Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
Field notes on working with triples

# *Operationalizing Linked Open Data*

Rob Warren[1] with much input from
S. Brown[2], A. Lemak[2], C. Faulkner[2], S. Hulan[5], C. Schwartz[3], J. Schellinck[4], D. Evans, M. Farrell[5], et al.

[1]@muninn_project warren@muninn-project.org
Adjunct, Math and Stats, Carleton Uni.
Muninn Project, Canadian Writing Research Collaboratory
[2]sbrown@uoguelph.ca - Uni. of Guelph and Uni. of Alberta
Canadian Writing Research Collaboratory
[3]Hiberdata [4]Sysabee [5]Uni. of Waterloo

Canadian Linked Data Initiative Summit 2016
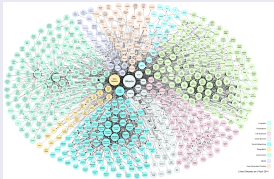https://github.com/rwarren2/cldisummit

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
Field notes on working with triples

## Who am I?

### LOD Cloud 2014



### Muninn WW1



YSA

Semantic
Quran

JITA

Er
M
N

ZDB

Muninn
World War I

Sudoc.fr

Pub
Bielefeld

Aspire
Qmul

MSC

RKB

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
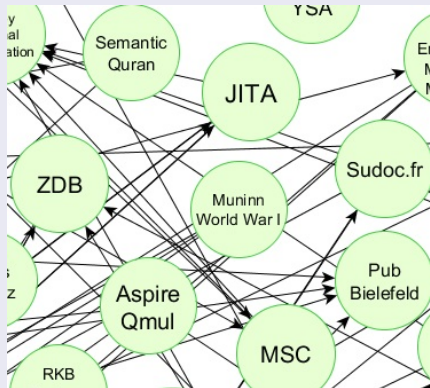Field notes on working with triples

## Who am I?

### CWRC

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
Field notes on working with triples

# First ★ ★ ★ ★ ★ data set on the Canada Open Data Portal

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
Field notes on working with triples

**Presentation Outline**

**1 Who am I?**

**2 Why Linked (Open) Data?**

**3 Field notes on vocabularies**

**4 Field notes on publishing data**

**5 Field notes on working with triples**

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
Field notes on working with triples

## The business value of LOD.

- Citations! If you can cite it, it exists!
- Externalize your costs to someone else.
- Document your data's idiosyncrasies.
- ~~Linked Data is just another fad.~~
- ~~It's already on my website.~~
- ~~People will steal my data.~~
- ~~There are errors is my data.~~

## Observations:

1. There is a bigger market for the individual pieces of your publication than the whole of it.

2. There is a bigger market for your data with people that don't share your alphabet.

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
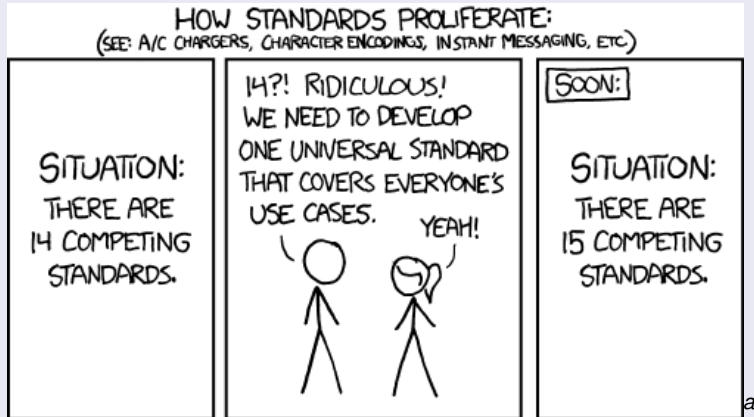Field notes on working with triples

## The propeller-head value of LOD.

- You have a machine readable URI to work with.
- You can support multiple serializations.
- You can still reference something, even if not "Open".
- You can annotate the data to the *n*th degree.
- Easy provenance and tracking of changes.
- You get multiple languages and Unicode for free.

## Observations:

1. Forces separation between the **data** and the **application**.
2. Your use cases for the data are never what people want out of the **application**.
3. *LOD engages with people by engaging their machines.*

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
Field notes on working with triples

## Vocabularies: Use a standard. (Which one!?)



[a] https://xkcd.com/927/

Who am I?
Why Linked (Open) Data?
**Field notes on vocabularies**
Field notes on publishing data
Field notes on working with triples

## Vocabulary use options:

① Create your own.

② Use one existing vocabulary.

③ Use multiple existing vocabularies.
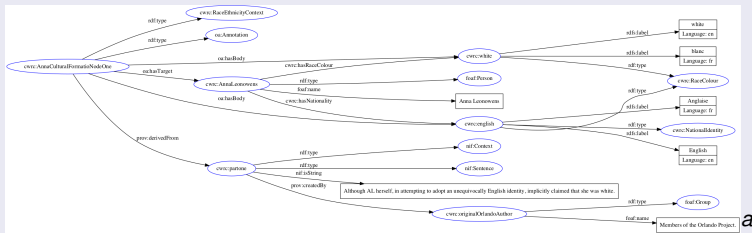
## The data consumer's perspective:

- Consumers want to know what to expect in vocabularies.
- Multiple vocabularies need relationships. (You build them).
- **The vast majority of data consumers cannot use ontology reasoning at query time.**

Who am I?
Why Linked (Open) Data?
**Field notes on vocabularies**
Field notes on publishing data
Field notes on working with triples

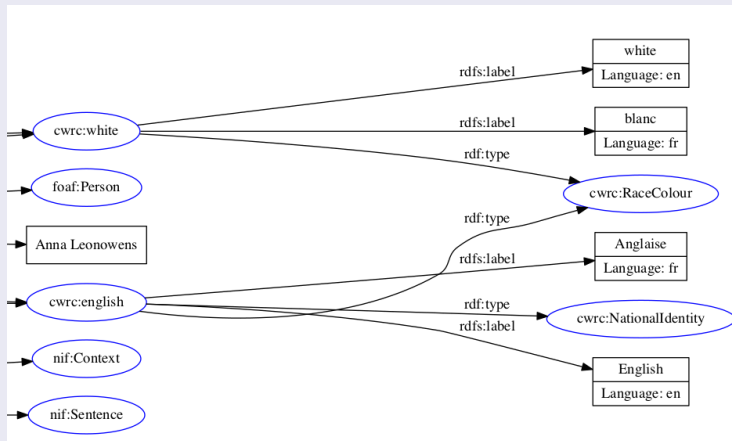**Case Studies:**

### Overview: CWRC (`http://www.cwrc.ca/`)

- Primarily Orlando TEI-style data.
- Schema definitions not ontologically sound.
- Custom ontology linked to other ontologies.
- Questions of ethnicity, race, skin colour alternate between vernacular and technical.

Who am I?
Why Linked (Open) Data?
**Field notes on vocabularies**
Field notes on publishing data
Field notes on working with triples

## Case Studies: CWRC Entry



[a] Anna LeonOwens

Who am I?
Why Linked (Open) Data?
**Field notes on vocabularies**
Field notes on publishing data
Field notes on working with triples

## Case Studies: CWRC Entry



[a]

[a] Anna LeonOwens

Who am I?
Why Linked (Open) Data?
**Field notes on vocabularies**
Field notes on publishing data
Field notes on working with triples

### Outcomes

- The Ontology is a data explanation tool. Initially (wrongly) seen as a controlled vocabulary.
- Much time is being spent on teasing out the intent of the data as written.
- The process is very demanding of the scholars.
- The CWRC ontology in its final form will have paradoxes. Acceptable because it explains data that was not built within an *ontologically rational* framework.
- This is good enough for partial data exchange.
- Massive ancillary linkages to other dataset.

Who am I?
Why Linked (Open) Data?
**Field notes on vocabularies**
Field notes on publishing data
Field notes on working with triples

**Case Studies:**

### Overview: Muninn (`http://rdf.muninn-project.org/`)

- Heterogeneous data sources: text, SQL, images, free form tabular.
- Erroneous, ambiguous and incomplete data.
- Multiple purpose built ontologies for specialized applications.
- Move to standardized ontologies as they become available. (re: Organization Ontology)
- No "single" truth, but you are free to decide for yourself.

Who am I?
Why Linked (Open) Data?
**Field notes on vocabularies**
Field notes on publishing data
Field notes on working with triples

## Private Peat, by Harold R. Peat

*I was sharing a box with a lad whom I heard the fellows call* **Bob**.

*"You're in the right direction-don't turn round!"*

## Private Peat



## Partial Information

$<$owl:oneOf rdf:parseType="Collection"$>$
$<$owl:Thing rdf:about="Bob #1"/$>$
$<$owl:Thing rdf:about="Bob #2"/$>$
$<$owl:Thing rdf:about="Bob #3"/$>$ ...
$<$/owl:oneOf$>$
$<$rel:knowsByReputation
rdf:resource="The Mad Major"/$>$

Who am I?
Why Linked (Open) Data?
**Field notes on vocabularies**
Field notes on publishing data
Field notes on working with triples

## Attestation Papers

DOB 1893-02-31 - February 31, 1893

## Partial Information

<owl:time rdf:about="Birth">
<time:hasDateTimeDescription>
<time:DateTimeDescription ...>
<time:year>1893</time:year>
<time:DateTimeDescription>
</time:hasDateTimeDescription>
<rdf:value>1893-02-31</rdf:value>
</owl:time>

## Harry Baird

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
Field notes on working with triples

## Case Studies Muninn:

## British Trench Map Coordinate Translation App

Who am I?
Why Linked (Open) Data?
**Field notes on vocabularies**
Field notes on publishing data
Field notes on working with triples

## Field notes on Vocabularies: Conclusions

1. **The public interacts with Applications not Data, but Data is why we are here.**
2. Do not ever design vocabulary for the application.
3. Old data is never clean, sensical or well behaved. The ontology / vocabulary has to say so and work with it.
4. Reuse vocabularies and create new ones on a case by case basis.
5. Great resource at
   https://lov.okfn.org/dataset/lov

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
**Field notes on publishing data**
Field notes on working with triples

**Publishing Linked Data:**

## Checklist:

- Dereferencable (URI's for everything)?
- Content negotiation (**The data format is dead.**)?
- Public facing SPARQL server?
- Machine and Human readable vocabulary definition?
- Machine and Human readable data set definition?
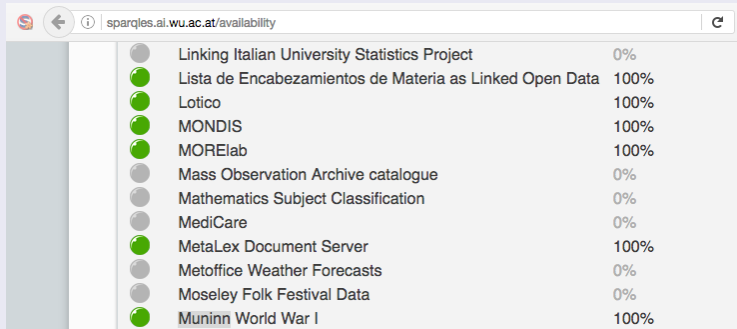- Production, in-house use of the SPARQL on day 1?

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
**Field notes on publishing data**
Field notes on working with triples

### Important Note: People write bad programs.

"If builders built buildings the way programmers make programs, the first woodpecker to come along would destroy civilization." - Gerald Weinberg

### Corollary:

Get someone who knows public facing infrastructure to look things over for you.

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
**Field notes on publishing data**
Field notes on working with triples

## SPARQL servers



SPARQL allows for custom retrieval queries over HTTP without having you involved.

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
**Field notes on publishing data**
Field notes on working with triples

### An important note about SPARQL

Run SPARQL queries through a reverse HTTP proxy: ngix,
polipo, etc.

### Why?

- Offending programmers can be safely ignored.
- Allows for light infrastructure abuse (auto-complete queries).
- Improves performance without heavy planning.

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
**Field notes on working with triples**

### Tracking data in large data stores:

- Generate more data as a byproduct of operations:
  It is easier to delete old triples than to rebuild triples that
  should have existed.

- Tracking provenance of node is trivial; consider building it
  into your work flow.

- Data and meta-data are merging.

- The most awesome use of your data is a use case you
  have not thought of.

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
Field notes on working with triples

## Dealing with contentious issues (1/2):



muninn-ww1:Military/Trench/5712bc467a2a3cf2e154b304adb4cc2f
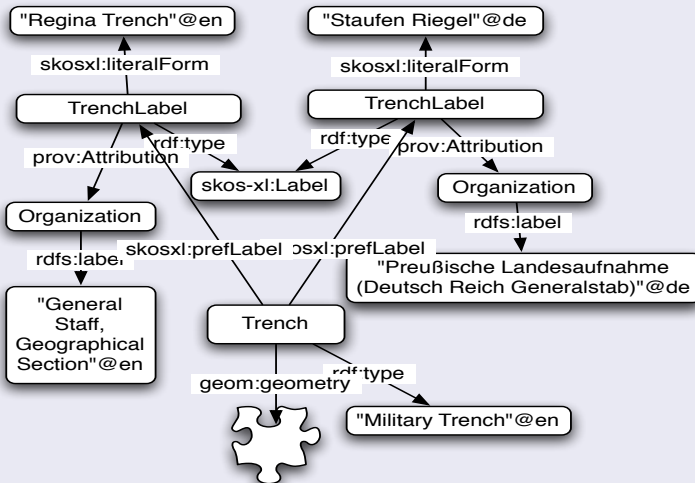  ⟶ rdf:type ⟶ mil:MilitaryTrench, time:TemporalEntity, http://geovocab.org/spatial#Feature
  ⟶ rdfs:label ⟶ "German held trench, Regina, Grandcourt Area"@en
  ⟶ owl:sameAs ⟶ dbpedia:Regina_Trench
  ⟶ time:hasDateTimeDescription ⟶ muninn-
ww1:DateTimeDescription/f48c39552b0c7d810f5a59ea7fb9f2de
  ⟶ foaf:name ⟶ "Regina"@en
  ⟶ prov:wasGeneratedBy ⟶ muninn-ww1:Process/ReginaTrenchExtraction
  ⟶ prov:hadPrimarySource ⟶ muninn-ww1:map/f48c39552b0c7d810f5a59ea7fb9f2de
  ⟶ void:inDataset ⟶ muninn-ww1:Dataset/ReginaTrench
  ⟶ geom:geometry ⟶ muninn-ww1:Military/Geometry/5712bc467a2a3cf2e154b304adb4cc2f
  ⟶ http://www.w3.org/2008/05/skos-xl#prefLabel ⟶ muninn-
ww1:AltLabel/5712bc467a2a3cf2e154b304adb4cc2f, muninn-
ww1:PrefLabel/5712bc467a2a3cf2e154b304gdb4cc2f

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
**Field notes on working with triples**

## Dealing with contentious issues(2/2):

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
**Field notes on working with triples**

## Important ontological note:

The *thing* and the *name of the thing* are not the same *thing*.

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
**Field notes on working with triples**

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
**Field notes on working with triples**

## Getting value out of low-value items:

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
**Field notes on working with triples**

## Print your own Battlefield

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
Field notes on working with triples

## Recap:

- Linked Open Data is about data, not applications.
- The *thing* and the *name of the thing* are not the same *thing*.
- The most awesome use of your data is a use case you have not thought of.
- Vocabulary use means something.
- LOD engages with people by engaging their machines.

## Further information

- http://www.cwrc.ca/
- http://www.muninn-project.org/
- https://www.youtube.com/watch?v=aJW16qFkGHU

Who am I?
Why Linked (Open) Data?
Field notes on vocabularies
Field notes on publishing data
**Field notes on working with triples**

# Questions?



CWRC / CSÉC

UNIVERSITY OF
ALBERTA

linked
modernisms

Canada
Foundation for
Innovation

UNIVERSITY
of GUELPH

Social Sciences and Humanities
Research Council of Canada
Conseil de recherches en
sciences humaines du Canada
Canadá

compute | calcul
canada | canada