

# Reusing Linguistic Resources

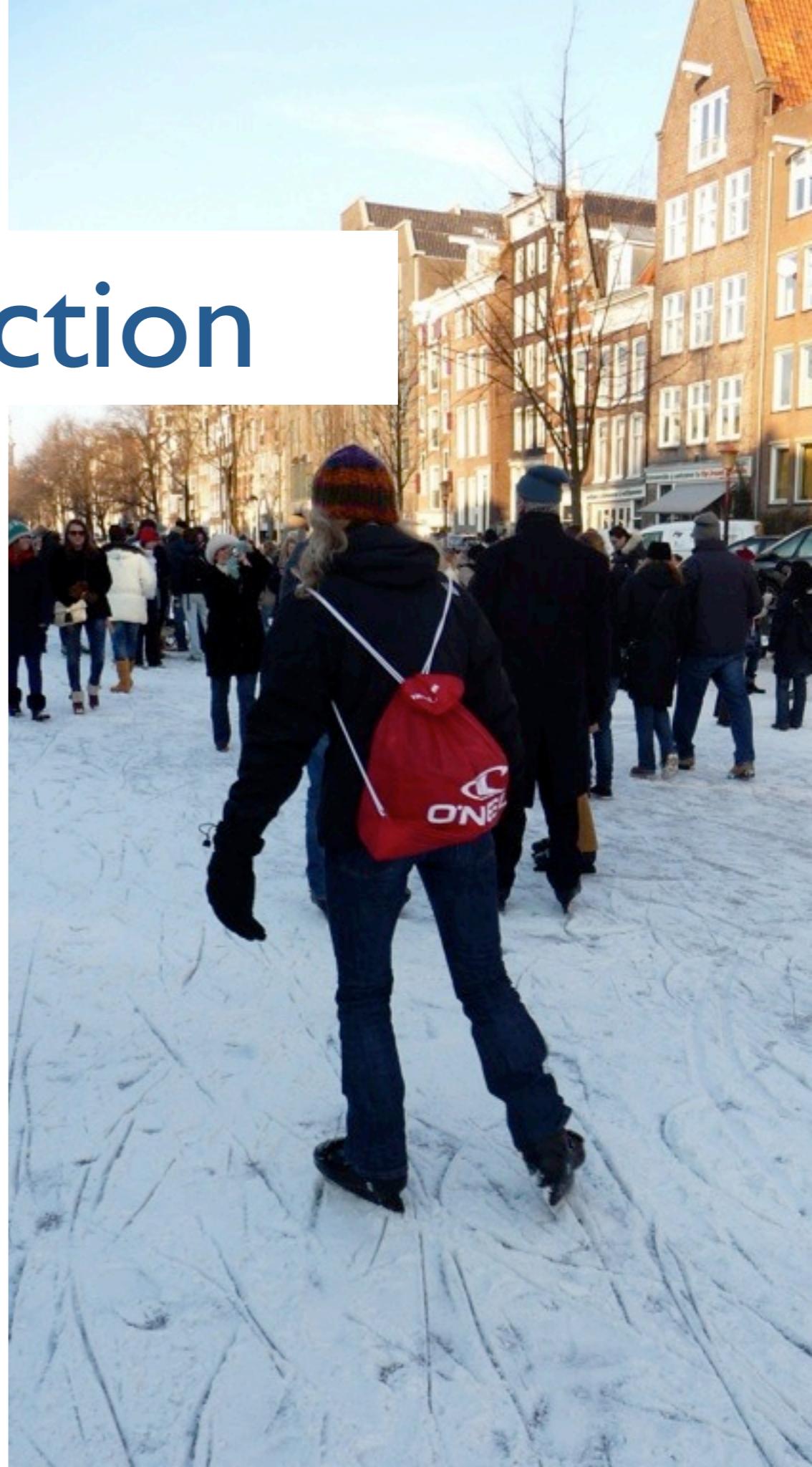
## Tasks and Goals for a Linked Data Approach

Marieke van Erp  
[marieke@cs.vu.nl](mailto:marieke@cs.vu.nl)



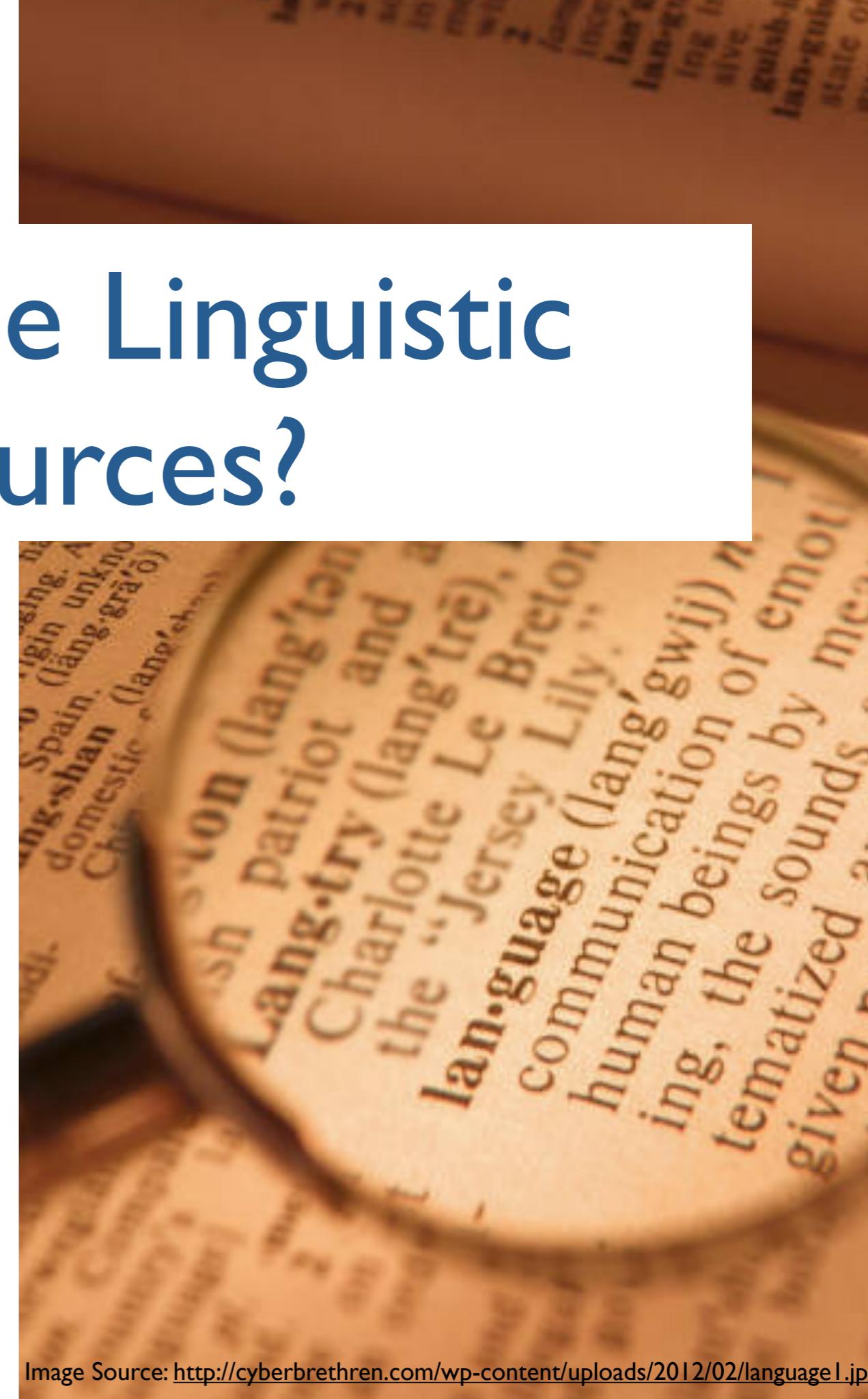
# Introduction

- BA, MA & PhD coupling/  
information extraction  
**@Tilburg University**
- Since 2009: SemWeb group  
**@VU University Amsterdam**



# Why Reuse Linguistic Resources?

- Linguistic resources are expensive to create
- ...and difficult to use for ‘outsiders’
- How can we reach out to the ‘outside world’?



# Make reuse easier!

- Increased visibility
- Social value:
  - stimulates collaboration
  - accelerates innovation
- External quality control



# What's holding us back?

- Fear?
- Habit?



# Practical Constraints

1. Task specificity
2. Formats
3. Different conceptual models
4. No machine-readable definitions
5. Lack of metadata



# I. Task-specificity

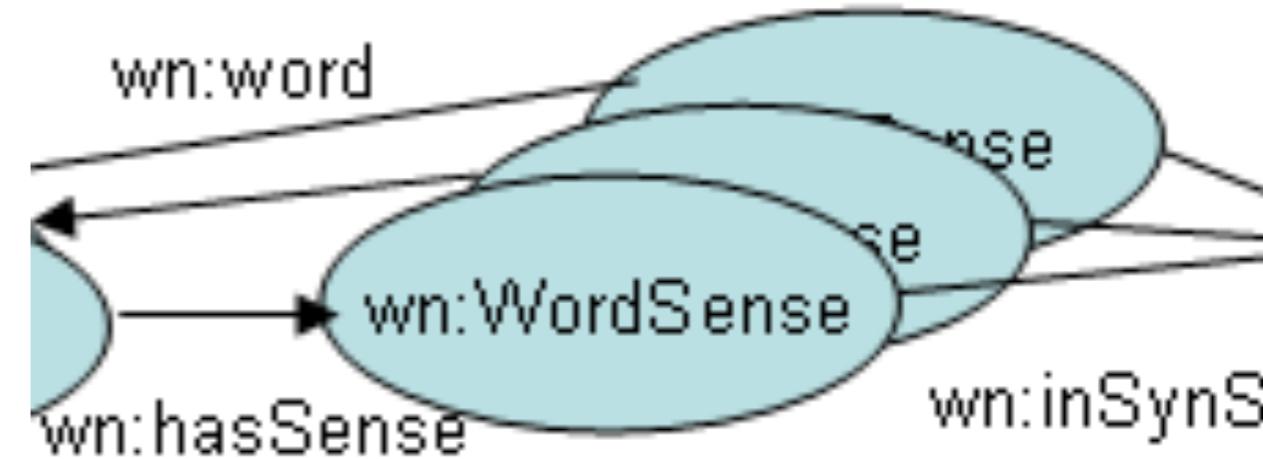
- Resources are often geared towards one specific task e.g., part-of-speech tagging, named entity recognition
- *How can we make our resources more flexible?*



## 2. Formats

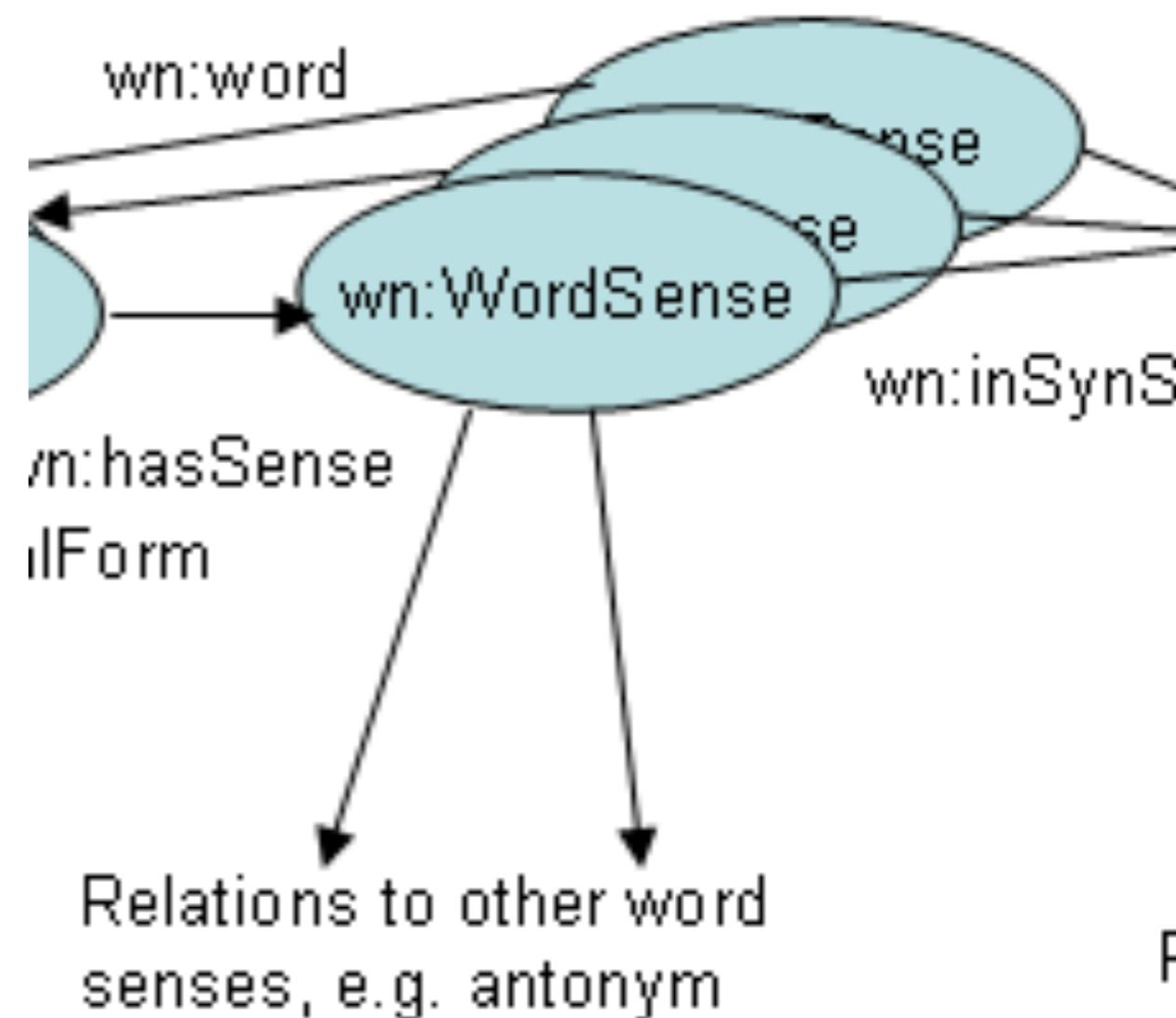
- XML, inline XML, CSV, one word per line, one sentence per line, slashtags, ARFF,





### 3. Conceptual Models

- An NP is an NP is an NP?
- “President Obama signed the National Defense Authorization Act after months of debate”
  - NE:“President Obama”?
  - NE:“Obama”?



## 4. Lack of Machine- Readable definitions

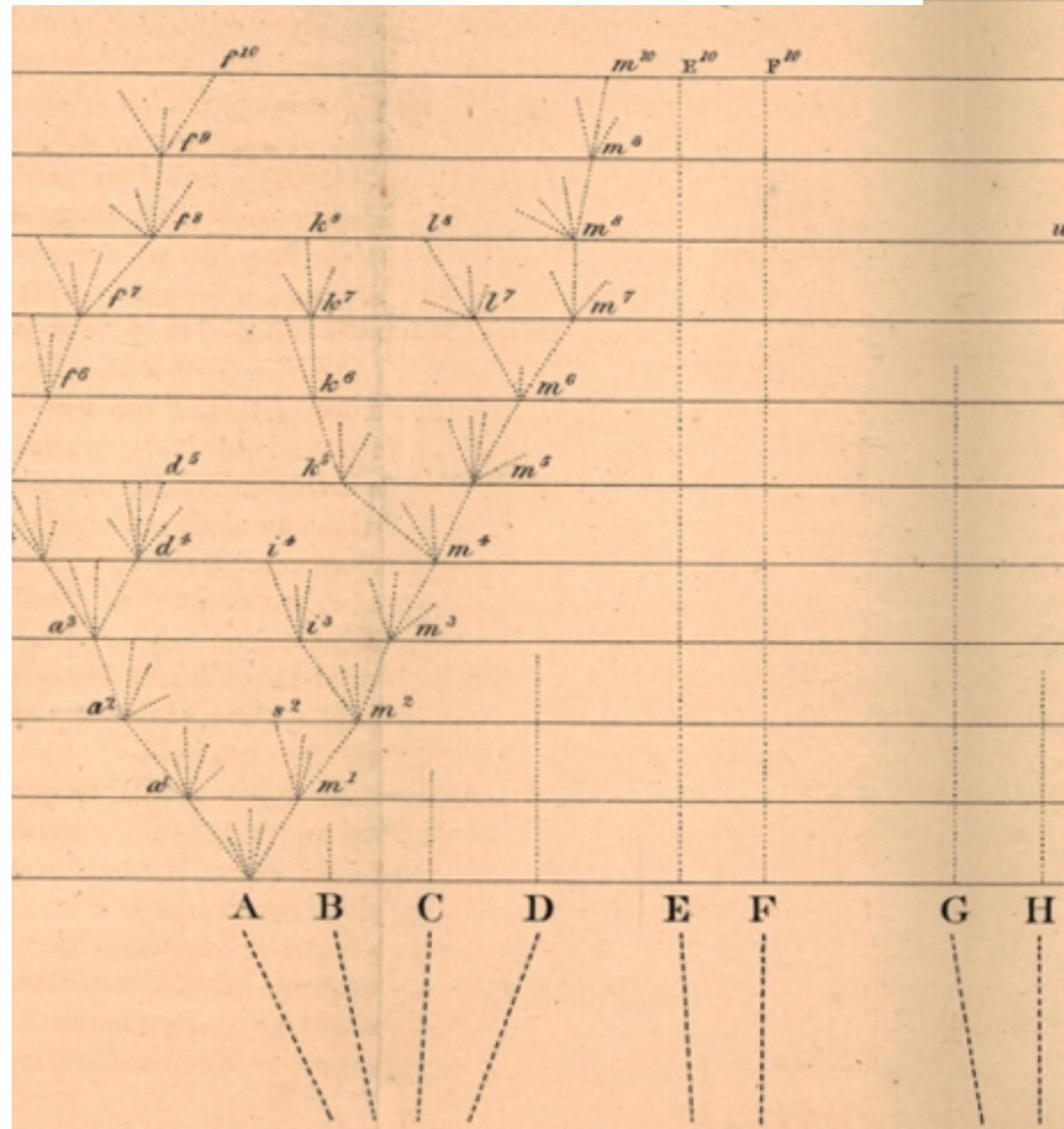
- For integration or reuse manual effort is needed
  - time consuming
  - difficult to track definitions
  - not scalable

A standard barcode graphic consisting of vertical black bars of varying widths on a white background.

3456 78

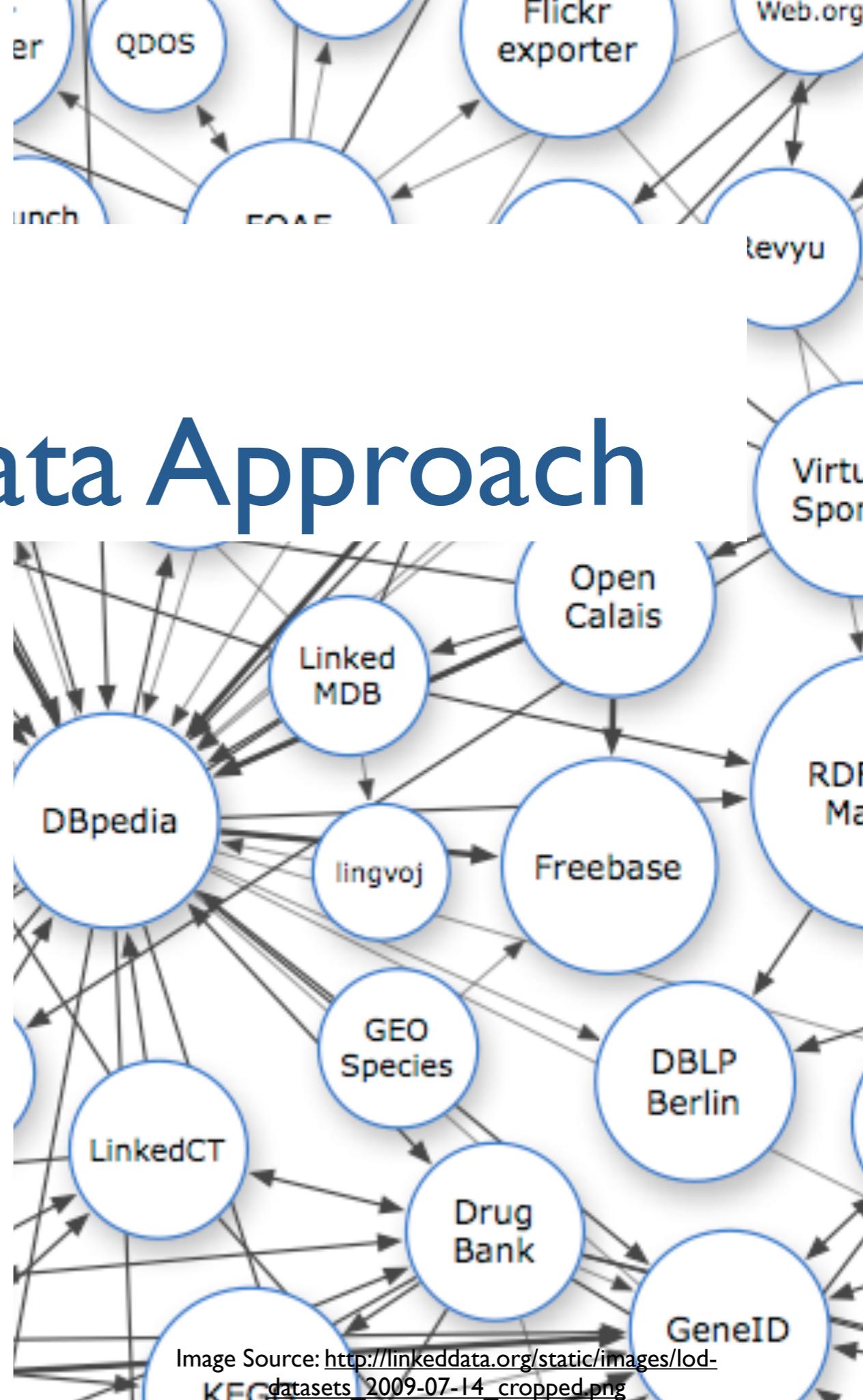
# 5. Lack of Metadata

- Can I trust this data provider?
- How was this data created?
- How many annotators?
  - for the entire data set?
  - per instance?
- If generated automatically,  
what were the parameters?



# A Linked Data Approach

- Linked Data is not a magic solution to all problems
- ...but it is better than what we've got at this moment

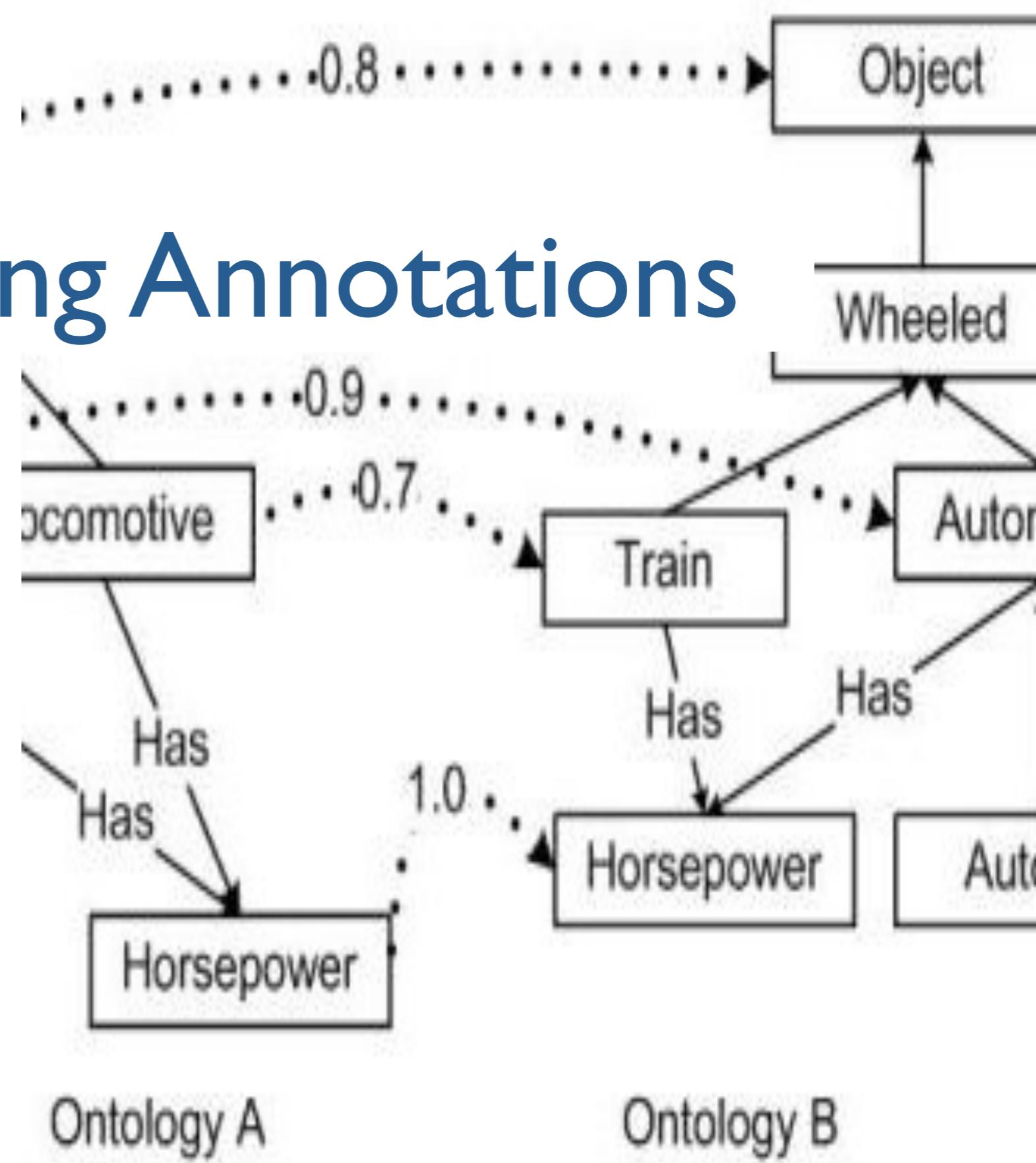


# I. Using RDF

- RDF is not inherently better than some other formats, but it is used by many
- + SPARQL makes it easy to retrieve data

## 2. Mapping Annotations

- A single conceptual model for all linguistic resources is not going to happen
- ...but can we spot the similarities between models and utilise that?



## 3. Grounding

- It's only linked data if you link it to other sources
- Added bonus: automatic sense disambiguation + access to a wealth of extra knowledge about your data item

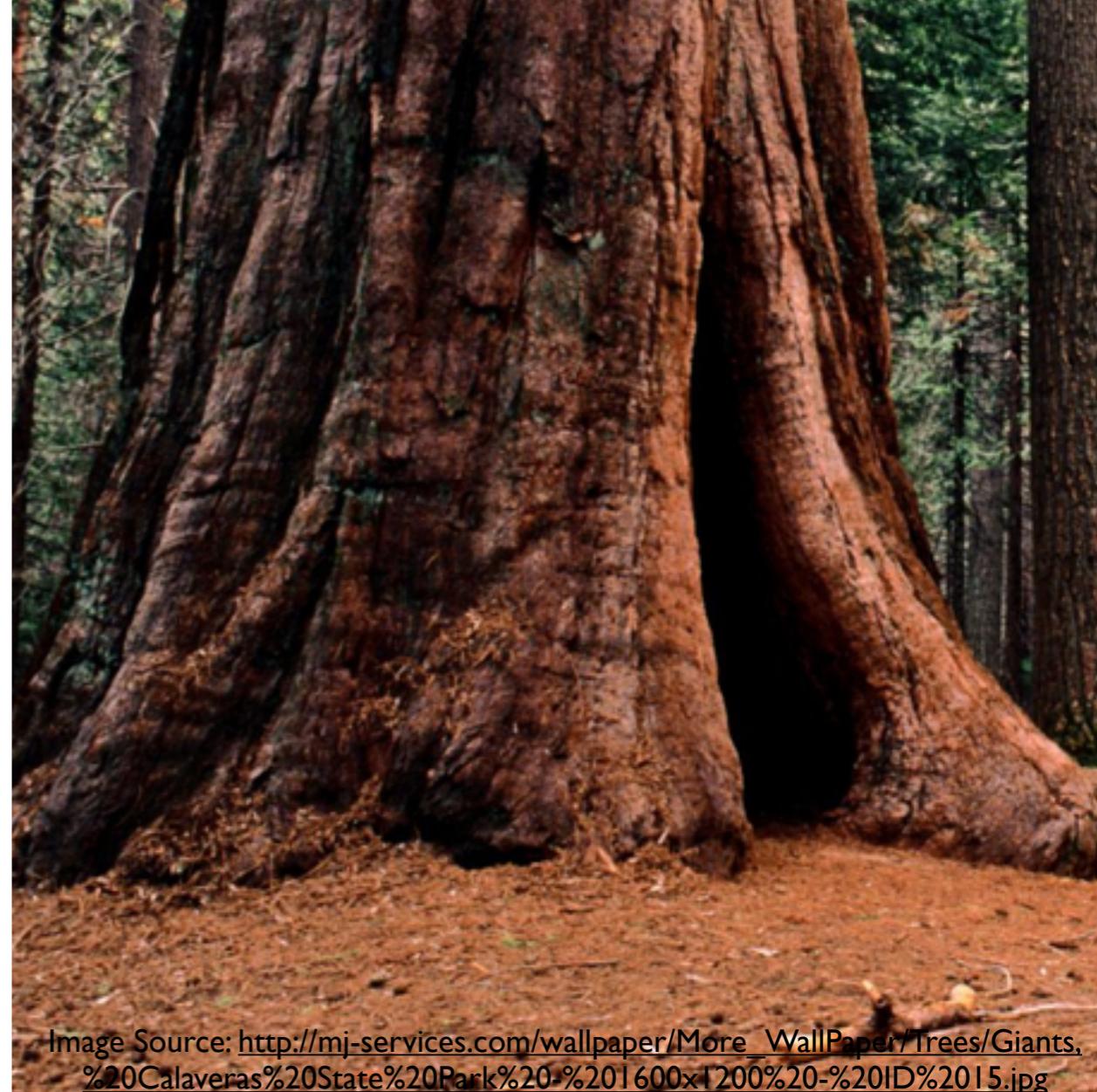
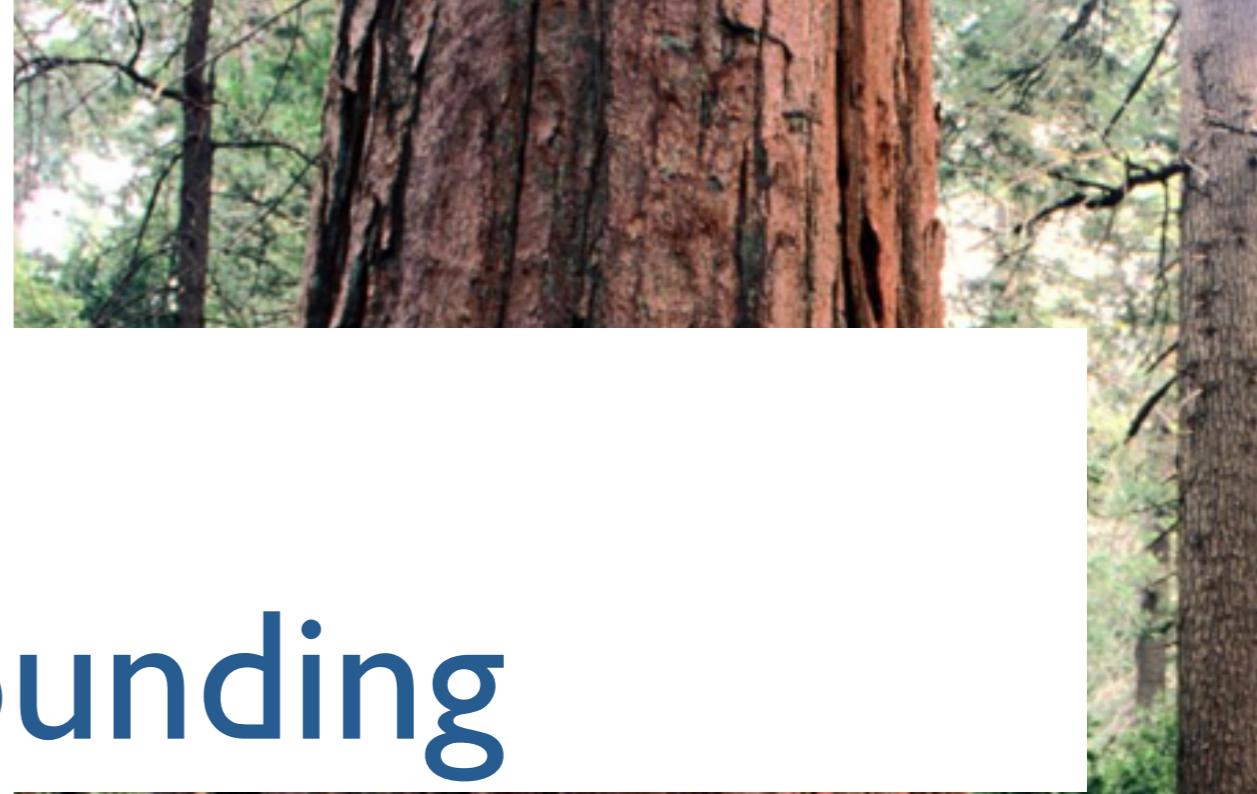


Image Source: [http://mj-services.com/wallpaper/More\\_WallPaper/Trees/Giants,%20Calaveras%20State%20Park%20-%201600x1200%20-%20ID%2015.jpg](http://mj-services.com/wallpaper/More_WallPaper/Trees/Giants,%20Calaveras%20State%20Park%20-%201600x1200%20-%20ID%2015.jpg)

## 4. Define Your Metadata

- Include your data model
- Preferably give each instance's provenance
  - collection
  - annotation/creation
  - previous versions
  - confidence



Image Source: <http://www.wineaustralia.com/australia/Portals/2/November%20E-news/Magazine%20of%20Provenance%20Final.jpg>

# Conclusions

- Look for similarities between resources
- Say where your resource comes from
- Use standards, or make it easy for others to convert your data to a standard
- Link to other data



# Questions?



[marieke@cs.vu.nl](mailto:marieke@cs.vu.nl)  
<http://www.cs.vu.nl/~marieke>

# Acknowledgment

- This work is funded by NWO in the CATCH programme, grant 640.004.801

