# Scalable RDF Management in the Web of Data

*Toward Efficient Interchange of RDF Data Streams*

Javier D. Fernández

Ontology Engineering Group
Departamento de Inteligencia Artificial
Facultad de Informática
Universidad Politécnica de Madrid

Index

1. Background
2. Motivation
3. What have I done so far?
4. What am I currently doing?
5. Next steps

**1.** Use URIs as names for things.

**2.** Use HTTP URIs so that people can look up those names.

**3.** When someone looks up a URI, provide useful information.

**4.** Include links to other URIs.
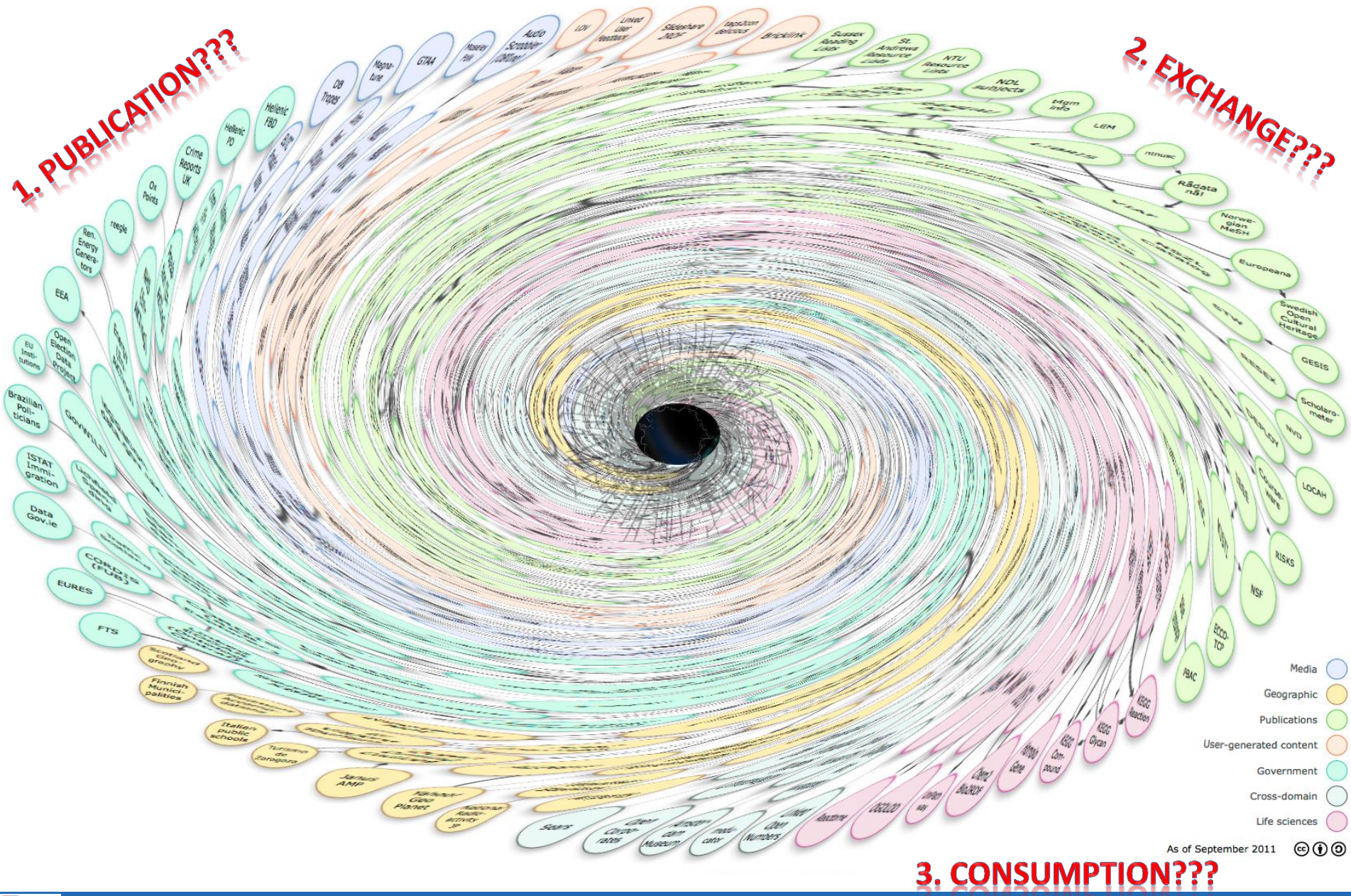
DBpedia

431 M.triples~ 63 GB

Media
Geographic
Publications
User-generated content
Government
Cross-domain
Life sciences

As of September 2011

**Linked Data cloud: > 62 billion triples**

1. PUBLICATION???

2. EXCHANGE???

3. CONSUMPTION???

Media
Geographic
Publications
User-generated content
Government
Cross-domain
Life sciences

As of September 2011

| Dataset | en | ca | de | es | eu |
|---|---|---|---|---|---|
| Mapping-based Types | nt ? | nt ? | nt ? | | |
| | nq ? | nq ? | nq ? | | |
| | ttl ? | ttl ? | ttl ? | | |
| Mapping-based Types (Heuristic) | nt ? | -- | -- | | |
| | -- | -- | -- | | |
| | ttl ? | -- | -- | | |
| Mapping-based Properties | nt ? | nt ? | nt ? | | |
| | nq ? | nq ? | nq ? | | |
| | ttl ? | ttl ? | ttl ? | | |
| Mapping-based Properties (Cleaned) | nt ? | -- | -- | | |
| | nq ? | -- | -- | | |
| | ttl ? | -- | -- | | |
| Mapping-based Properties (Specific) | nt ? | nt ? | nt ? | | |
| | nq ? | nq ? | nq ? | | |
| | ttl ? | ttl ? | ttl ? | | |
| **Dataset** | **en** | **ca** | **de** | **es** | **eu** |
| Titles | nt ? | nt ? | nt ? | nt ? | nt ? |
| | nq ? | nq ? | nq ? | nq ? | nq ? |
| | ttl ? | ttl ? | ttl ? | ttl ? | ttl ? |

431 M.triples~ 63 GB

> "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"

Claudio Gutiérrez

- Very verbose
  - Designed for human readability (not for machines)
  - HUGE → text compression/decompression

- Lack of (standard) metadata
  - "What is this?" phenomenon

- Search offline
  - Scan the whole exchanged dump.
  - (decompress)+ index the file + search

•Very verbose

(for machines)

•Lac

•"Wh

In this RDF deluge,
**if RDF is meant to be machine processable,**
why are we using plain RDF??

dex the file

Given an RDF dataset, potentially huge, a lightweight binary RDF can encode the data leveraging the skewed structure of RDF graphs for the purposes of

1. Large spatial savings,
2. easy and modular data-centric publication and parsing and
3. data retrieval.

Applications:
- Publish a large dataset on the Web.
- Transfer between two servers.
- Distributed RDF Data Management.
- Fast In-Memory Query Engine.

- Binary Serialization of RDF
- Highly compact
- Includes indexes to solve SPARQL Triple Patterns once it is loaded in main memory
- W3C Submission. http://www.w3.org/Submission/2011/03/

**RDF**

**H**eader
- metadata describing the RDF dataset

**D**ictionary
1 aa..
2 ab..
3 bu
...
- Mapping between IDs ←→elements in the dataset

**T**riples
- Structure of the data after the ID replacement

- Mapping of strings to correlative IDs. {1..n}
- Lexicographically sorted, no duplicates.
- Front Coding for each section.

Triples

1 2 6

1 3 2

2 1 3

2 2 4

2 2 5

2 4 1

3 3 2

S

P

O

| Array Y | 2 | 3 | 1 | 2 | 4 | 3 | |
|---------|---|---|---|---|---|---|---|
| Bitmap Y | 1 | 0 | 1 | 0 | 0 | 1 | |

| Array Z | 6 | 2 | 3 | 4 | 5 | 1 | 2 |
|---------|---|---|---|---|---|---|---|
| Bitmap Z | 1 | 1 | 1 | 1 | 0 | 1 | 1 |

# FOUR current ways of consuming HDT

## rdfhdt.org

- [1] Command line Tool:
  - Export/Import
    - $ rdf2hdt file.nt output.hdt
    - $ hdt2rdf file.hdt output.nt
  - Query
    - $ hdtsearch file.hdt

# FOUR current ways of consuming HDT

- [2] C++/Java Library -> Use within Jena!

```java
// Load HDT file
QueryableHDT hdt = HDTFactory.createQueryableHDT();
hdt.loadFromHDT("data/example.hdt", null);
hdt.loadOrCreateIndex(null);

// Search pattern: Empty string means "any"
IteratorTripleString it = hdt.search("", "", "");
while (it.hasNext()) {
    TripleString ts = it.next();
    System.out.println(ts);
}

// Create Jena Model on top of HDT.
HDTGraph graph = new HDTGraph(hdt);
Model model = new ModelCom(graph);
```

# FOUR current ways of consuming HDT

- [3] Web Service:
  - Import into HDT
    - http://srvgal85.deri.ie/hdt-online/
    - Thanks to Michael Hausenblas.

# FOUR current ways of consuming HDT

- [4] Desktop tool HDT-it!
  - Thanks to Mario Arias (DERI)

- Data is ready to be consumed 10-15x faster.
  - Exchange time reduced.
  - Indexing burden on server = Lightweight client processing.
- Competitive query performance.
  - Very fast on triple patterns.
  - Joins on the same scale of existing solutions.
- This is useful for applications that…
  - need a fast, compact read-only in-memory RDF store.
  - consider a static view of RDF datasets
  - want to share self-queryable RDF dumps.
  - need fast download & query.

# *Toward Efficient Interchange of RDF*
# Data Streams

- 468 stations
- 4.3 M users/day

**PH1.-** Given a set of RDF data streams, it is possible to define an RDF interchange format that optimizes the space and time for the data exchange and parsing.

**PH2.-** Given an RDF streaming engine, and a set of SPARQL queries, a RDF interchange format can be tuned to offer better performance in data exchange among processing nodes and query resolution.

**RQ1.-** Is HDT a good solution for these dynamic data?   NO
**RQ2.-** Which are the particularities of RDF data streams?
**RQ3.-** Can an RDF interchange format be parallelizable for compression and decompression (parsing)
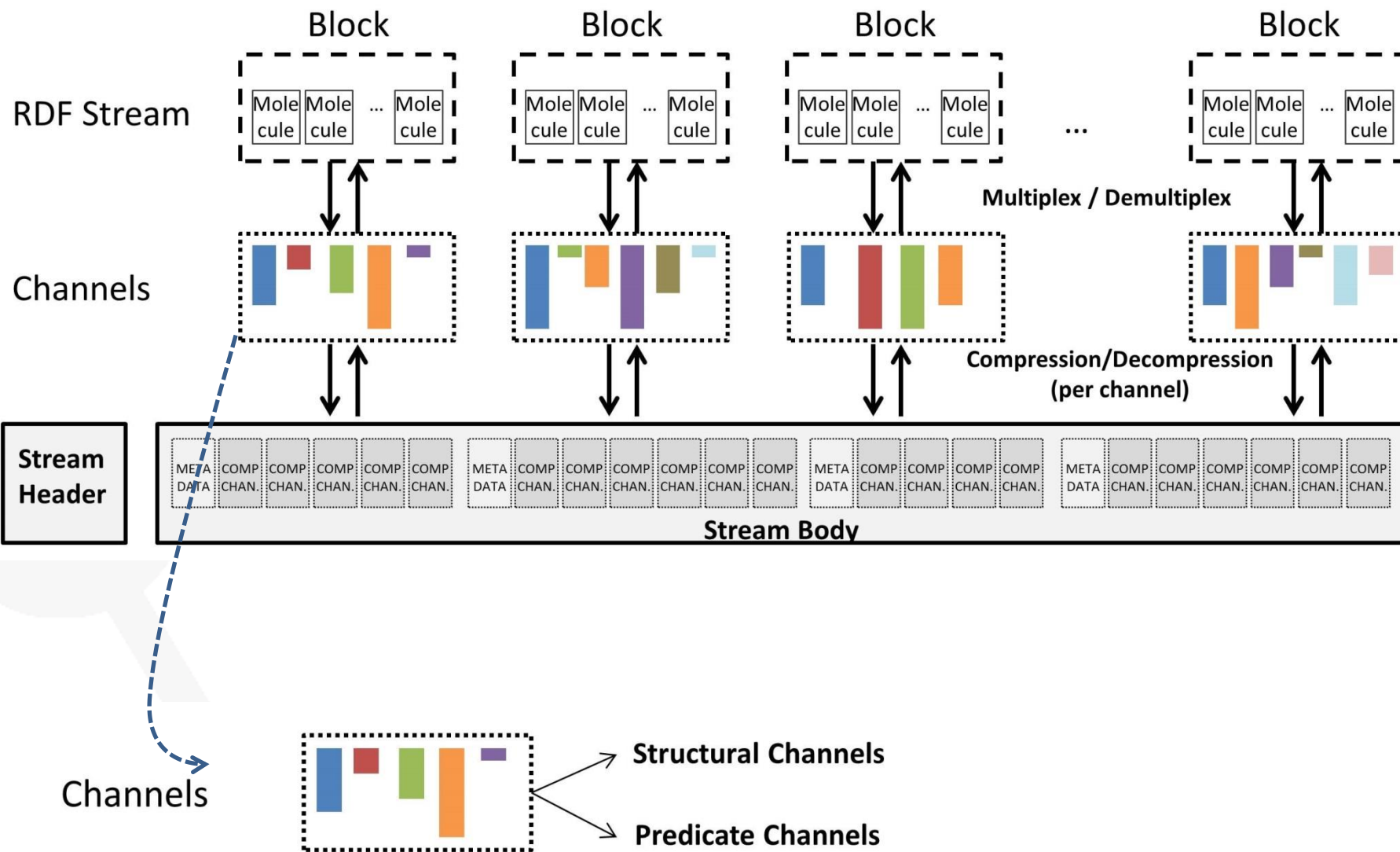
**sens-obs:Observation_AirTemperature_4UT01_2003_3_31_22_15_00**
    a    weather:TemperatureObservation ;
    om-owl:observedProperty weather:_AirTemperature ;
    om-owl:procedure sens-obs:System_4UT01 ;
    om-owl:result sens-obs:MeasureData_AirTemperature_4UT01_2003_3_31_22_15_00 ;
    om-owl:samplingTime sens-obs:Instant_2003_3_31_22_15_00 .

**sens-obs:Observation_WindGust_4UT01_2003_3_31_18_25_00**
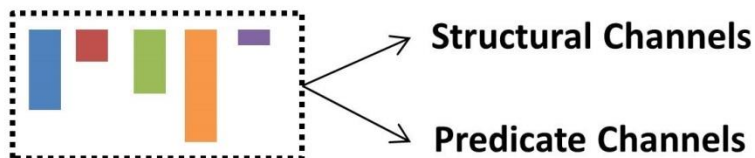    a    weather:WindSpeedObservation ;
om-owl:observedProperty weather:_WindGust ;
om-owl:procedure sens-obs:System_4UT01 ;
om-owl:result sens-obs:MeasureData_WindGust_4UT01_2003_3_31_18_25_00 ;
om-owl:samplingTime sens-obs:Instant_2003_3_31_18_25_00 .

Based on: Efficient XML Interchange (EXI) format

# Channels



→ **Structural Channels**

→ **Predicate Channels**

## Structural Channels

| Structure of the Molecule | New Structure (if needed) | New Predicate (if needed) | Main Term of the Molecule | New Term (if needed) |
|---|---|---|---|---|
| 5<br>16<br>5<br>8<br>5<br>0 | *2;4;5,0 | <http://example.org/ontology/Prop> | 8<br>1<br>0<br>8<br>3<br>0 | <http://example.org/resource/25><br><http://example.org/resource/90> |
| [IDs of Structures] | [Encoded Structures of Predicates] | [Strings] | [IDs of Terms] | [Strings] |

**Compression possibilites**

| Differential | None (negligible) | Prefix comp.<br>Zlib<br>Snappy<br>gzip | Differential | Prefix comp.<br>Zlib<br>Snappy<br>gzip |
|---|---|---|---|---|

## Predicate Channels

**{One Channel per different predicate in the related structures}**

| Predicate 2 | Predicate 4 | Predicate 5 | New Terms (if needed) | Predicate 6 |
|---|---|---|---|---|
| 15.86<br>18.78<br>19.3<br>20.5<br>24.5<br>.... | "free text…."<br>"comment…"<br>"another text…" | 101<br>245<br>0<br>284 | <http://example.org/resource/52><br><http://example.org/resource/98> | 284<br>345<br>0 |
| [Object Values]<br>[Meta: xsd:double] | [Object Values]<br>[Meta: strings] | [Term IDs]<br>[Meta: IDs] | [Strings] | [Term IDs]<br>[Meta: IDs] |

**Compression possibilites**

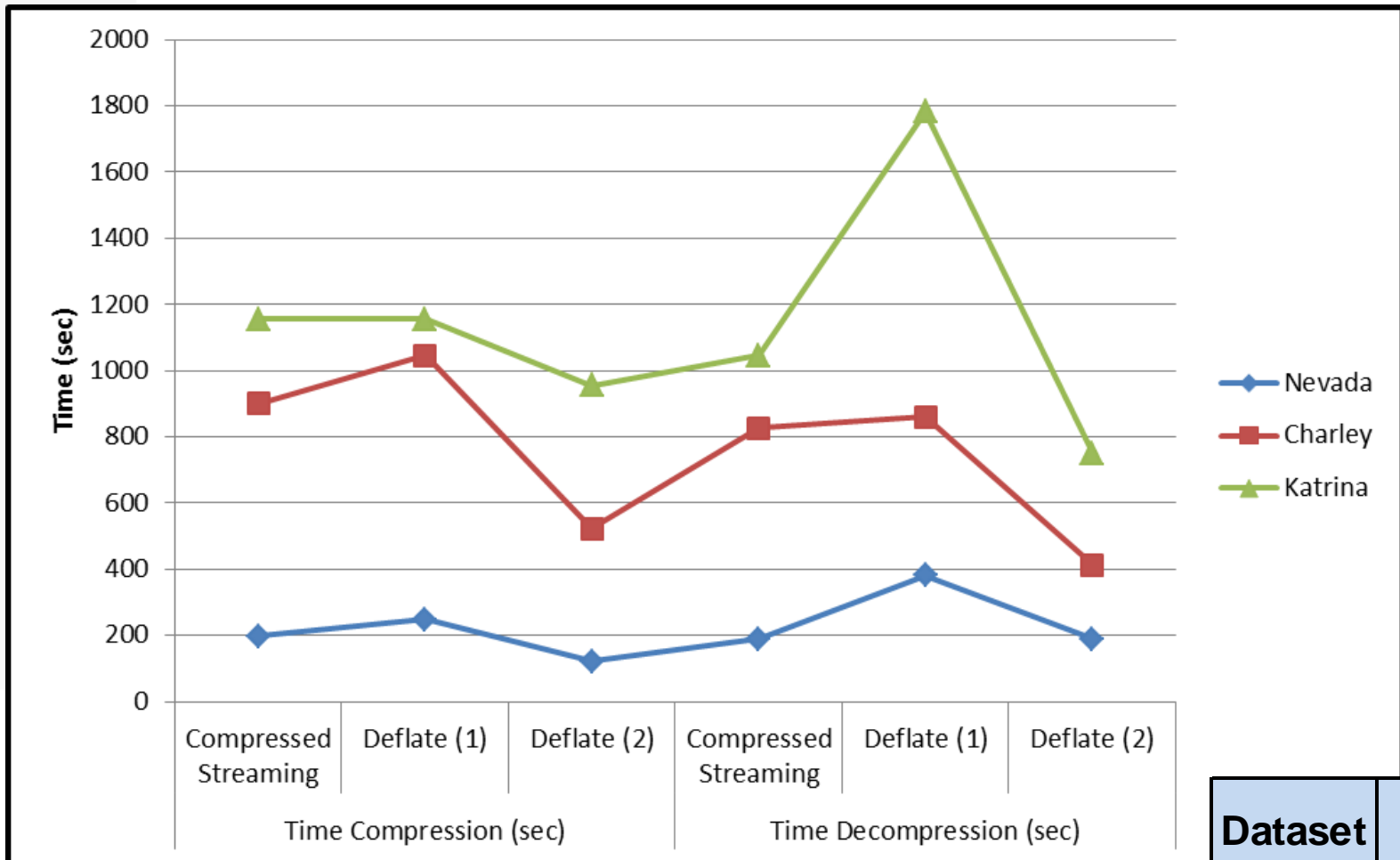| Differential | Zlib<br>Snappy<br>gzip | Differential | Prefix comp.<br>Zlib<br>Snappy<br>gzip | Differential |
|---|---|---|---|---|

Compression ratio ( New Size / Original Size)

| Dataset | Triples | Size (MB) |
|---------|---------|-----------|
| Nevada | 56,566,688 | 7,495 |
| Charley | 108,644,568 | 21,470 |
| Katrina | 179,128,407 | 35,548 |

| Dataset | Triples |
|---------|-------------|
| Nevada | 56,566,688 |
| Charley | 108,644,568 |
| Katrina | 179,128,407 |

- Finish and test the proposal with different data streams
  - Sensor and other data streams → data is welcome!
  - Release the library (Java) → feedback is welcome!
  - ISWC paper

- Parallel compression/decompression
  - preliminary proposal on Storm

- Integration within RDF streaming Engines
  - e.g. morph-streams, CQELS Cloud
  - 3 purposes:
    - scaling to higher input data rates
    - minimizing the data exchange among processing nodes
    - serving a small set of operators on the compressed data

# Scalable RDF Management in the Web of Data

*Toward Efficient Interchange of RDF Data Streams*

Javier D. Fernández

Ontology Engineering Group
Departamento de Inteligencia Artificial
Facultad de Informática
Universidad Politécnica de Madrid