# Ontology Virtualization for Smart Data

## A Semantics Perspective on Open Data Sharing

Krzysztof Janowicz

**STKO Lab**, University of California, Santa Barbara, USA

**RDA Metadata and Semantics Workshop, Feb. 23, 2015**

# Big Data As The New Natural Resource



http://www.ibmbigdatahub.com/infographic/big-data-new-natural-resource

# Big Data As The New Natural Resource



http://www.ibmbigdatahub.com/infographic/big-data-new-natural-resource

**A suitable analogy?**

- **Natural**, i.e., not man-made
- **Exhaustible**, finite quantity
- **Renewable**, replenishable
- **Consumed**, altered
- **Building block**

If we don't even **understand** what data are, how should we make **sense** out of them?

# The Data Retrieval Problem Is Real



Even the major data hubs such as Data.gov still rely on keyword-based search and have **unreliable, incomplete, and missing metadata**. For this type of retrieval problems, even '**a little semantics goes a long way**' (Hendler 1997).

## Sensemaking is Difficult – Fitness for Puspose is Key



- There is no shortage of data, but finding data that is **fit for a certain purpose** is difficult.
- Data as **statements** not as truth.
- **Heterogeneity** is caused by cultural differences, progress in science, viewpoints, granularity, ...
- **semantics does not come for free.**
- Lack of **provenance** information
- **Sensemaking** requires more powerful semantic technologies and ontologies (compared to IR).

## Meaningful Analysis and Synthesis is Difficult



- Ensuring that data is **analyzed and combined** in a meaningful way is far from trivial.

- **What if the information on how to use the data would come together with these data?**

- Focus on **smart data** instead of (merely on) smart applications.

- The purpose of ontologies is not to agree on the meaning of terms but to make the data provider's **intended meaning explicit**.

**A little experiment**: The statement *all rivers flow into other water bodies* is not useful because it is **'true'**[1], but because...?

---

[1] It is not; rivers can flow into the ground or just dry up entirely before reaching another water body.

# The Semantic Cube



**http://goo.gl/fBHie6**

A Misunderstanding     Semantic Web Value Proposition     Towards Ontology Virtualization
○         ○○○○○●○         ○○○○○○○○○○

The Semantic Value Proposition

## Value Proposition

**Why use Semantic Web, Linked Data, and Ontologies?**

- **Federated queries** over multiple data sources
- **Unique global identifiers** easy conflation and deduplication
- **Transparent data model**; reduces the need for guessing
- **No data silos**, no API restrictions
- Many pre-defined **lightweight vocabularies** (ontologies)
- **Smart data** reduces the need for smart applications
- **Machine reasoning** support
- **No** need for agreement, ontologies make hidden assumptions explicit
- Does away with the **data – metadata** distinction!

The Smart Data Argument

*One of the key arguments underlying the Linked Data paradigm is to **make data smart**, not applications. Instead of developing increasingly complex software, the so-called business logic should be moved to the (meta)data. The rationale is that smart data will make all future applications more usable, flexible, and robust, while **smarter applications fail to improve data** along the same dimensions.*

(**http://goo.gl/fBHie6**)

# Smart Data enables Semantic Search



Creating enriched, semantically-lifted **Linked Data** on top of Esri's ArcGIS Online

# Flexible Pattern-based Data Exploration



An user-friendly interface on top of the **DBpedia** SPARQL endpoints.

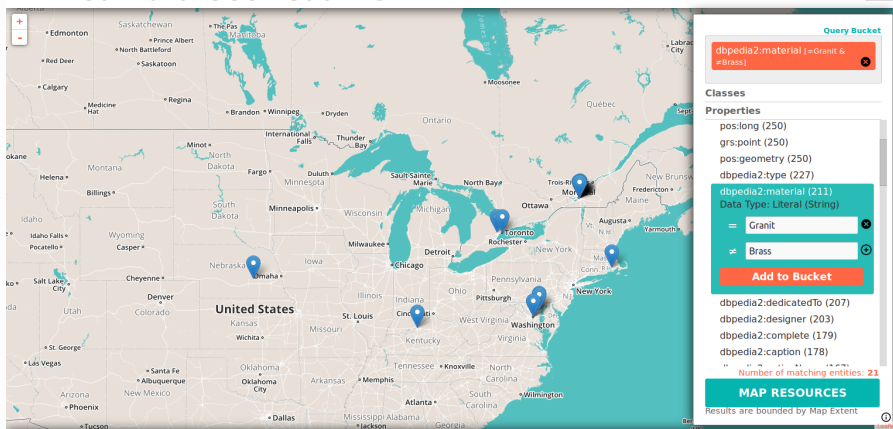# Flexible Pattern-based Data Exploration



An user-friendly interface on top of the **ADL Gazetteer** SPARQL endpoints.

## Ontology Design Pattern in a Nutshell



- **Modular** but **self-contained** building blocks
- Some patterns are **strategies**
- **Reusable** and **extendible**
- Even huge ontologies can be modularized using ODP (for example **DOLCE**)
- **No need** to import **full** ontology and **all** *ontological* **commitments**
- Different **types** of patterns, e.g. **content** vs. logical
- How **many patterns** are there?

A Misunderstanding
○

Semantic Web Value Proposition
○○○○○○

Towards Ontology Virtualization
○○○○●○○○○○

Ontology Design Patterns

# A (More Complex) Semantic Trajectory Pattern



A pattern for **discrete** trajectories of people, wildlife, vessels, and so forth.

# Ontology Design Patterns Can Be Specialized



Figure 13.2: ⟨Trajectory⟩ pattern specialised for cruises
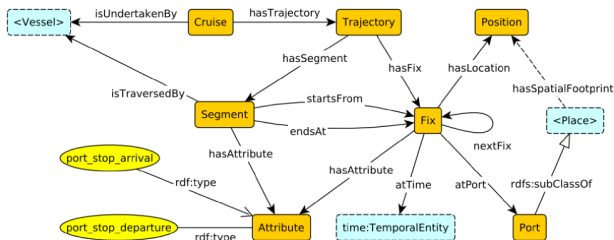
Fix ⊑ ∃hasLocation.Position ⊓ ∃atTime.time:TemporalEntity ⊓ (=1 hasFix⁻.Trajectory)
       ⊓ (⩽1 nextFix.Fix) ⊓ ¬∃nextFix.Self

    Segment ⊑ (=1 startsFrom.Fix) ⊓ (=1 endsAt.Fix) ⊓ (=1 hasSegment⁻.Trajectory)
  ∃nextFix.Fix ⊑ (=1 startsFrom⁻.Segment)
∃nextFix⁻.Fix ⊑ (=1 endsAt⁻.Segment)
    startsFrom ∘ nextFix ⊑ endsAt
        hasFix ∘ startsFrom⁻ ⊑ hasSegment

Trajectories that model the research **cruises** of scientific vessels

# A Micro-Ontology for Cruises



**Combining** the InformationObject, Event, Vessel, and Trajectory patterns

# W3C Semantic Sensor Network XG Ontology

## The Ontology Standartization Argument

*Given the early success of data format standardization, we assume that standardizing meaning (via ontologies) is less difficult and more persistent than aligning and translating local micro- ontologies. What if standardization is the more difficult task ?*

(**http://goo.gl/2e751**)

## TOWARDS ONTOLOGY VIRTUALIZATION

**In analogy to hardware virtualization**: given a set of ontology design patterns and their combination into micro-ontologies, we can abstract from the underlying axiomatization by:

- Dynamically reconfiguring patterns in a **plug&play** style
- **Bridging** between different patterns an micro-theories
- Providing ontological **views** and **semantic shortcuts** that suit particular provider, user, and use case needs by highlighting or hiding certain aspects of the underlying ontological model
- Map between major **modeling styles**,e.g., the use of instances versus classes

How do we handle different **ontological commitments?**

Quine: TO BE IS TO BE THE VALUE OF A (BOUND) VARIABLE

**Example**: Transportation is moving goods from one location to another.