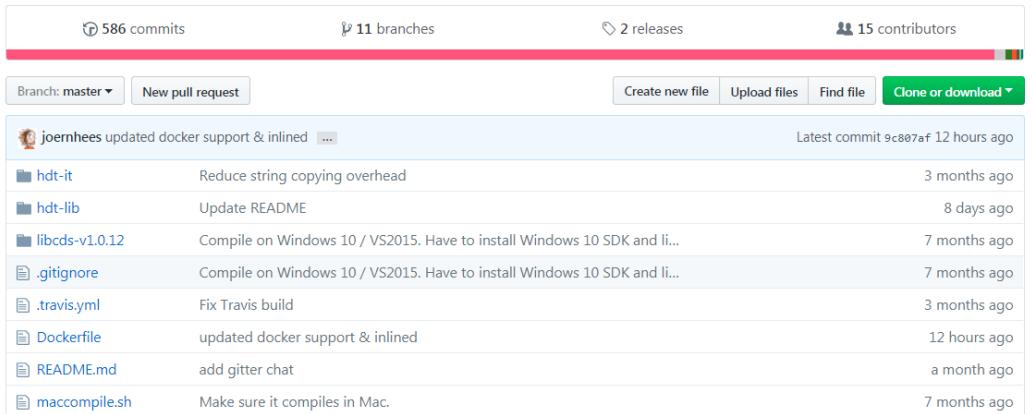
[read more](#)



RDF/HDT is a project funded by the Spanish Ministry of Economy and Competitiveness (TIN2009-14009-C02-02); Science Foundation Ireland Grant No. SFI/08/CE/I/1380; Lion-II; Chilean Fondecyt 1110267 and 1-110656.



# <https://github.com/rdfhdt>, C++ and Java tools



Only in the last two weeks...



# Get more than 650K HDT datasets from LOD Laundromat...



## LOD Wardrobe

This is where the cleaned data is stored. You can download both clean and dirty (i.e. the original) data. Each data document contains a meta-data description that includes all the stains that were detected.

Url	filter	Services				Download as			Triples	
										filter
<a href="http://csarven.ca/">http://csarven.ca/</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		54	
<a href="http://www.example.com/swc.rdf">http://www.example.com/swc.rdf</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		0	
<a href="http://rdf.freebase.com/ns/m.012_1_">http://rdf.freebase.com/ns/m.012_1_</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		16	
<a href="http://datendienst.dnb.de/cgi-bin/mabit.pl?cmd=fetch&amp;userId=openData&amp;pass=openData&amp;mabheft=ZDBTitel.ttl.gz">http://datendienst.dnb.de/cgi-bin/mabit.pl?cmd=fetch&amp;userId=openData&amp;pass=openData&amp;mabheft=ZDBTitel.ttl.gz</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		25.650.502	
<a href="http://prologmoo.com/downloads/mud.ttl">http://prologmoo.com/downloads/mud.ttl</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		4.707	
<a href="http://franz.com/franz/download/allegrograph/healthcare.nt.gz">http://franz.com/franz/download/allegrograph/healthcare.nt.gz</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		0	
<a href="http://almere.piolt.nl/sparql?default-graph-uri=&amp;query=CONSTRUCT%0D%0A%7B%0D%0A++++%3Fa+geo%3Ageometry+%3Fc+.%0D%0A%7D%0D%0AWHERE%0D%0A%7B+%3Fa+geo%3Ageometry+%3Fc+.%7D%0D%0A&amp;should-sponge=&amp;format=text%2Fturtle&amp;timeout=0&amp;debug=on">http://almere.piolt.nl/sparql?default-graph-uri=&amp;query=CONSTRUCT%0D%0A%7B%0D%0A++++%3Fa+geo%3Ageometry+%3Fc+.%0D%0A%7D%0D%0AWHERE%0D%0A%7B+%3Fa+geo%3Ageometry+%3Fc+.%7D%0D%0A&amp;should-sponge=&amp;format=text%2Fturtle&amp;timeout=0&amp;debug=on</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		10.000	
<a href="http://almere.piolt.nl/sparql?default-graph-uri=&amp;query=construct%0D%0A%7Bgraph+%3Fg-%7B%3Fs+%3Fp+%3Fo%7D%0D%0A%7Bgraph+%3Fg%0D%0A%7B%0D%0A%3Fs+%3Fp+%3Fo%0D%0A%7D%0D%0AFILTER+regex%28str%28%3Fg%29%2C%22%5Enttp%3A%2F%2Fbag_kadaster.nif%2F%22%29%0D%0A%7D+LIMIT+10000&amp;should-sponge=&amp;format=text%2Fturtle&amp;timeout=0&amp;debug=on">http://almere.piolt.nl/sparql?default-graph-uri=&amp;query=construct%0D%0A%7Bgraph+%3Fg-%7B%3Fs+%3Fp+%3Fo%7D%0D%0A%7Bgraph+%3Fg%0D%0A%7B%0D%0A%3Fs+%3Fp+%3Fo%0D%0A%7D%0D%0AFILTER+regex%28str%28%3Fg%29%2C%22%5Enttp%3A%2F%2Fbag_kadaster.nif%2F%22%29%0D%0A%7D+LIMIT+10000&amp;should-sponge=&amp;format=text%2Fturtle&amp;timeout=0&amp;debug=on</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		8.367	
<a href="http://nxp.dydra.com/nxp/public-data.ttl">http://nxp.dydra.com/nxp/public-data.ttl</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		1.630.760	
<a href="https://www.dropbox.com/s/okzyej25j2aypk/BEL2RDFexample.owl?dl=1">https://www.dropbox.com/s/okzyej25j2aypk/BEL2RDFexample.owl?dl=1</a>		<a href="#">#LD Browse</a>	<a href="#">#LD Query</a>	<a href="#">Metadata</a>	<a href="#">GZIP</a>	<a href="#">HDT</a>	<a href="#">Original Dirty File</a>		33.150	

Showing 1 to 10 of 658,045 entries (filtered from 657,902 total entries)

Previous 1 2 3 4 5 ... 65805 Next

# And query them online

**LOD Laundromat Linked Data Fragments server**

#LD  
Linked Data Fragments

4b3f9dd9aee2ad2cba453e23a8f4ae39

Query 4b3f9dd9aee2ad2cba453e23a8f4ae39 by triple pattern

subject: \_\_\_\_\_  
predicate: \_\_\_\_\_  
object: \_\_\_\_\_

[Find matching triples](#)

Matches in 4b3f9dd9aee2ad2cba453e23a8f4ae39 for { ?s ?p ?o }  
Showing triples 1 to 100 of >25,650,502 with 100 triples per page. [next](#)

```

10-3 identifier "(DE-101)010000046".
10-3 identifier "(DE-600)10-3".
10-3 identifier "(OColc)183215696".
10-3 title "Bibliografia polarografica".
10-3 isPartOf 205482-6.
10-3 issued "1949-1964".
10-3 medium 1044.
10-3 title "Bibliografia polarografica / P. 1. Elenco dei lavori e indice degli autori".
10-3 placeOfPublication "Spoleto".
10-3 succeededBy 3564-6.
10-3 type Periodical.
100-4 identifier "(DE-101)01000078X".
    
```

**Query Linked Data on the Web**  
Live in your browser, powered by Triple Pattern Fragments.

#LD  
Linked Data Fragments

Choose datasources:

http://ldf.lodlaundromat.org/4b3f9dd9aee2ad2cba453e23a8f4ae39

Directors of movies starring Brad Pitt

Type or pick a query:

```

SELECT ?movie ?title ?name
WHERE {
?movie dbpedia-owl:starring [ rdfs:label "Brad Pitt"@en ];
      rdfs:label ?title;
      dbpedia-owl:director [ rdfs:label ?name ].
FILTER LANGMATCHES(LANG(?title), "EN")
FILTER LANGMATCHES(LANG(?name), "EN")
}
    
```

[Execute query](#)

0 results in 1.6s

# Still room for optimization

- *Dynamicity:*
  - **Most compact data structures are “static”, but data may evolve**
  - **Tradeoff between compression and fast generation**
- *Advanced capabilities:*
  - **Reasoning (entailment)**
- *Query optimization:*
  - **Leverage precomputed statistics (stored in the Header)**
  - **Leverage on-demand statistics (Triple Pattern cardinality)**

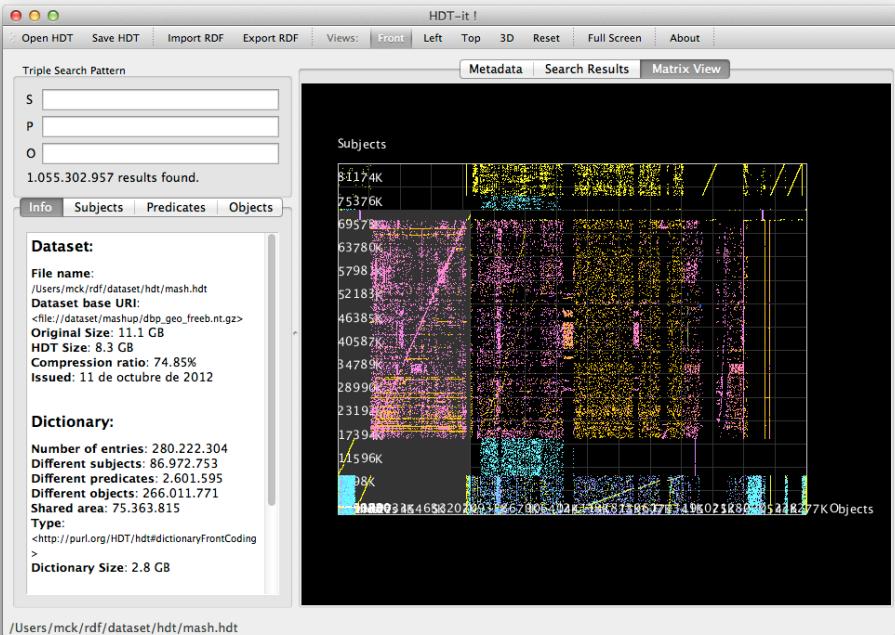
# How can I consume HDT



- GUI Tool
- C++/Java Libraries
  - Basic command line tools
  - Use them in your projects
- Set up a SPARQL Endpoint with HDT and Fuseki
- Set up a Linked Data Fragments Endpoint with HDT

# Consuming HDT

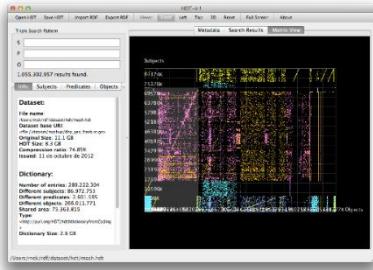
- 1) Desktop tool HDT-it!
- Thanks to Mario Arias



# Consuming HDT

- 1) Desktop tool HDT-it!
  - Download the tool for your OS:
    - <http://www.rdfhdt.org/downloads/>
  - Get an HDT dataset from the web
    - <http://www.rdfhdt.org/datasets/>
    - <http://lodlaundromat.org/wardrobe/>

Try it yourself!



OR convert your RDF dataset with the tool.

- As a suggestion of small datasets:
  - SWDF (242K triples) or the bigger DBLP (55M triples)

# Consuming HDT

- 2) C++/Java Libraries: <https://github.com/rdfhdt>
  - Command line Tools

	HDT-C++	HDT-Java
Basic command line tools	X	X
TP search	X	X
Full SPARQL	-	with Jena
Parametrizable Compression	X	-
Full text support	X	-
Practical Uses	LDF	Jena, Fuseki

# Consuming HDT

- 2) C++/Java Libraries: <https://github.com/rdfhdt>
  - Command line Tools
    - For simplicity, we will use **Java** here
    - Download hdt-java library from <https://github.com/rdfhdt/hdt-java/>
      - git clone <https://github.com/rdfhdt/hdt-java.git>
    - Install the library with maven:
      - mvn install
    - Query an HDT file:
      - Go to HDT-cli and execute:
        - ./bin/hdhSearch.sh /path/to/your/hdt
        - This will open a simple console where you can query triple patterns
    - Export/Import
      - \$> rdf2hdt file.nt output.hdt
      - \$> hdt2rdf file.hdt output.nt

Try it yourself!

# Consuming HDT

- 2) C++/Java Libraries: <https://github.com/rdfhdt>
  - Use in your projects
    - E.g. in java, include dependencies in pom.xml

Try it yourself!

```
<project xmlns="http://maven.apache.org/POM/4.0.0" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://maven.apache.org/xsd/maven-4.0.0.xsd">4.0.0example.orgtestHDT0.0.1-SNAPSHOTorg.rdfhdthdt-java-core2.0-SNAPSHOT <!-- 1.1 is the published one, 2.0-SNAPSHOT is the one we have installed locally --&gt;</dependency>org.rdfhdthdt-jena2.0-SNAPSHOT <!-- 1.1 is the published one, 2.0-SNAPSHOT is the one we have installed locally --&gt;</dependency>
```

# Consuming HDT

- 2) C++/Java Libraries: <https://github.com/rdfhdt>
- Use in your projects
  - E.g. in java, include dependencies in pom.xml
    - Example 1: create a new Class and query our HDT

```

1 package testHDT;
2
3 import java.io.IOException;
4
5 import org.rdfhdt.hdt.exceptions.NotFoundException;
6 import org.rdfhdt.hdt.HDT;
7 import org.rdfhdt.hdt.HDTManager;
8 import org.rdfhdt.hdt.triples.IteratorTripleString;
9 import org.rdfhdt.hdt.triples.TripleString;
10
11 public class FirstTest {
12
13     public static void main(String[] args) throws NotFoundException,
14     IOException {
15
16         HDT myhdt = HDTManager.mapIndexedHDT("/path/to/your.hdt", null);
17             // use loadIndexedHDT if you have enough RAM
18
19         IteratorTripleString it = myhdt.search("", "", "");
20             // any S,P,O. Use "" for a variable.
21
22         while (it.hasNext()){
23             TripleString triple = it.next();
24             System.out.println(triple);
25         }
26     }
27
28 }
```

Try it yourself!

# Consuming HDT

- 2) C++/Java Libraries: <https://github.com/rdfhdt>
  - Use in your projects
    - E.g. in java, include dependencies in pom.xml
      - Example 2: Load the model in Jena and query full SPARQL

Try it yourself!

```

18
19 public class JenaTest {
20
21     public static void main(String[] args) throws NotFoundException, IOException {
22
23         HDT myhdt = HDTManager.mapIndexedHDT("/home/javi/LearningWeek/data/swdf-2012-11-28.hdt", null);
24         // use loadIndexedHDT if you have enough RAM
25
26         // Create Jena Model on top of HDT.
27         HDTGraph graph = new HDTGraph(myhdt);
28         Model model = ModelFactory.createModelForGraph(graph);
29
30         String queryString = "SELECT ?s ?p ?o WHERE {?s <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <http://xmlns.com/foaf/0.1/Person> . "
31             + "?s ?p ?o }";
32         Query query = QueryFactory.create(queryString);
33         QueryExecution qexec = QueryExecutionFactory.create(query, model);
34         ResultSet results = qexec.execSelect();
35         for (; results.hasNext(); {
36             QuerySolution soln = results.nextSolution();
37             RDFNode s = soln.get("s"); // Get a result variable by name.
38             RDFNode p = soln.get("p"); // Get a result variable by name.
39             RDFNode o = soln.get("o"); // Get a result variable by name.
40
41             System.out.println("Solution " + s + " " + p + " " + o);
42         }
43     }
44 }
```

# Consuming HDT

- 3) Set up a SPARQL Endpoint with HDT and Fuseki
  - Go to hdt-fuseki and compile adding the dependencies:
    - mvn package dependency:copy-dependencies
  - Run fuseki
    - ./bin/hdtEndpoint.sh --hdt path/to/dataset.hdt /mydataset
  - Open your Web Browser and go to: <http://localhost:3030>
    - Select Control Panel / Dataset / myDataset and click Select
    - Type your SPARQL Query and see the results.
  - Be careful with the number of results, here there is no limitation in the number of results such as in e.g. virtuoso:
    - select \* WHERE{ ?s ?p ?o} LIMIT 400

Try it yourself!

# Consuming HDT

- 4) Set up a Linked Data Fragments Endpoint with HDT
  - Download LDF Server (Node.js is the best one but we will use java for simplicity in the installation).
    - `git clone https://github.com/LinkedDataFragments/Server.Java.git` or download <https://github.com/LinkedDataFragments/Server.Java/archive/master.zip>
    - Install the server, avoid the test (it fails :)
      - `mvn install -Dmaven.test.skip=true`
    - Open the file config-example.json and modify the settings to point to your hdt, e.g.
      - `"settings": { "file": "/home/user/myfile.hdt" }`
    - Run the server
      - `java -jar target/ldf-server.jar`
    - Access <http://localhost:8080>



# ACKs



#LD  
Linked Data Fragments



**DataWebResearch**  
<http://dataweb.infor.uva.es>

 GHENT  
UNIVERSITY

**VU** UNIVERSITY AMSTERDAM



Knowledge Representation  
and Reasoning

  
Universidad de Valladolid

 Insight

  
UNIVERSIDAD DE CHILE

 Hasso  
Plattner  
Institut  
IT Systems Engineering | Universität Potsdam