# About me:

Valladolid

Pablo de la Fuente
Miguel A. Martínez-Prieto

Santiago

Claudio Gutiérrez

Madrid

Óscar Corcho
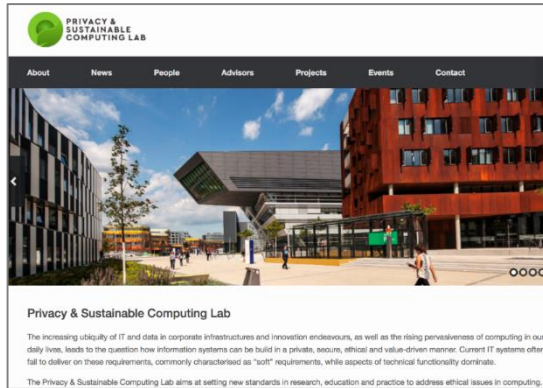
Rome

Maurizio Lenzerini

Vienna

Axel Polleres

- **Research interest:** Semantic Web, Open Data, Big (Semantic) Data Management, Databases, Data Compression, Privacy and Security
  - https://www.wu.ac.at/en/infobiz/team/fernandez/

# Where I am coming from

**Privacy & Sustainable Computing Lab**

- http://www.privacylab.at/
- Launched September 2016 with various important stakeholders: technologists, standardization, activists…
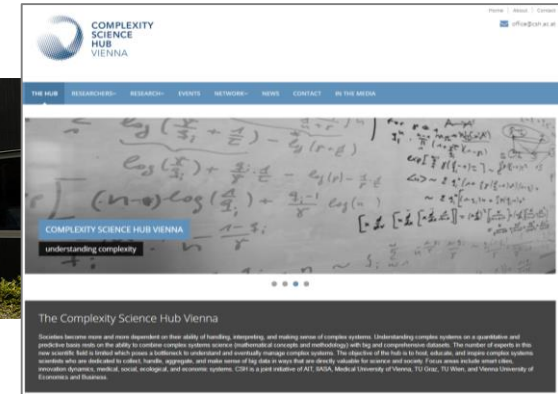- Goal: setting new standards in research, education and practice to address **ethical issues in computing**.

**Complexity Science Hub Vienna**

- http://csh.ac.at
- Launched June 2016 with Austrian stakeholders (TU, WU, Medical University of Vienna, TU Graz, AIT, IIASA)
- Goal: host, educate, and inspire complex systems scientists who are dedicated to collect, handle, aggregate, and **make sense of big data** in ways that are directly valuable for science and society.

**Prof. Axel Polleres**

Institute for Information Business

# Agenda

- **What I have done**

- **What it's in my plate**

- **Open issues**

img: zurb.com

# Agenda

- **What I have done**

- **What it's in my plate**
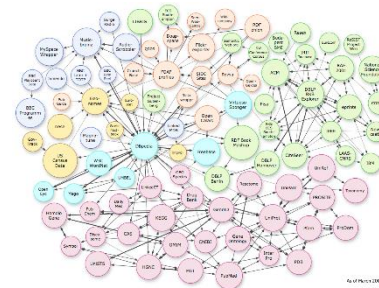
- **Open issues**

img: zurb.com

# 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
  - Hugh resources to download (large) datasets
    - + deal with the SPARQL Endpoints (zombies)
  - Messiness of the data
  - Hugh resources to index (large) datasets locally
  - Hugh resources to query (large) datasets locally and to serve them online
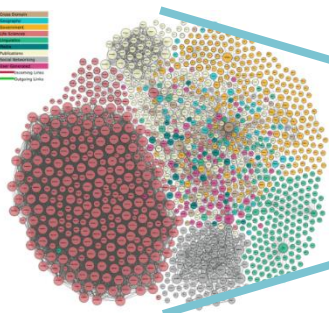
Is it much better now in 2018 ??

2009





img: Beth Scupham

# My main contribution

- Compressing and Indexing of Big Semantic Data

RDF/**HDT**

- Highly **compact serialization** of RDF (slightly more than gzip, half size of LZO)
- 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, Jena, RDF3X.
    - Main backend of Triple Pattern Fragments (**TPF**)
    - Supports FULL SPARQL as the compressed backend store of **Jena**, with an efficiency on the same scale as current more optimized solutions

LOD-a-lot

| 28 | 524 | 15.7 | 144 |
|---|---|---|---|
| Billions | GBs | GBs | seconds |
| Triples | Size | Memory Footprint | Loading Time |

http://purl.org/HDT/lod-a-lot

# SOLID architecture: Big Semantic Data in Real Time

- Based on the Lambda architecture

Martínez-Prieto, M. A., Cuesta, C. E., Arias, M., & Fernández, J. D. (2015). The solid architecture for real-time management of big semantic data. *Future Generation Computer Systems*, *47*, 62-79.

# Efficient RDF Interchange (ERI)



weather: TemperatureObservation — rdf:type

weather: AirTemperature — ssn:observedProperty

ex:CelsiusValue — ???

Light ID-31 · Humidity ID-32 · temperature ID-30 · wind ID-33

Structural Dictionary

1.- Learn patterns from the stream
2.- Sender sends the ID of the pattern and the data that differ from the pattern

- Remains efficient in performance (similar to DEFLATE)
  - Time overheads are relatively low and can be assumed in many scenarios.
- Operations on the compressed information
  - E.g. Discard all info except predicate ex:CelsiusValue

# And RDF archiving/versioning

*Querying Archives of Dynamic Linked Open Data*

FWF
Der Wissenschaftsfonds.

DIACHRON

Managing the Evolution and Preservation of the Data Web (FP7)

PRELIDA

Preserving Linked Data (FP7)

Research projects

The Dynamic Linked Data Observatory

DBpedia
WayBackMachine
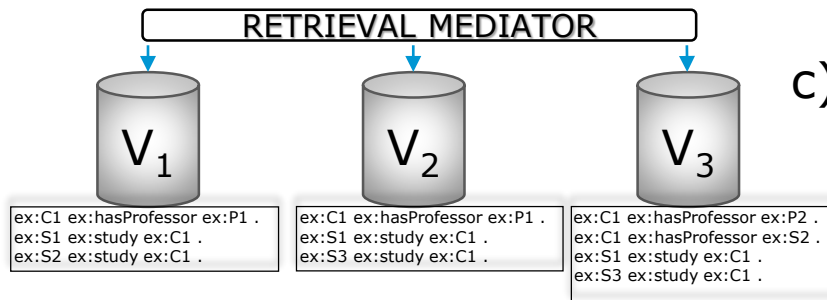
Archives

v-RDFCSA

DYDRA

Apache Marmotta

Stardog

LEDS

memento

#LD
Linked Data Fragments

Tools

HOBBIT

BEnchmark of RDF ARchives

Benchmarking

10

# RDF Archiving. Archiving policies

## a) Independent Copies/Snapshots (IC)

RETRIEVAL MEDIATOR

$V_1$

$V_2$

$V_3$

ex:C1 ex:hasProfessor ex:P1 .
ex:S1 ex:study ex:C1 .
ex:S2 ex:study ex:C1 .

ex:C1 ex:hasProfessor ex:P1 .
ex:S1 ex:study ex:C1 .
ex:S3 ex:study ex:C1 .

ex:C1 ex:hasProfessor ex:P2 .
ex:C1 ex:hasProfessor ex:S2 .
ex:S1 ex:study ex:C1 .
ex:S3 ex:study ex:C1 .

## c) Timestamp-based approach (TB)

RETRIEVAL MEDIATOR

$V_{1,2,3}$

ex:C1 ex:hasProfessor ex:P1 [$V_1$,$V_2$].
ex:C1 ex:hasProfessor ex:P2 [$V_3$].
ex:C1 ex:hasProfessor ex:S2 [$V_3$].
ex:S1 ex:study ex:C1 [$V_1$,$V_2$,$V_3$].
ex:S2 ex:study ex:C1 [$V_1$].
ex:S3 ex:study ex:C1 [$V_2$,$V_3$].

## b) Change-based approach (CB)

$\Delta_2^+$

ex:S2 ex:study ex:C1 .

$\Delta_3^+$

ex:C1 ex:hasProfessor ex:P2 .
ex:C1 ex:hasProfessor ex:S2 .

$V_1$

ex:C1 ex:hasProfessor ex:P1 .
ex:S1 ex:study ex:C1 .
ex:S2 ex:study ex:C1 .

RETRIEVAL MEDIATOR

$\Delta_2^-$

ex:S3 ex:study ex:C1 .

$\Delta_3^-$

ex:C1 ex:hasProfessor ex:P1 .

a) Independent Copies/Snapshots (IC)

c) Timestamp-based approach (TB)



| RAW DATA(GZIP) | DIFF (GZIP) | JENA TDB | | | HDT | |
|---|---|---|---|---|---|---|
| | | IC | CB | TB | IC | CB |
| 23 GB | 14 GB | 225 GB | 196 GB | 83 GB | 46 GB | 26 GB |

## Materialize (s,?,? ; version)

| IC | CB | HB$_4$ | HB8 | HB$_{16}$ |
|---|---|---|---|---|
| 48 GB | 28 GB | 34 GB | 31 GB | 29 GB |



*Hybrid approach*

# Agenda

- **What I have done**

- **What it's in my plate**

- **Open issues**

img: zurb.com

# CitySPIN project: Cyber-Physical Social Systems for City-wide Infrastructures



❖ Provide a scalable data integration framework for Cyber-Physical Social Systems (CPSSs) based on Linked Data technologies

Funding body:
- Austrian Federal Ministry of **Transport**, Innovation and Technology (BMVIT) and the Austrian **Research** Promotion Agency (FFG)

Project Duration:
- 30 months; 1.10.2017-31.3.2020

Technical coordination:
-  Marta Sabou (TU Vienna)

# What is a CPSS?



M. Z. C. Candra, H.L. Truong, "*Reliable coordination patterns in Cyber-Physical-Social Systems*," 2016 International Conference on Data and Software Engineering (ICoDSE), 2016.

ACK: Marta Sabou

# CitySPIN Use Cases

***UC Energy***: *Smart energy planning*

   **Goal**: optimize energy network and pricing
             2 M people + 230K businesses

   **How?**: understand who needs energy, when, where, how often, how
   happy they are with current services

   CitySPIN provides methods to collect and integrate customer data from:
   • Sensors
   • Internal customer legacy systems
   • Third party data: open data, social data

   … and derive customer behavioral patterns

***UC2* Mobility**: *Customer- focused Budgeting of Transport Infrastructure
Maintenance*

# CitySPIN model

# Process Discovery on Linked-Data streams

- Enriched event streams with Knowledge Graphs.



Semantic process mining: basic elements

[deMedeiros2007]

**Stock Market Events and Background Knowledge about Company Dependencies**

[Teymourian2012]

# SPECIAL (EU Horizon 2020)



https://www.specialprivacy.eu/

Taken from CNIL's twitter account

# Use Case:

An example scenario:

Sue

**Policy layer** allows Sue to revoke her consent

All her data gets **automatically** deleted from the Gym's and the BeFit's database

**Transparency ledger** reveals that the ad was sent according to Sue's consent

Scalable due to a distributed architecture based on encrypted, compressed Linked Data

SPECIAL

Sue got a second job and cannot exercise for a while

BeFit sends unwanted ad by a local Gym to Sue
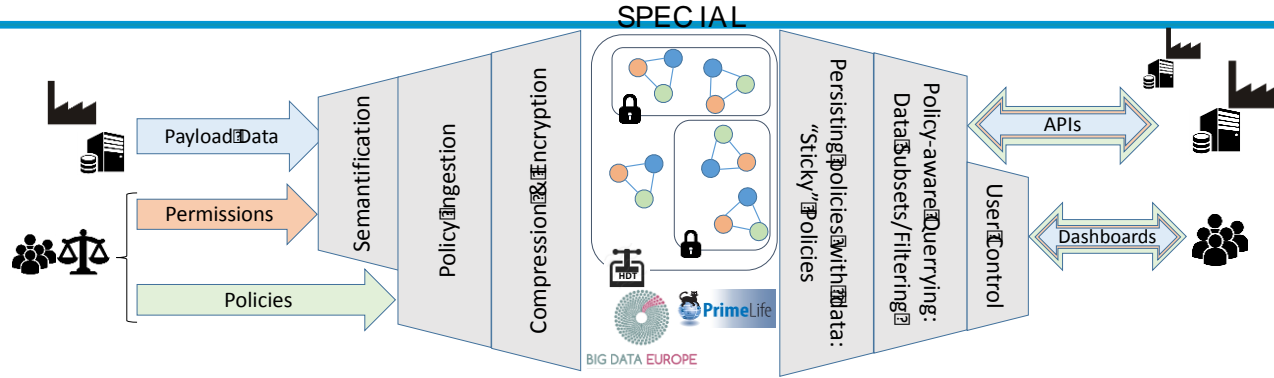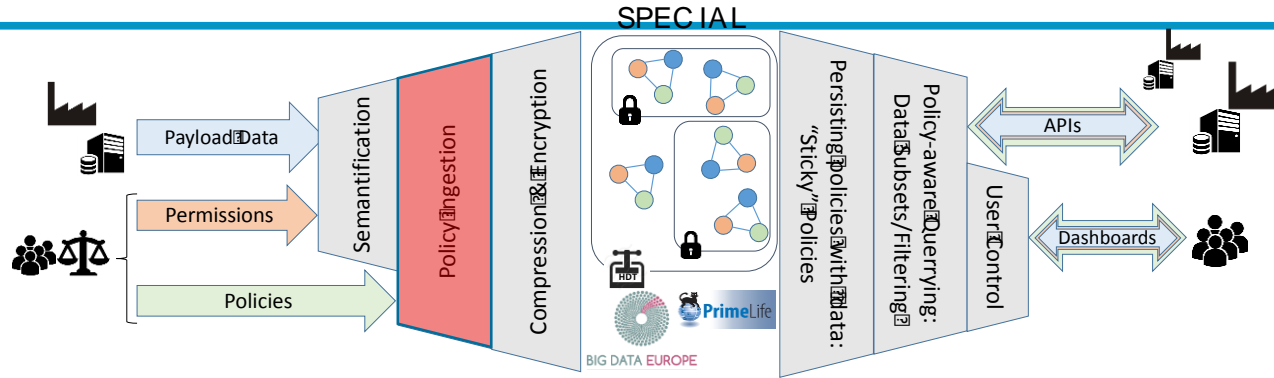
# Objectives:

- Policy management framework
  - ❖ Gives **users control** of their personal data
  - ❖ Represents and integrates **access/usage policies** and **legislative requirements** in a **machine readable format (vocabulary)**

- Transparency and compliance framework
  - ❖ Provides information on how data is **processed** and with whom it is **shared**
  - ❖ Allows data subjects to take **corrective action**

- Scalable policy-aware Big Data architecture
  - ❖ Build on top of the **Big Data Europe (BDE)** platform **scalability and elasticity mechanisms**
  - ❖ Extended BDE with **robust policy, transparency** and **compliance protocols**
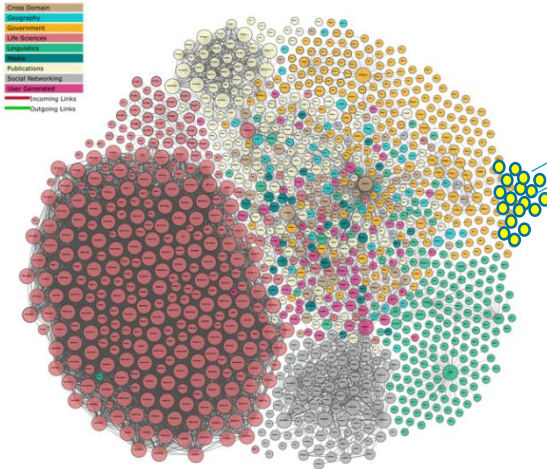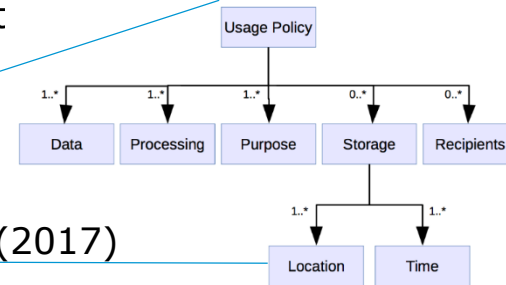
# SPECIAL Technical Components

# SPECIAL Technical Components



SPECIAL

Payload Data → Semantification → Policy ingestion → Compression & Encryption → Persisting policies with data: "Sticky" Policies → Policy-aware Querying: Data Subsets/Filtering → APIs

Permissions

Policies

User Control → Dashboards

BIG DATA EUROPE · PrimeLife · HDT

- Record **context information and access/usage** const

e.g.
W3C ODRL/POE (2017)
W3C PROV (2013)
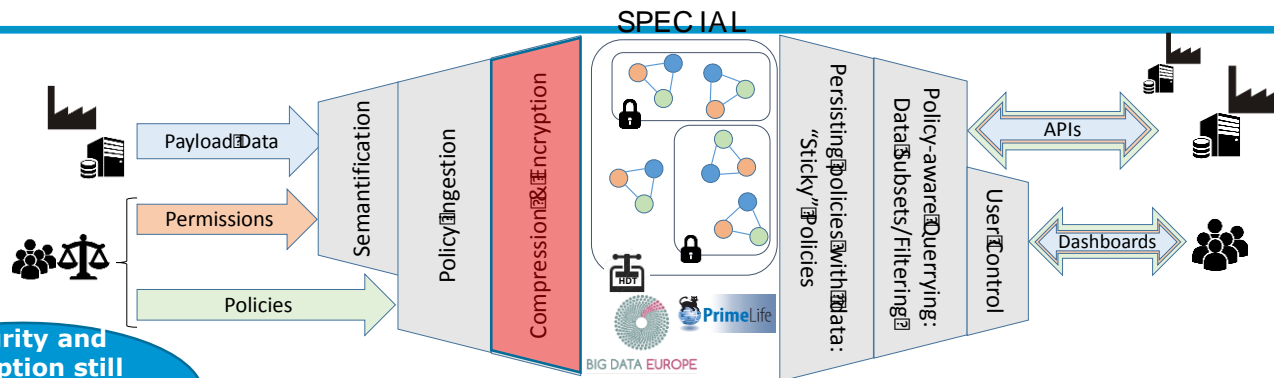Time Ontology in OWL (2017)



Usage Policy

| Data | Processing | Purpose | Storage | Recipients |

Location | Time

**Data Privacy Controls and Vocabularies**

A W3C Workshop on Privacy and Linked Data

17-18 April 2018, WU Wien, Vienna, Austria, Europe
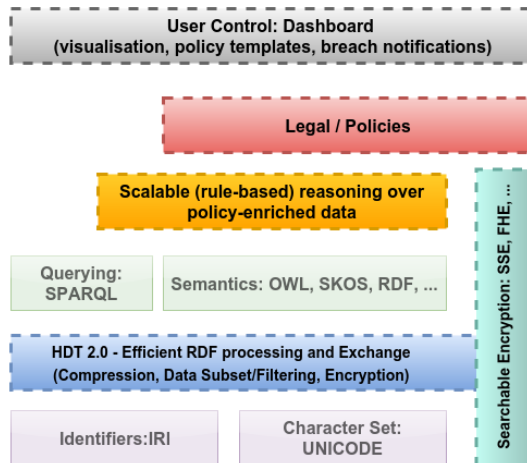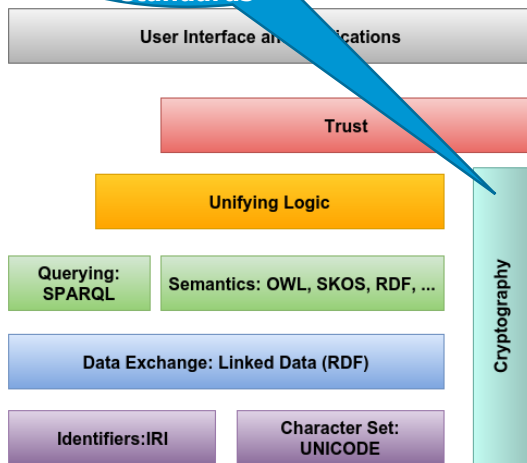
https://www.w3.org/2018/vocabws/

# SPECIAL Technical Components



Security and encryption still missing in the Linked Data standards

Storing consent, transparency records in RDF requires technology to harness RDF with:

- Queryable encryption

- Access control

- Compression (build on top of HDT)

**Self-Enforcing Access Control for Encrypted Linked Data.** Javier D. Fernández, Sabrina Kirrane, Axel Polleres, and Simon Steyskal. Extended Semantic Web Conference (ESWC 2017), May 2017

# Agenda

- **What I have done**

- **What it's in my plate**

- **Open issues**

img: zurb.com

# Expectations

democratizes the access to Big Linked Data
= Cheap, scalable consumers

*LOD-a-lot*

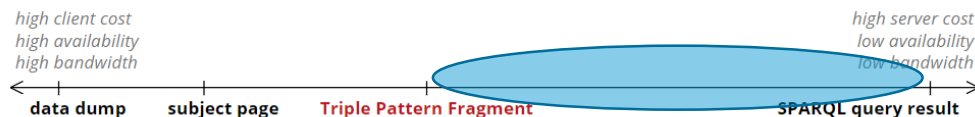#LD
Linked Data Fragments

**Reality**

# (some) Open issues

- ***"Low-cost" Backends***

  - *Compression vs. dynamicity:*

    - **Most compact data structures are "static", but data may evolve**
    - **Tradeoff between compression and fast generation**

  - *Advanced capabilities:*

    - **Reasoning (entailment)**
    - **Graph navigations (besides SPARQL)**
      - E.g. shortest path, random walk

- **Clients.** *Thin->Fat->Smart*

  - *Adaptability*
    - *E.g. Share load*
  - *Query planning (LOD-a-lot based?)*
  - *Question answering (on bigger graphs)*



high client cost
high availability
high bandwidth

high server cost
low availability
low bandwidth

data dump   subject page   Triple Pattern Fragment   SPARQL query result

**Thank you!**

*javier.fernandez@wu.ac.at*