# Databases

## Graph Databases
### Motivation

André Santanchè e Patrícia Cavoto
Institute of Computing – UNICAMP
August 2015

# Graph Database

## Graph as the basic data model of a database
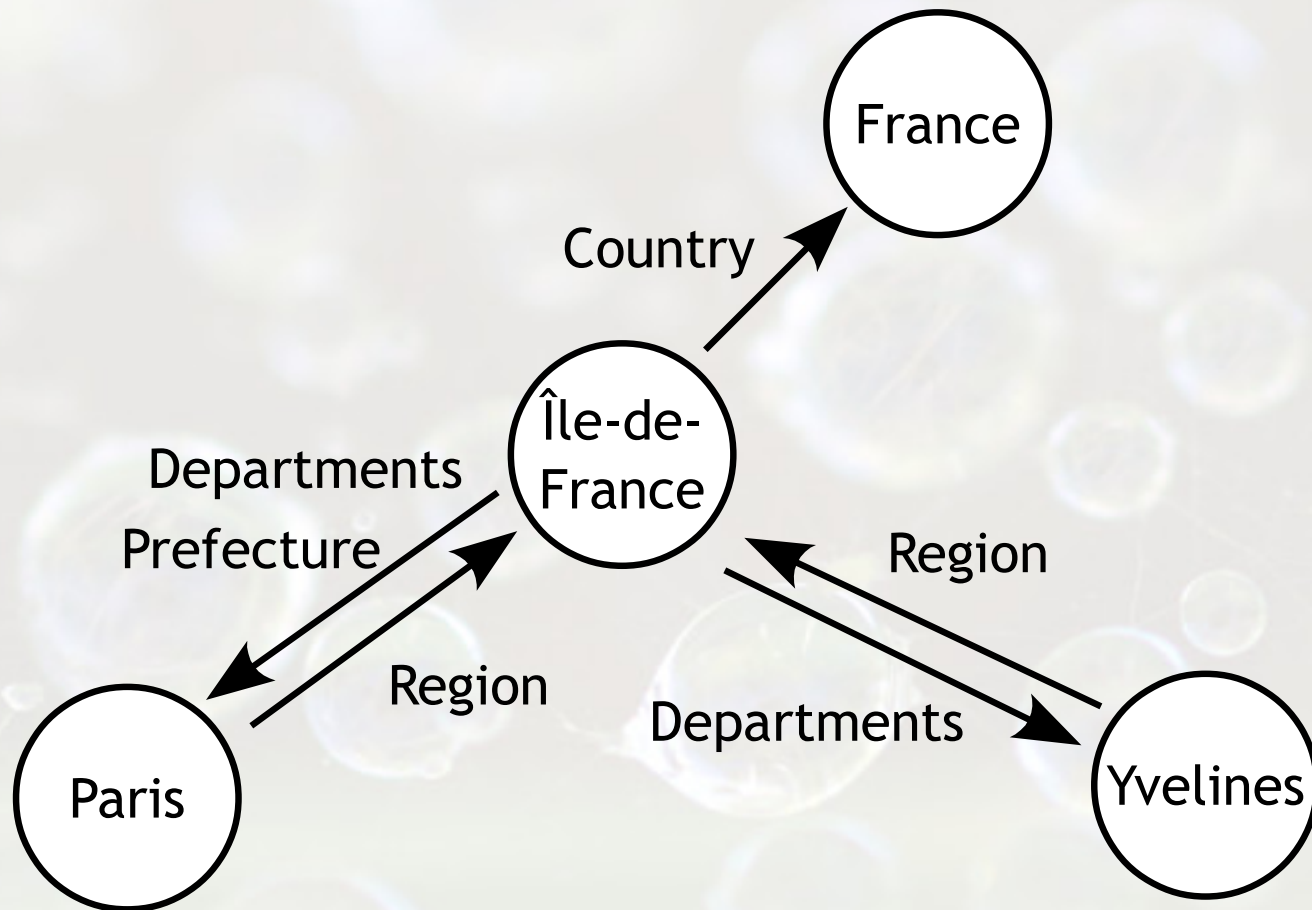
# Why Graphs?

# Why Graphs?

- The Web Effect
  - Linked Data
  - Social Networks and Social Content
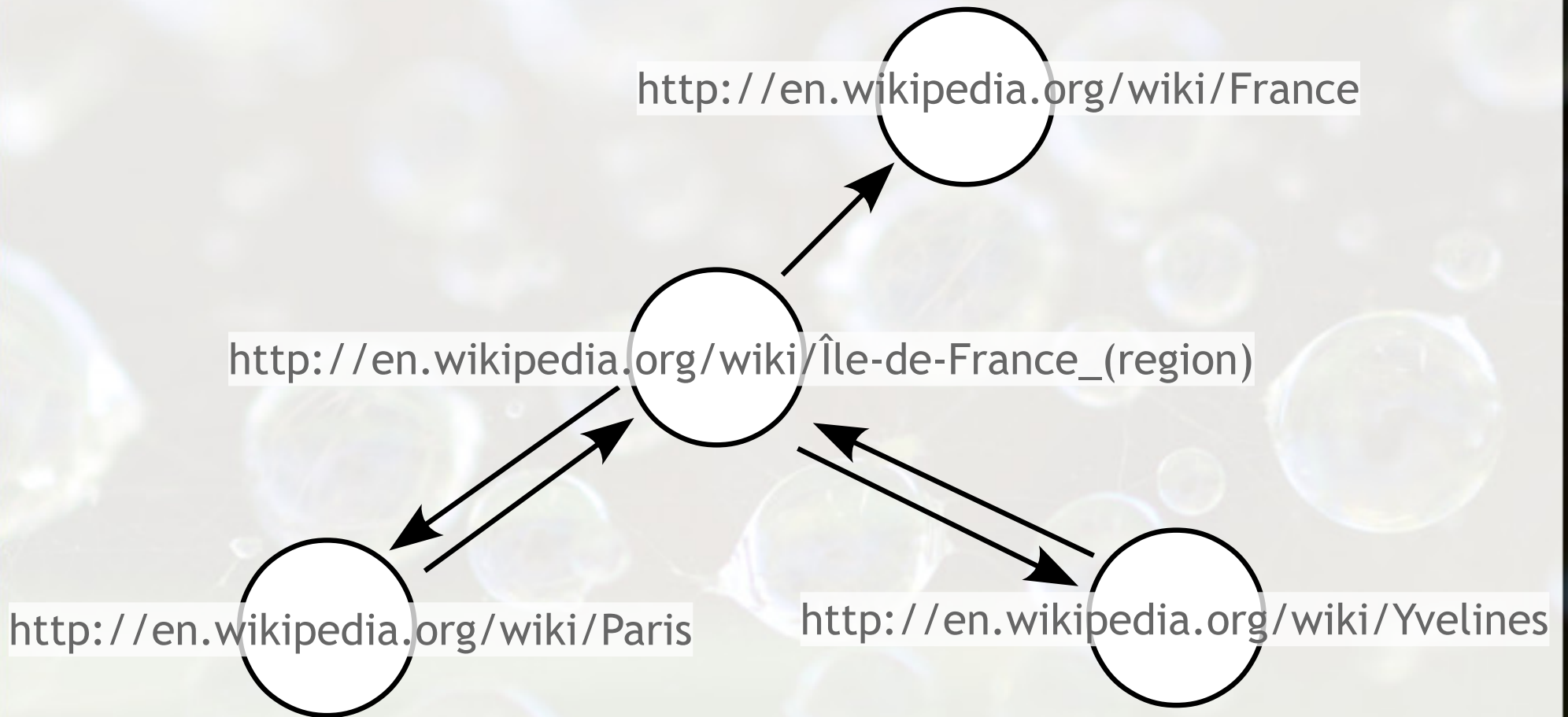- Sharing and interconnecting
- Complex Networks

Why Graphs?
# Linked Data

# DBPedia

# DBPedia (URIs)

http://en.wikipedia.org/wiki/France

http://en.wikipedia.org/wiki/Île-de-France_(region)

http://en.wikipedia.org/wiki/Paris

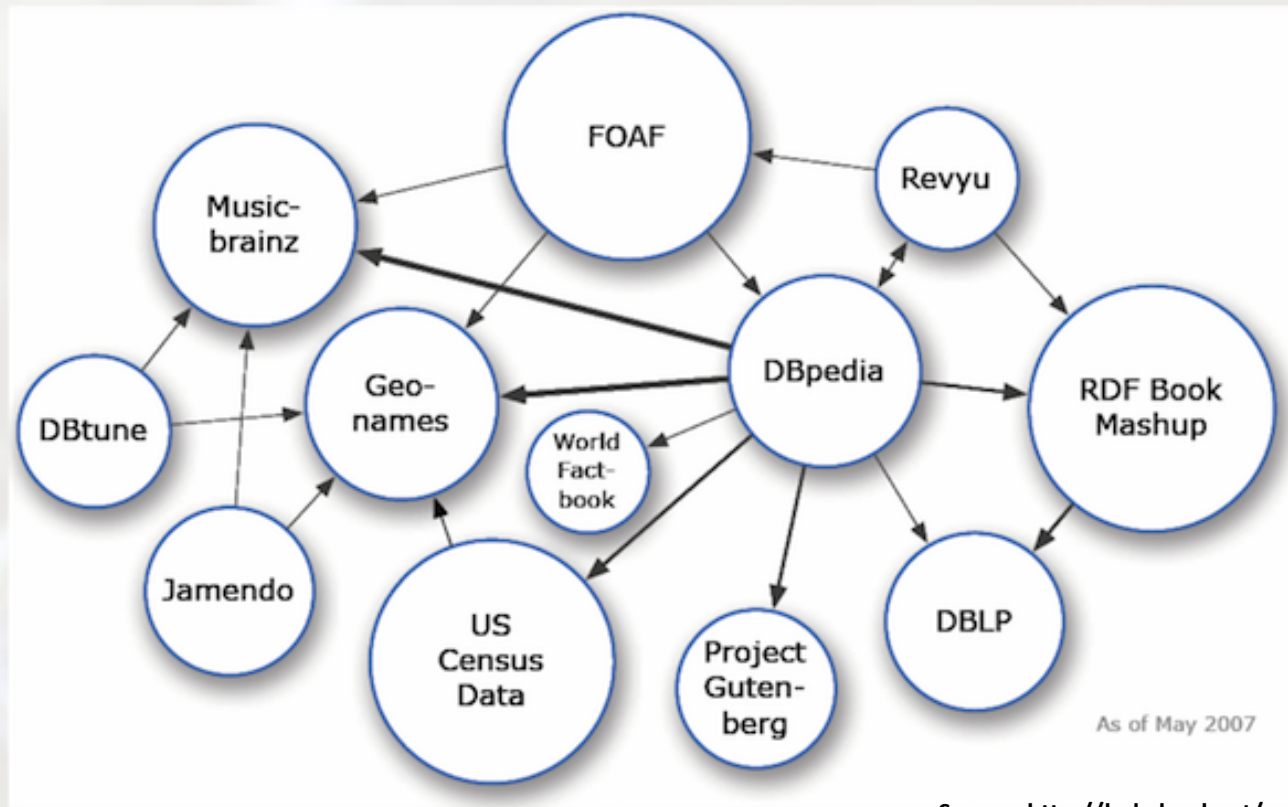http://en.wikipedia.org/wiki/Yvelines

# DBPedia – English

- **4 million things**
- **3.22 million classified in a consistent ontology**
  - 832,000 persons
  - 639,000 places (427,000 populated)
  - 372,000 creative works
    - 116,000 music albums; 78,000 films; 18,500 video games
  - 209,000 organizations
  - 226,000 species
  - 5,600 diseases.

# DBPedia – International

- 119 languages
- 24.9 million things
- 16.8 million interlinked with English
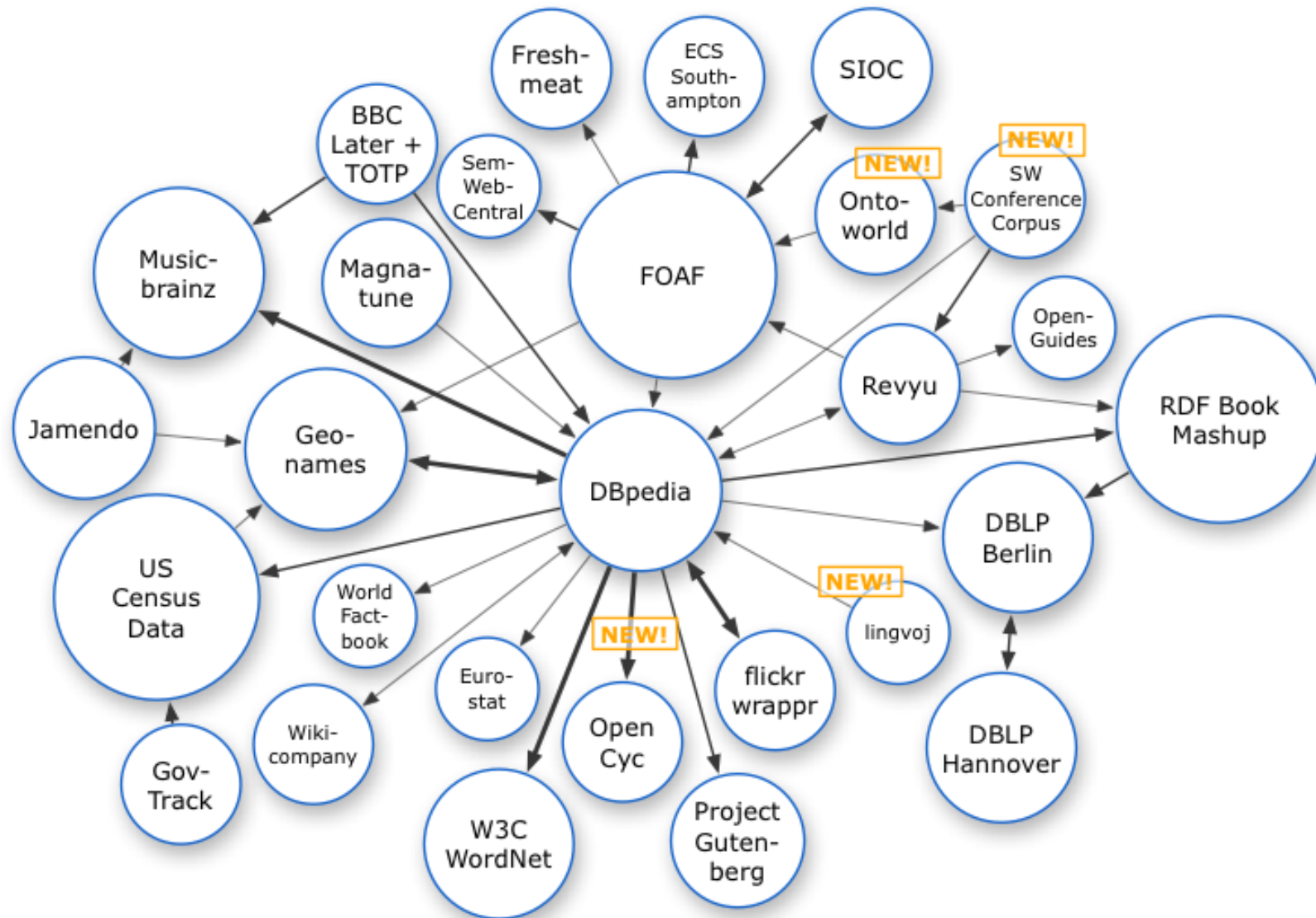- 12.6 million unique things

# Linked Data



Source: http://lod-cloud.net/

Datasets published following Linked Data 'format':  **05/2007**

# Linked Data



Source: http://lod-cloud.net/

Datasets published following Linked Data 'format':  **11/2007**

# Linked Data



Source: http://lod-cloud.net/

As of September 2008

Datasets published following Linked Data 'format':  **2008**

# Linked Data

Datasets published following Linked Data 'format': **2009**

# Linked Data

Datasets published following Linked Data 'format':  **2010**

# Linked Data



Datasets published following Linked Data 'format':  **2011**

# Linked Data
## 04/2014



Source: http://lod-cloud.net/

Linked Datasets as of August 2014

Publications
Life Sciences
Cross-Domain
Social Networking
Geographic
Government
Media
User-Generated Content
Linguistics

Why Graphs?
# Social Networks and Social Content

# Searching Pet

# Pet Shop Boys

Anirudh Koul



pet shop boys, pet, concert, tour

cat, kittens, eyes,
ears, pet, animal

dog, pet, animal,
funny, glasses

dog, pet, alaskan malamute

cat, kitty, eyes,
pretty

cat, kitten, garden,
pet

recycle, pet, plastic bottle,
polyethylene terephthalate

recycle, pet, plastic bottle,
polyethylene terephthalate

wine, pet, bottle

**sfroehlich1121**

cat, kittens, eyes,
ears, pet, animal

**Edward Corpuz**

dog, pet, alaskan malamute

**shorty_nz_2000**

dog, pet, animal,
funny, glasses

**sfroehlich1121**

cat, kitty, eyes,
pretty

**Jay Woodworth**

cat, kitten, garden,
pet

## Nemo's great uncle



recycle, pet, plastic bottle,
polyethylene terephthalate

## FaceMePLS



recycle, pet, plastic bottle,
polyethylene terephthalate

## Karl Baron



wine, pet, bottle

# Graph users/tags/resources

# Graph users/tags/resources

# Co-ocurrences and Latent Semantics

# Co-ocurrences and Latent Semantics

# Co-ocurrences and Latent Semantics

# Co-ocurrences and Latent Semantics

# Social Effect
# Suggested/Reinforced Tags

sfroehlich1121



```
black, cat, kitty, katze, long, hair, blue, eyes, pretty,
    canon, t1i, 500d, ef 100mm f/2.8 usm macro
```

# M1
# co-occurrences in a folksonomy

*tags / co-occurrences*

# Tagsets

# M2
# tagsets and their co-occurrences

*tagsets/co-occurrences*

# M4
# Social Ontology

*tagsets / typed relations*



*is-a*

# M3
# simplified model of an ontology

*concepts/relations*

# M5
# Folksonomized Ontology (FO)

*concepts / typed relations (ontology)*

*+*

*information content / co-occurrences (folksonomy)*

# Meaning of Pet needs meaning of other tags?

# Matrix

## What is the meaning of the word Love?

Why Graphs?
# Sharing and Interconnecting

# Models to Describe

# Describing Prehistoric Animals



| MNHN A. C. 8592 | | |
|---|---|---|
| **Is a** | Plesiosaurus dolichodeirus | |
| **Origin** | Lyme Regis | England |
| **Recognized** | 1824 | |
| **Size** | 5 | |

# Describing Prehistoric Animals

**SIPB R 90**

| Is a | Plesiosaurus dolichodeirus | |
|---|---|---|
| Origin | Lyme Regis | England |
| Recognized | 1830 | |
| Size | 5 | |

**STC223**

| Is a | Plesiosaurus gurgitis | |
|---|---|---|
| Origin | St. Croix | Switzerland |
| Recognized | 1964 | |
| Size | 3.5 | |

**MNHN 1912.20**

| Is a | Triceratops horridus | |
|---|---|---|
| Origin | Lance Creek | EUA |
| Recognized | 1889 | |
| Size | 9 | |

**FMNH PR2081**

| Is a | Tyrannosaurus rex | |
|---|---|---|
| Origin | Hell Creek | EUA |
| Recognized | 1990 | |
| Size | 12.3 | |

**Sue**

# Table

| Id | Is a | Origin Place | Origin Country | Recognized | Size |
|---|---|---|---|---|---|
| **MNHN A. C. 8592** | Plesiosaurus dolichodeirus | Lyme Regis | England | 1824 | 5 |
| **SIPB R 90** | Plesiosaurus dolichodeirus | Lyme Regis | England | 1830 | 5 |
| **STC223** | Plesiosaurus gurgitis | St. Croix | Switzerland | 1964 | 3.5 |
| **MNHN 1912.20** | Triceratops horridus | Lance Creek | EUA | 1889 | 9 |
| **FMNH PR2081** | Tyrannosaurus rex | Hell Creek | EUA | 1990 | 12.3 |

# Table

Excellent to Manage Data with Predictable Static Schema

Sharing?

# Documents and XML

# Documents and XML

# Documents and XML

# How much hierarchical is your data?

# Back to the Table

**Le Muséum national d'Histoire naturelle**

_includes_

| | Id | Is a | Origin Place | Origin Country | Recognized | Size |
|---|---|---|---|---|---|---|
| | **MNHN A. C. 8592** | Plesiosaurus dolichodeirus | Lyme Regis | England | 1824 | 5 |
| | **SIPB R 90** | Plesiosaurus dolichodeirus | Lyme Regis | England | 1830 | 5 |
| | **STC223** | Plesiosaurus gurgitis | St. Croix | Switzerland | 1964 | 3.5 |
| | **MNHN 1912.20** | Triceratops calicornis | Lance Creek | EUA | 1888 | 9 |
| | **MNHN 1912.20b** | Triceratops horridus | Lance Creek | EUA | 1889 | 9 |
| | **FMNH PR2081** | Tyrannosaurus rex | Hell Creek | EUA | 1990 | 12.3 |

_renamed_

_near_

# The tradeoff of Static Schemas

# XML Schema
# Not designed to mix

# DataSpaces
## Pay-as-you-go Integration



Franklin, M., Halevy, A., & Maier, D. (2005). From databases to dataspaces: a new abstraction for information management. SIGMOD Rec., 34(4), 27–33.

# DataSpaces and Linked Data

Why Graphs?
# Complex Networks

# Complex Networks

- Developed steadily since 1999

- Discrete systems are represented in terms of entities and relationships

(Luciano da F. Costa, 2013)

# Patient x Diagnosis



(Gomes-Jr, 2013)

# Complex Networks

- Physical relationships
  - neurons – nodes; connections – edges
- Force relationships
  - grains – nodes; force vectors – edges
- Social relationships
- Conceptual relationships

# Complex Networks Examples

http://www-personal.umich.edu/~mejn/networks/

# Les Miserables
## (Neo4j example)

# Cadeia Alimentar FishBase

# Freshwater food web



Freshwater food web: Neo Martinez and Richard Williams.

# Contagion of TB



Contagion of TB, books on politics: Valdis Krebs, www.orgnet.com.

# Yeast proteins



Yeast proteins: Sergei Maslov and Kim Sneppen,
Specificity and stability in topology of protein networks,
Science 296, 910-913 (2002).

# What can I ask to a Graph?

# Triceratops in a Graph



| MNHN 1912.20 | | |
|---|---|---|
| **Is a** | Triceratops horridus | |
| **Origin** | Lance Creek | EUA |
| **Recognized** | 1889 | |
| **Size** | 9 | |

# Triceratops in a Graph

# Triceratops in a Graph

# Tyrannosaurus in a Graph

# Analyzing in the Space



OpenStreetMap

# Analyzing in the Space

EUA

South Dakota

Wyoming

Cheyenne River Indian Reservation

Converse County

Hell Creek

Lance Creek

# Analyzing in the Space

EUA

South Dakota

Wyoming

Cheyenne River Indian Reservation

Converse County

Hell Creek

Lance Creek

# GeoNames

MNHN 1912.20
- is a → Triceratops horridus
- size → 9
- recognized → 1889
- origin → Lance Creek http://www.geonames.org/5829995/

United States http://www.geonames.org/6252001

Wyoming http://www.geonames.org/5843591

Niobrara County http://www.geonames.org/5833446

Lance Creek http://www.geonames.org/5829995/

# Processing the Query

- **Process by pattern**
  - Find [species] whose [origin] → <u>(0..*) part of</u> → EUA

- **Process by inference**
  - Find [species] whose [origin] → EUA
  - Rules:
    - If (A) origin (B) and (B) part of (C) => (A) origin (C)
    - If (A) part of (B) and (B) part of (C) => (A) part of (C)

# Looking the Topology

# Graph Topology

# Graph Topology

# Patient x Diagnosis



- Retrieve candidate **diagnosis** relevant to a given **patient** based on her **symptoms** and correlations with other similar patients.

(Gomes-Jr, 2013)

# From Dataspaces to Ontologies

From Dataspaces to Ontologies

Ontology

Linked Graph

Dataspace

# Building Patterns

# Data Exploration Prototype

# Integrating phenotypes
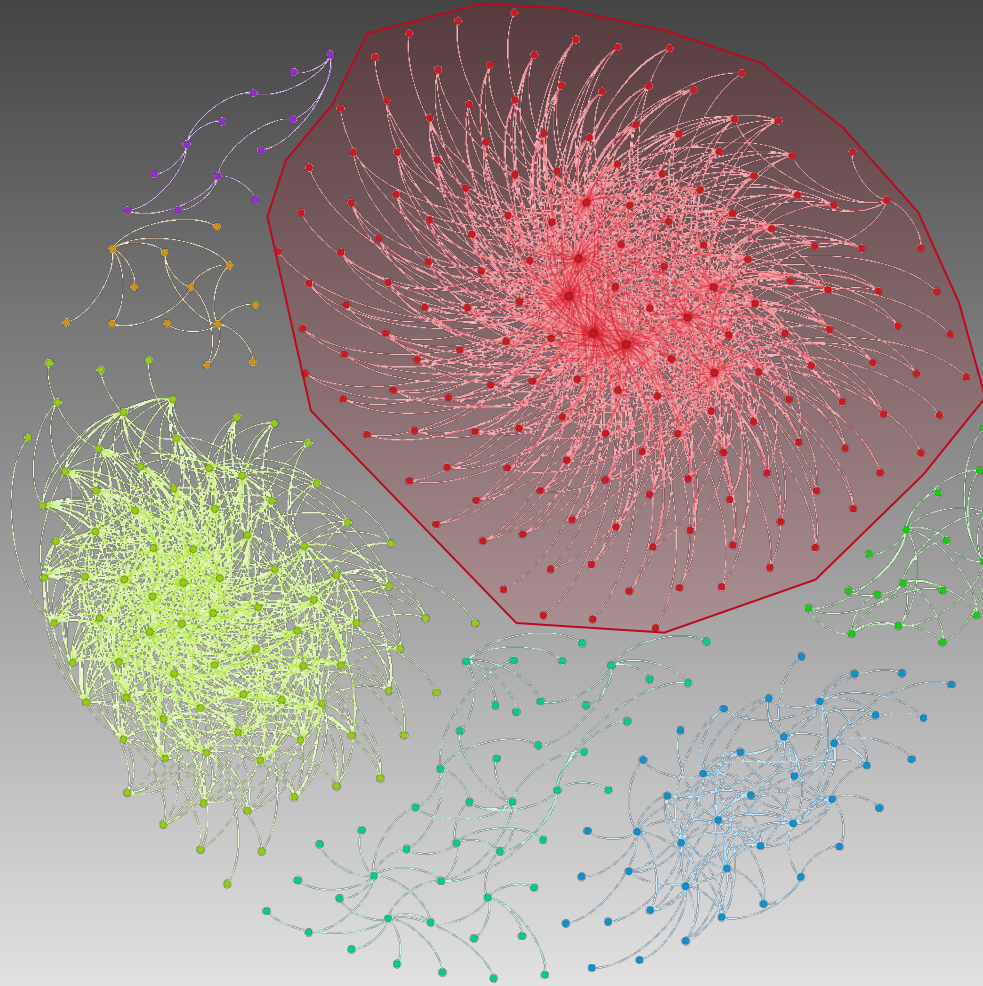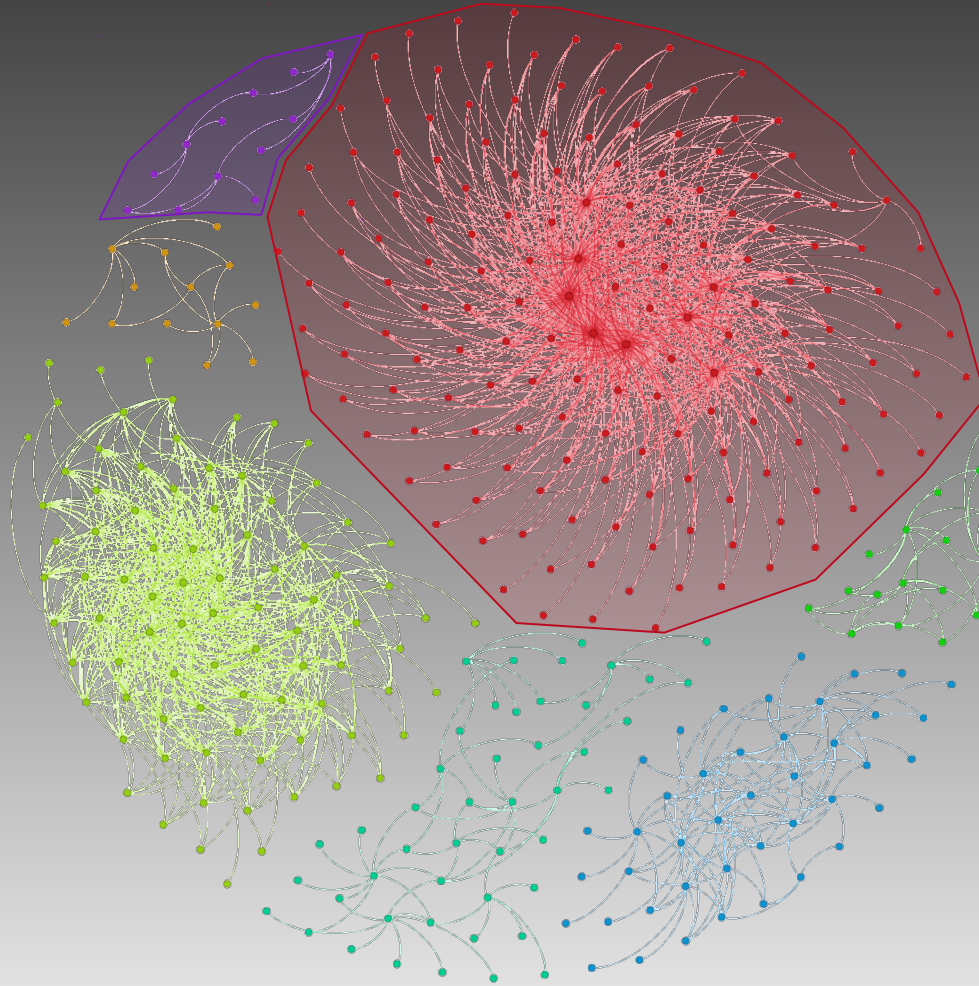
7 distinct morphological descriptions:
- genus Varanus
- *Varanus gouldii*
- *Varanus indicus*
- *Varanus prasinus*
- *Varanus salvator*
- *Varanus spiny*
- *Varanus timorensis*

# Our approach

7 distinct morphological descriptions:

· genus Varanus
· *Varanus gouldii*
· *Varanus indicus*
· *Varanus prasinus*
· *Varanus salvator*
· *Varanus spiny*
· *Varanus timorensis*

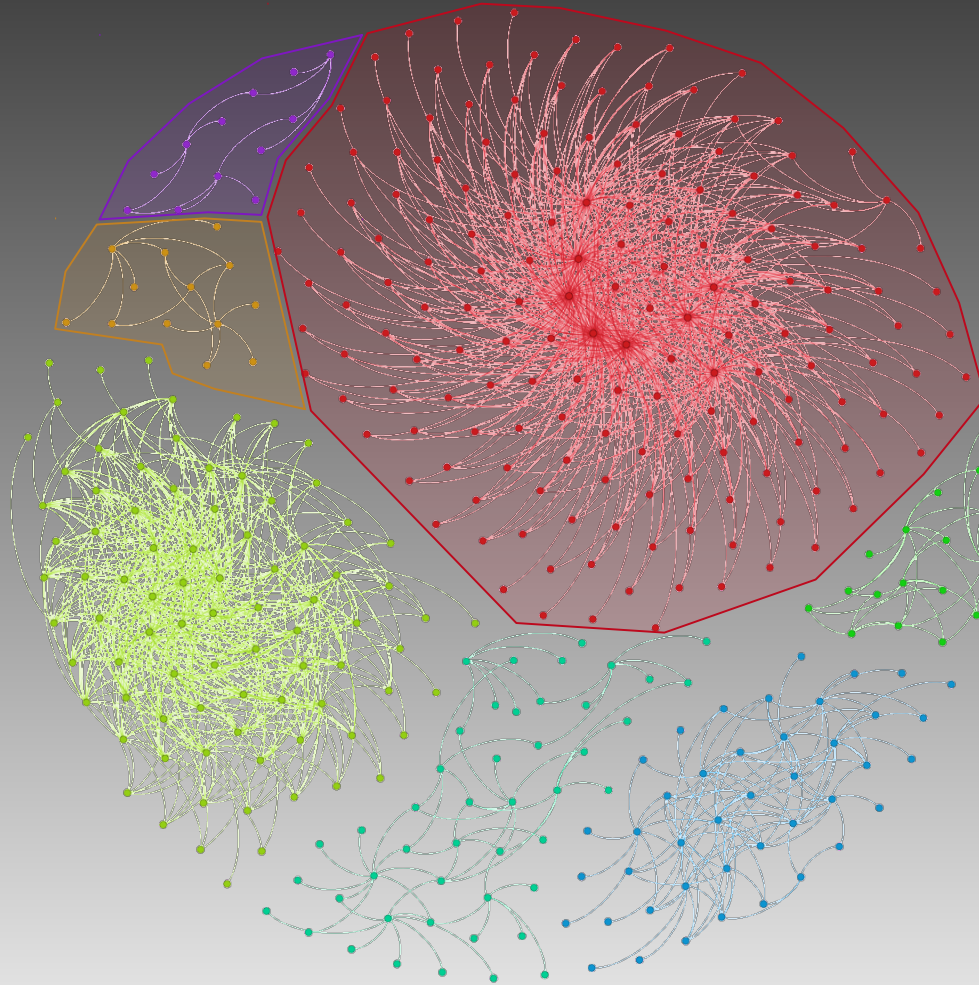# Our approach

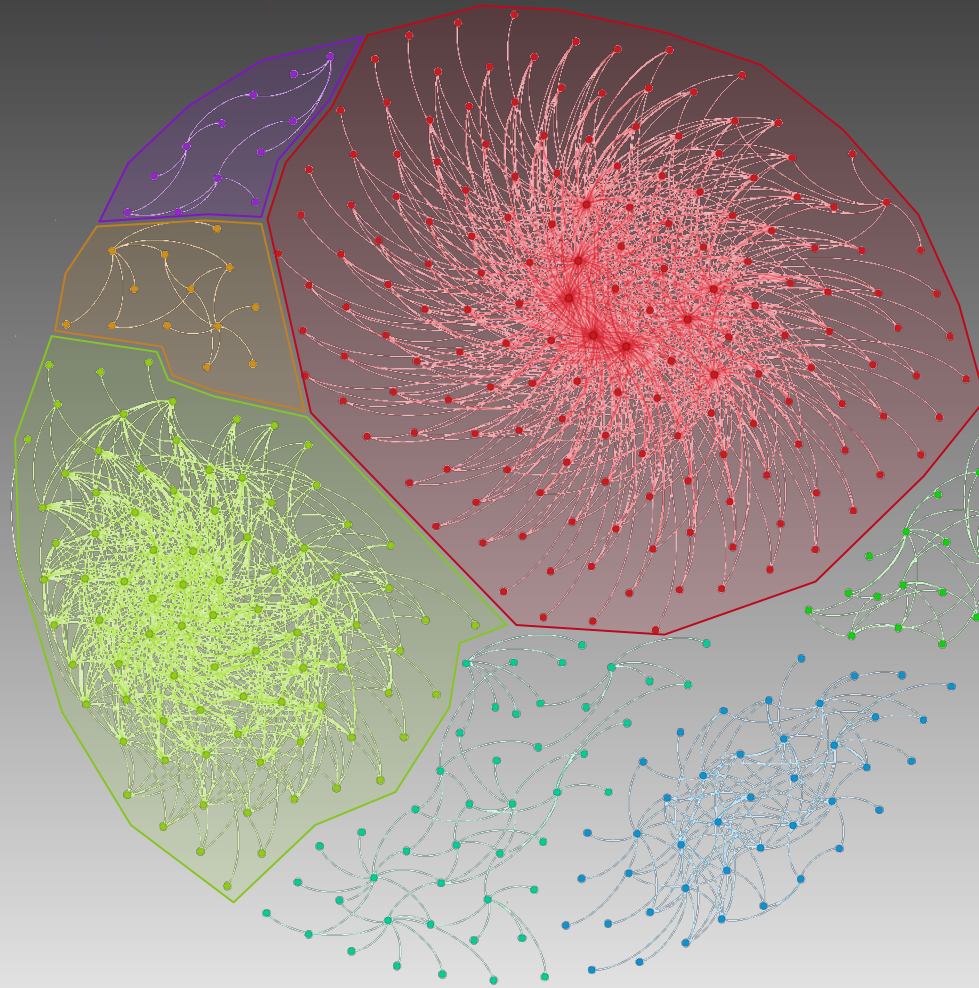7 distinct morphological descriptions:
- genus Varanus
- *Varanus gouldii*
- *Varanus indicus*
- *Varanus prasinus*
- *Varanus salvator*
- *Varanus spiny*
- *Varanus timorensis*

# Our approach

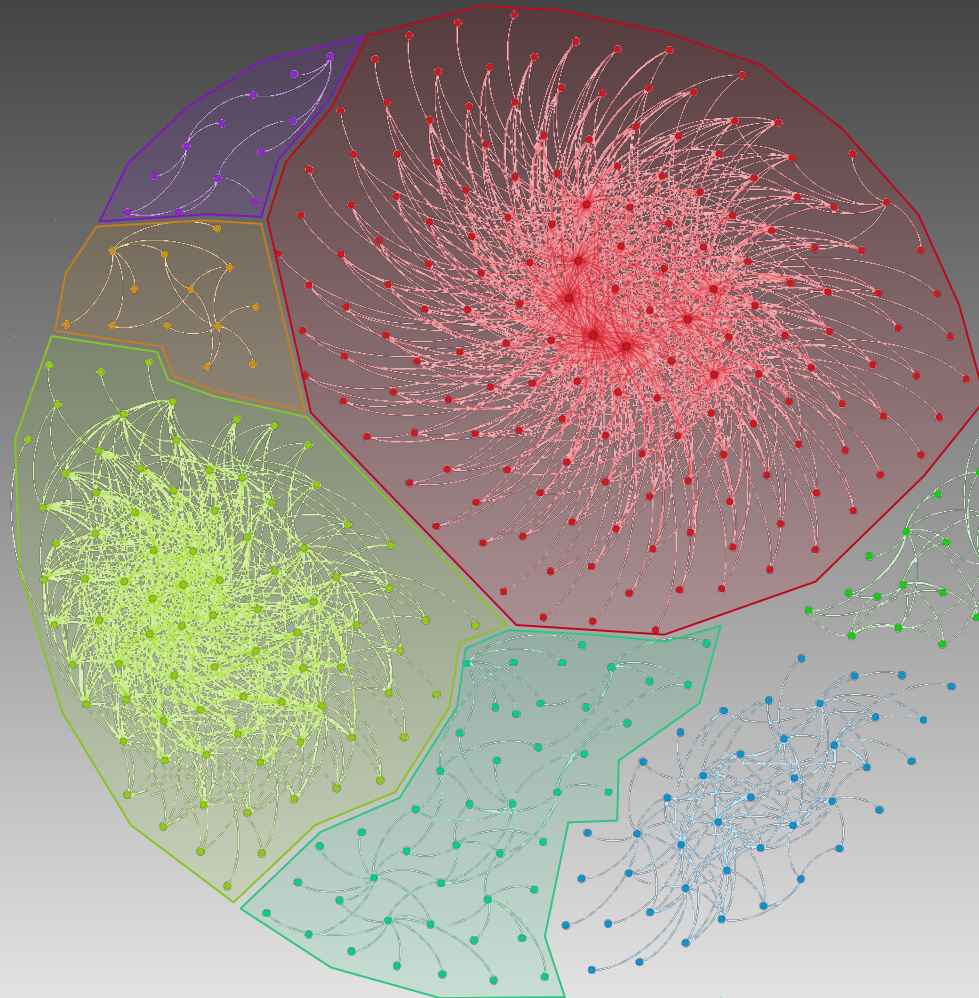7 distinct morphological descriptions:

- genus Varanus
- *Varanus gouldii*
- *Varanus indicus*
- *Varanus prasinus*
- *Varanus salvator*
- *Varanus spiny*
- *Varanus timorensis*

# Our approach

7 distinct morphological descriptions:
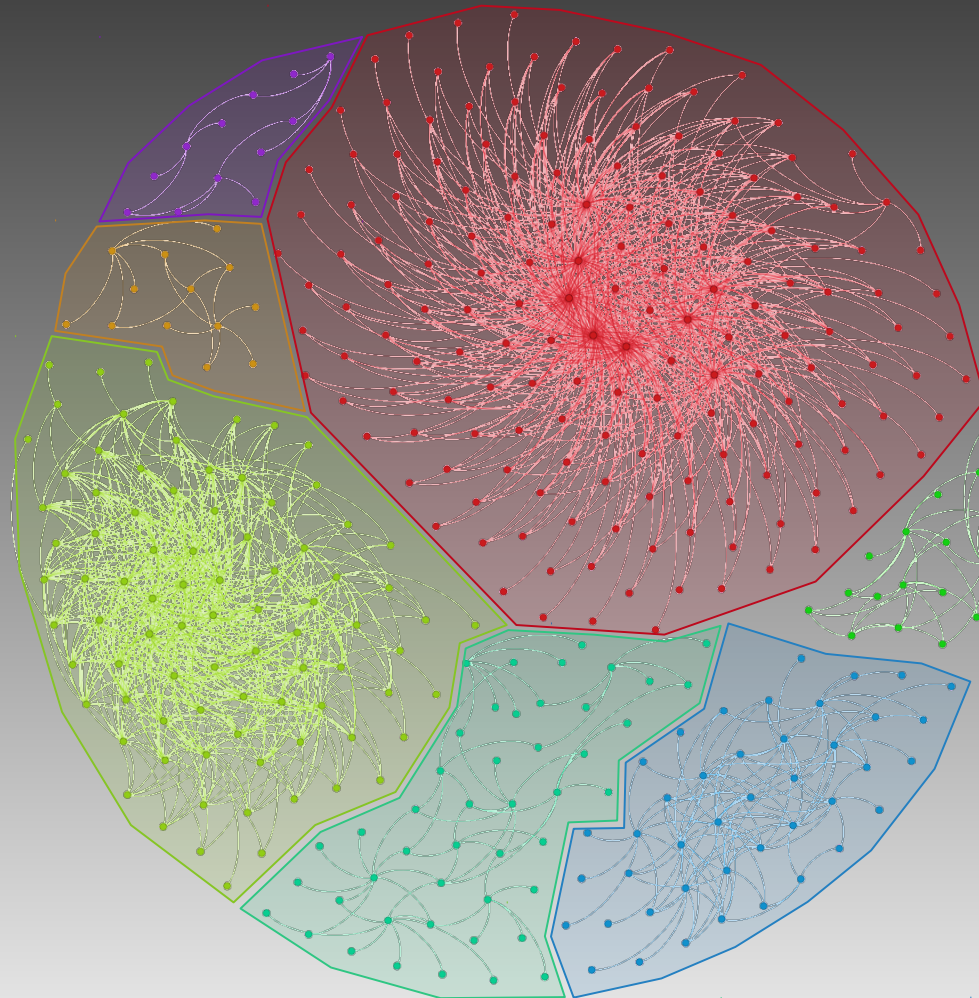- genus Varanus
- *Varanus gouldii*
- *Varanus indicus*
- *Varanus prasinus*
- *Varanus salvator*
- *Varanus spiny*
- *Varanus timorensis*

# Our approach

7 distinct morphological descriptions:
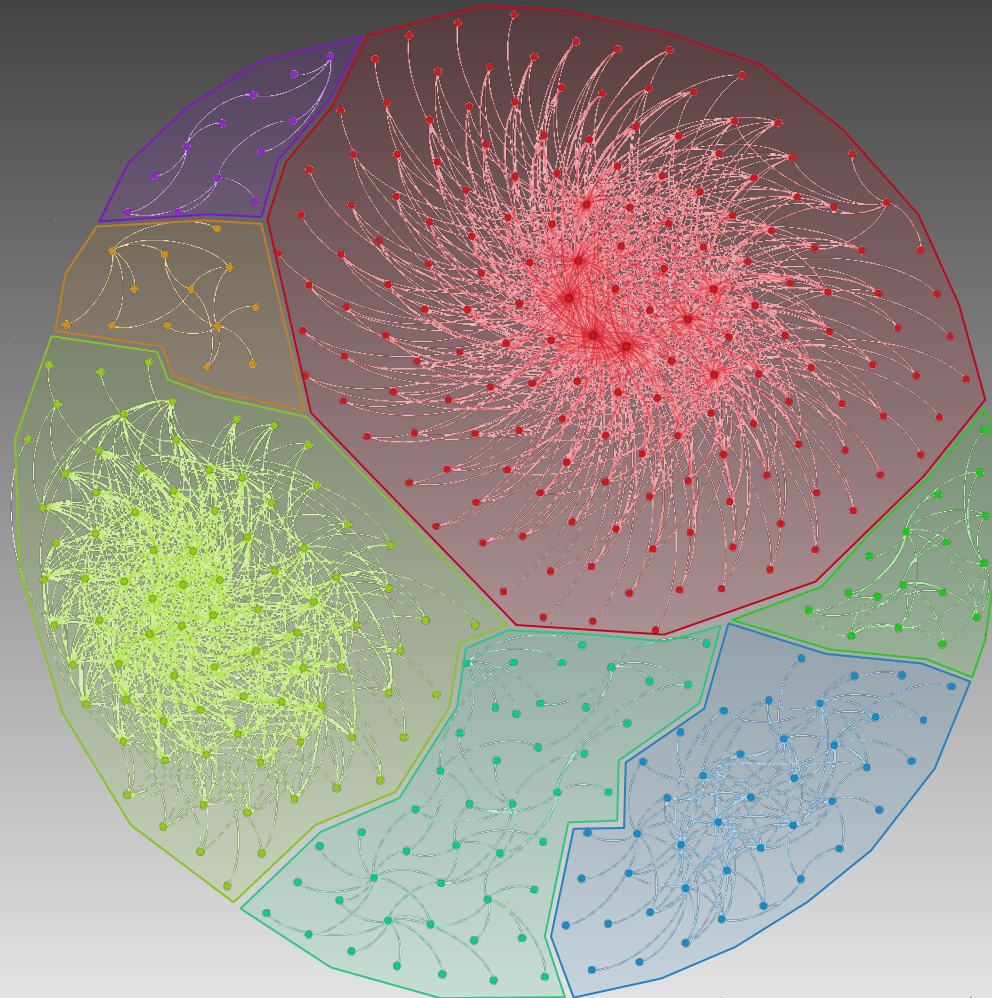
- genus Varanus
- *Varanus gouldii*
- *Varanus indicus*
- *Varanus prasinus*
- *Varanus salvator*
- *Varanus spiny*
- *Varanus timorensis*

# Our approach

7 distinct morphological descriptions:
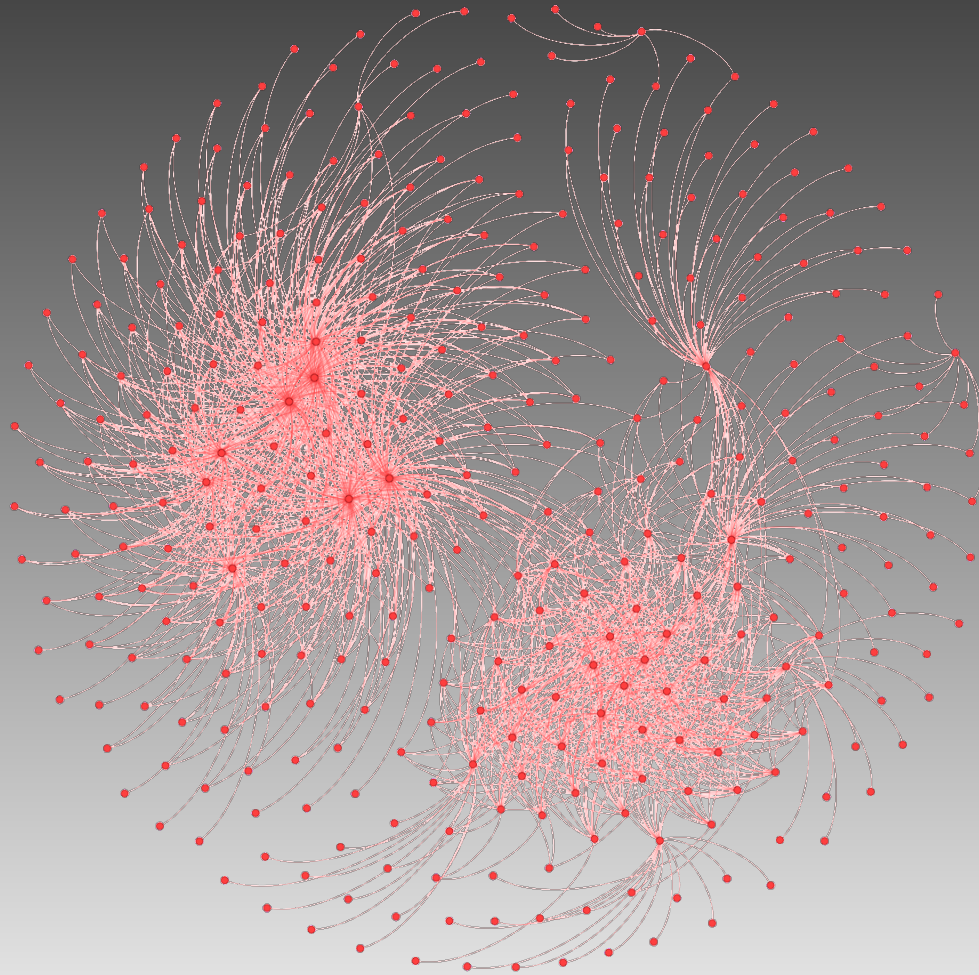- genus Varanus
- *Varanus gouldii*
- *Varanus indicus*
- *Varanus prasinus*
- *Varanus salvator*
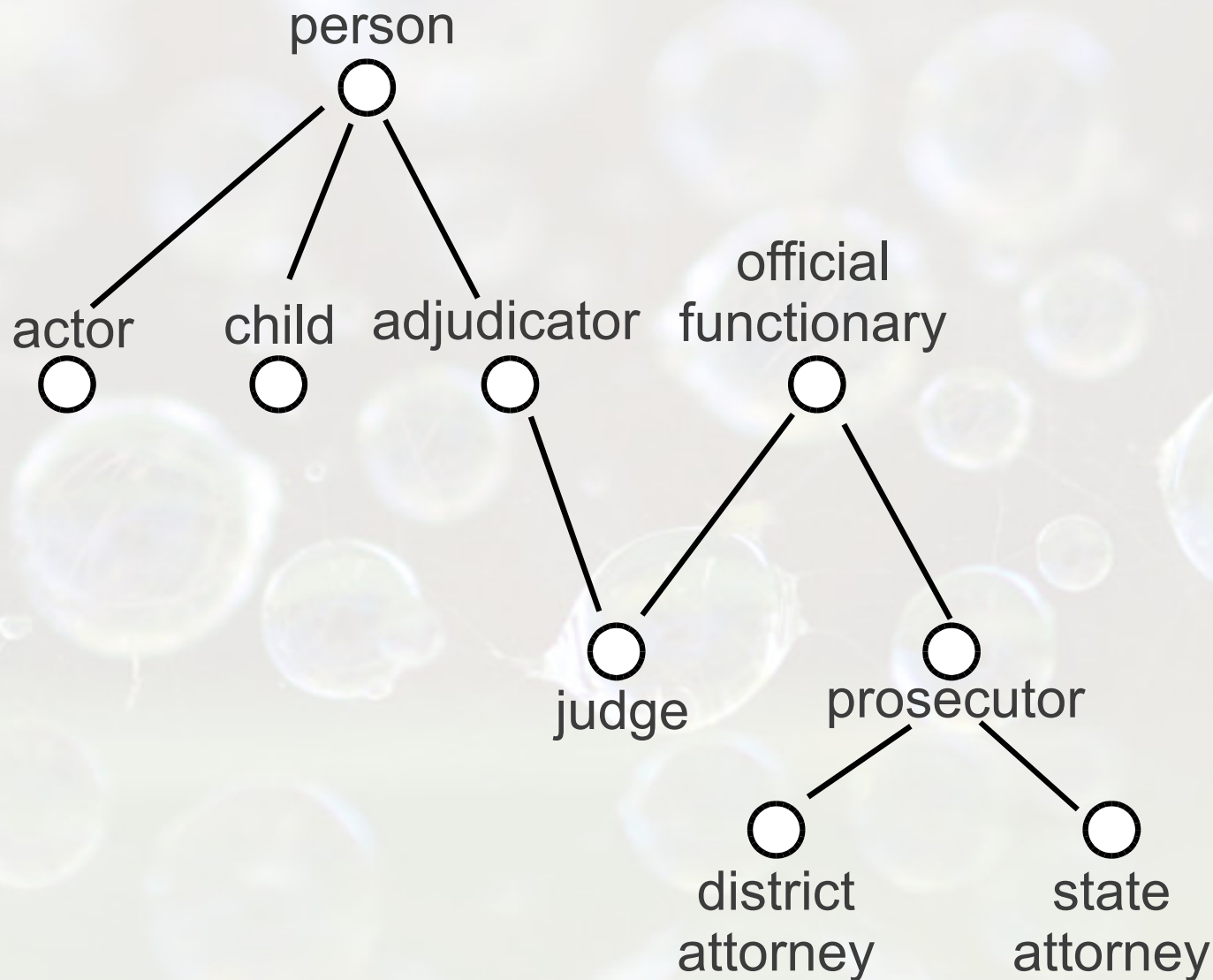- *Varanus spiny*
- *Varanus timorensis*

# Our approach

7 distinct morphological descriptions:
- genus Varanus
- *Varanus gouldii*
- *Varanus indicus*
- *Varanus prasinus*
- *Varanus salvator*
- *Varanus spiny*
- *Varanus timorensis*

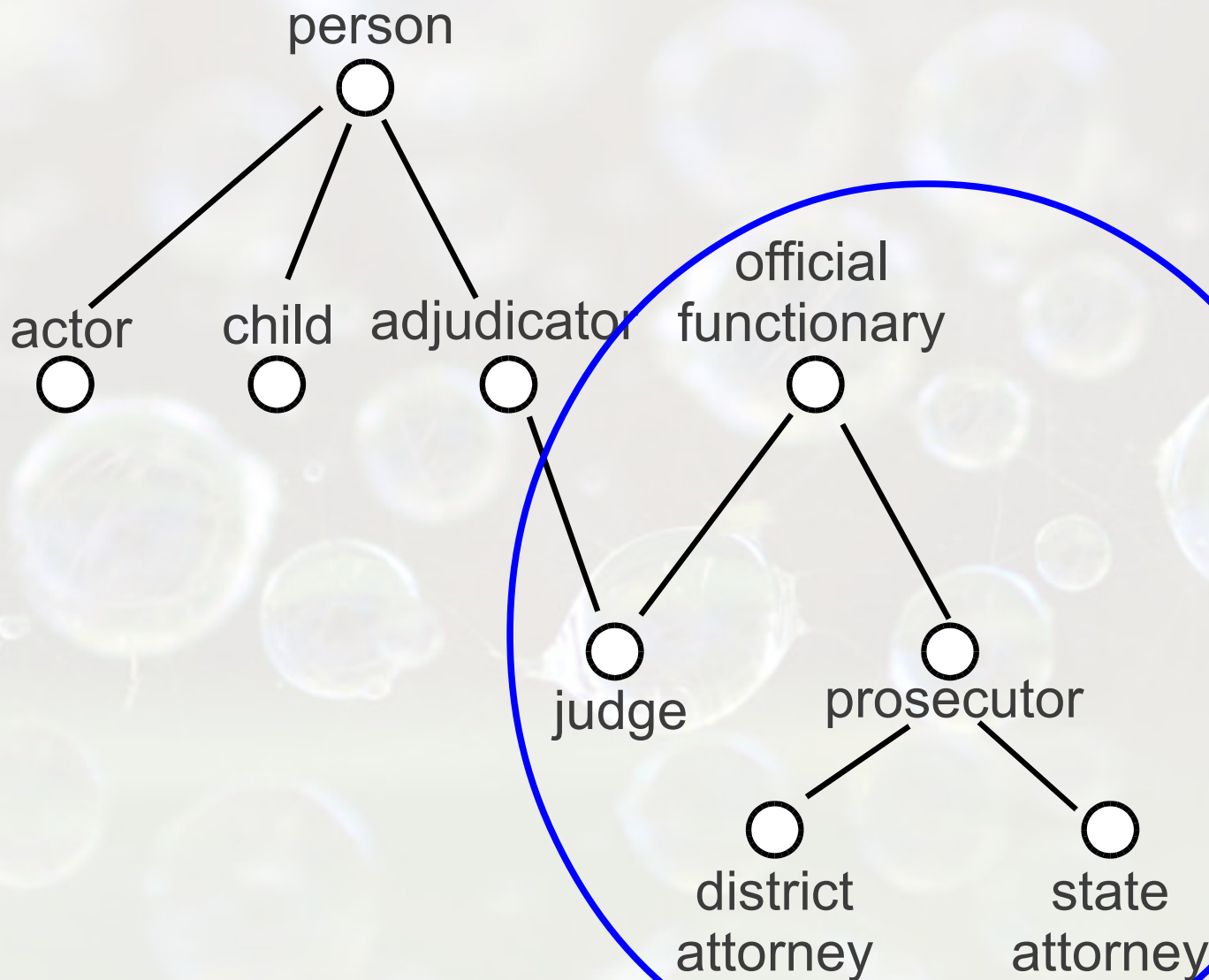# Our approach

7 distinct morphological descriptions:

- genus Varanus
- *Varanus gouldii*
- *Varanus indicus*
- *Varanus prasinus*
- *Varanus salvator*
- *Varanus spiny*
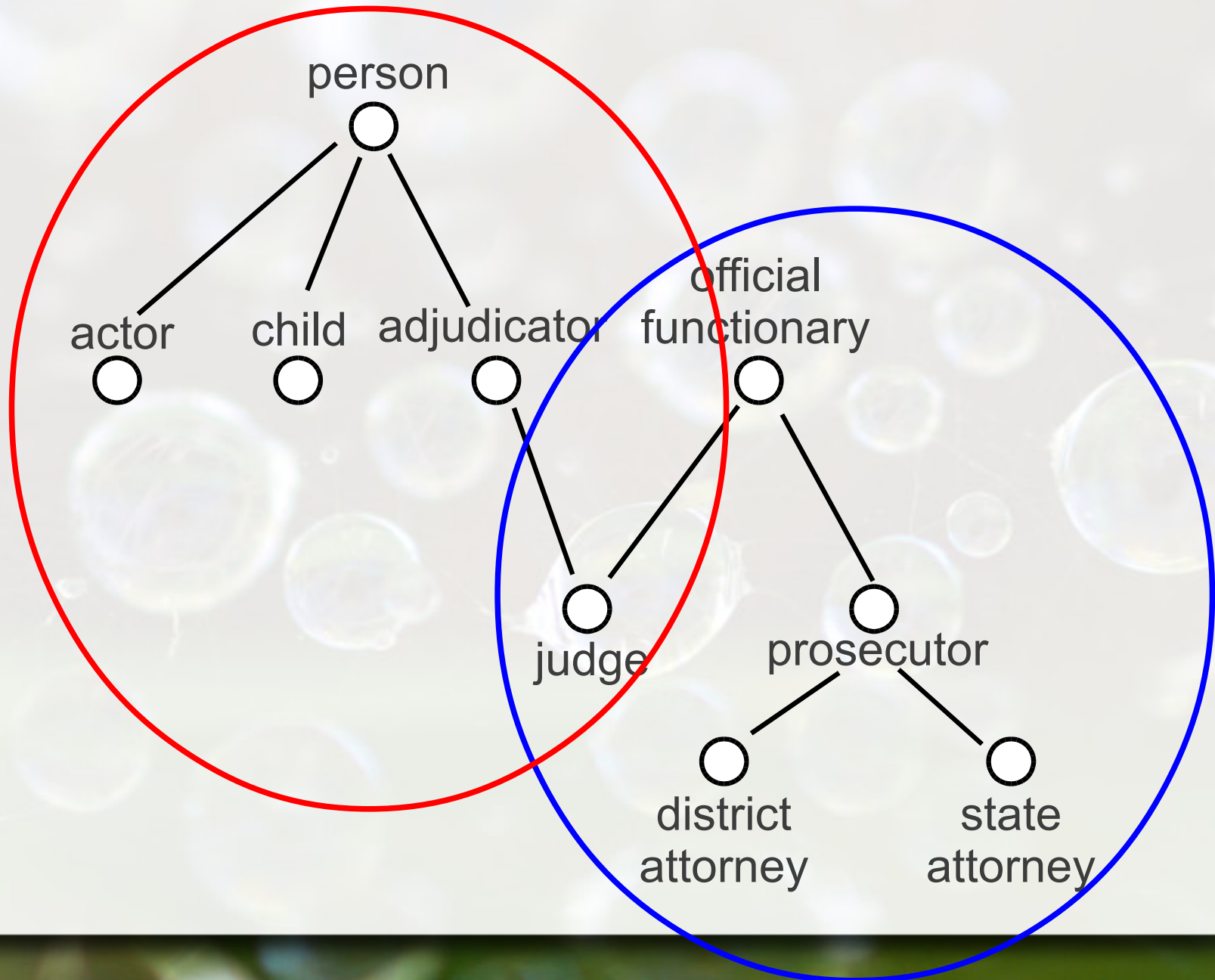- *Varanus timorensis*

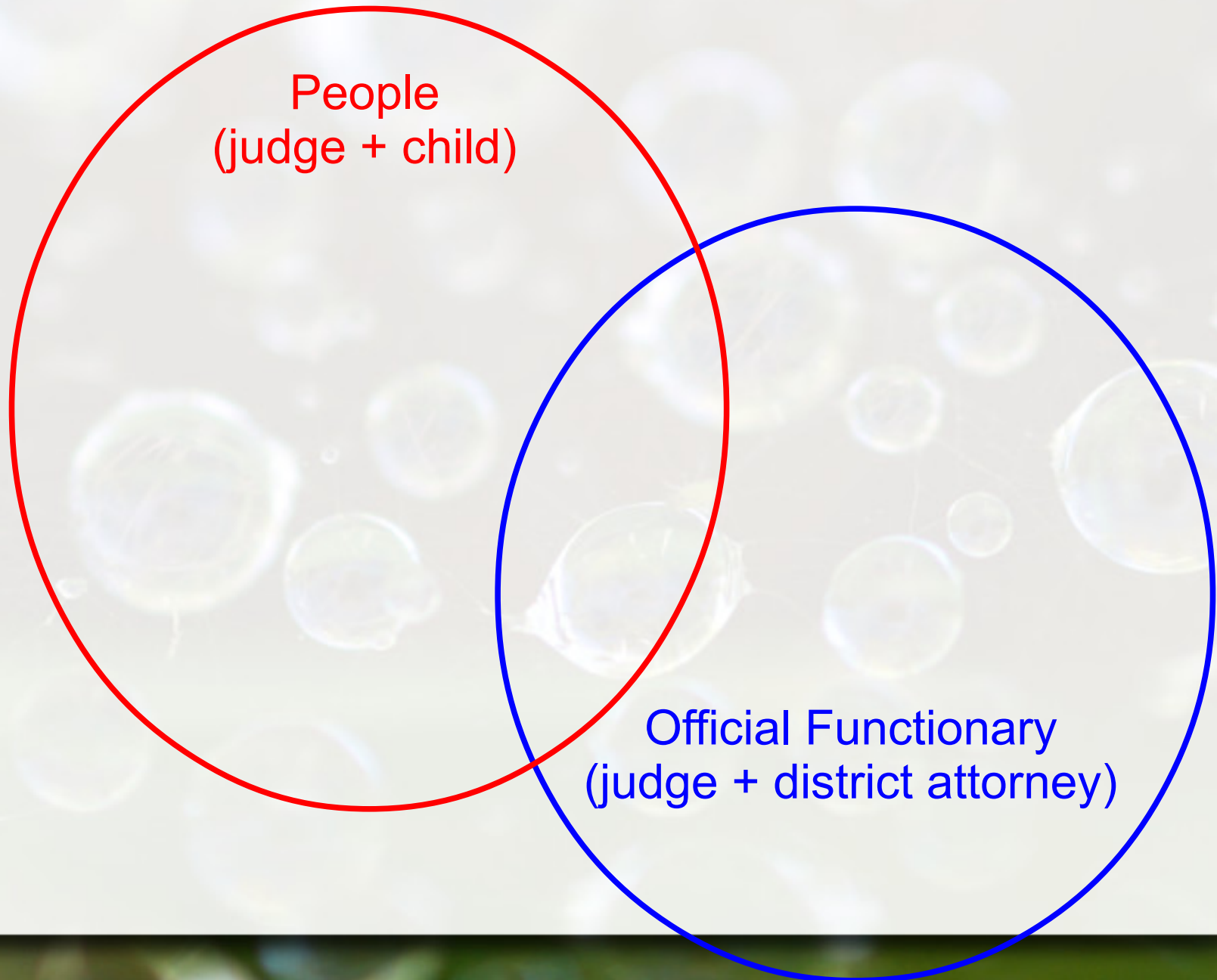# Evaluating Similarity

# Size of the Shared Universe

# Size of the Shared Universe

# Size of the Shared Universe

person

actor    child    adjudicator    official
functionary

judge    prosecutor

district
attorney

state
attorney

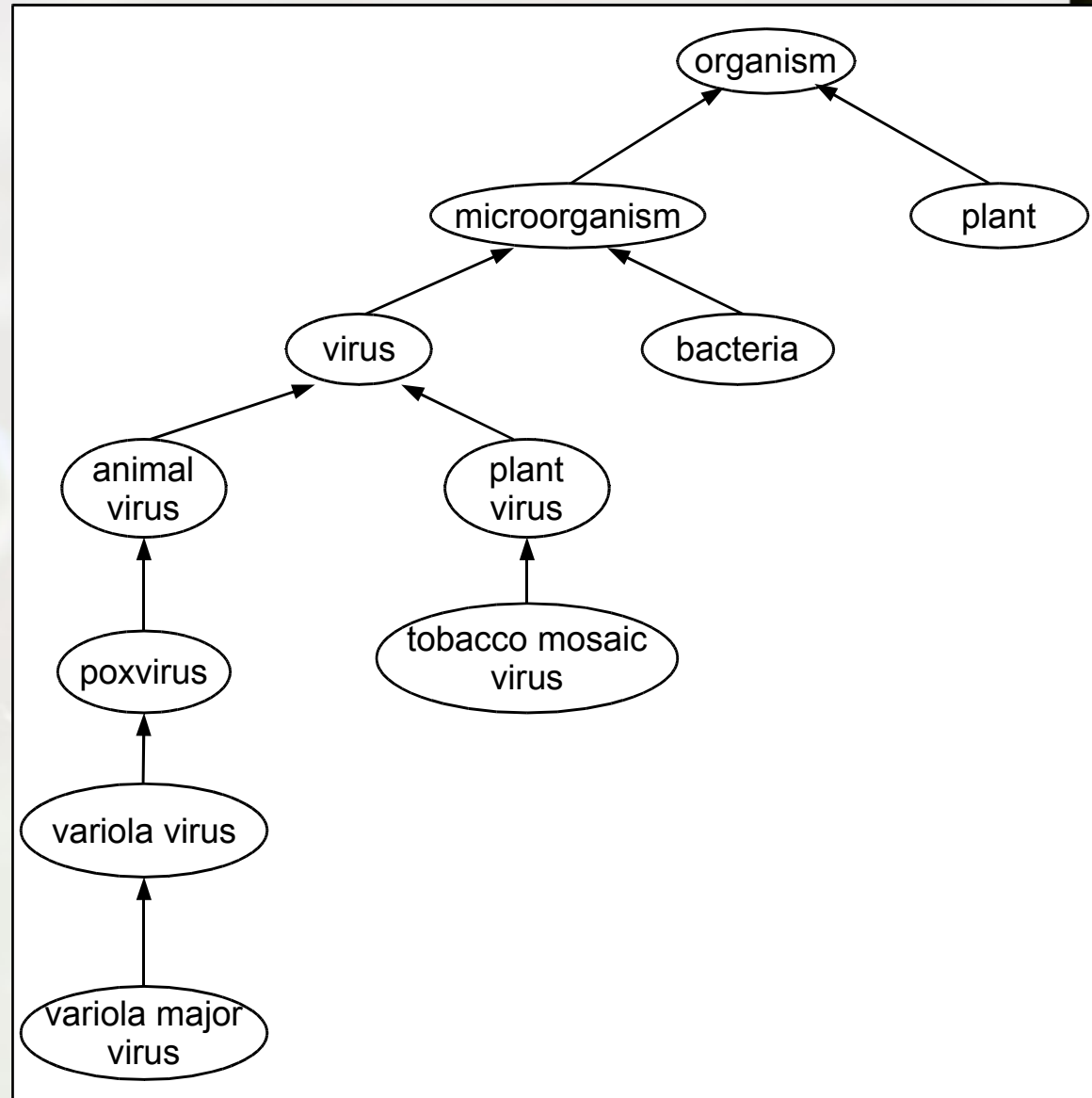# Size of the Shared Universe



People
(judge + child)

Official Functionary
(judge + district attorney)

# Similarity

- What is more similar of a **tobacco mosaic virus?**
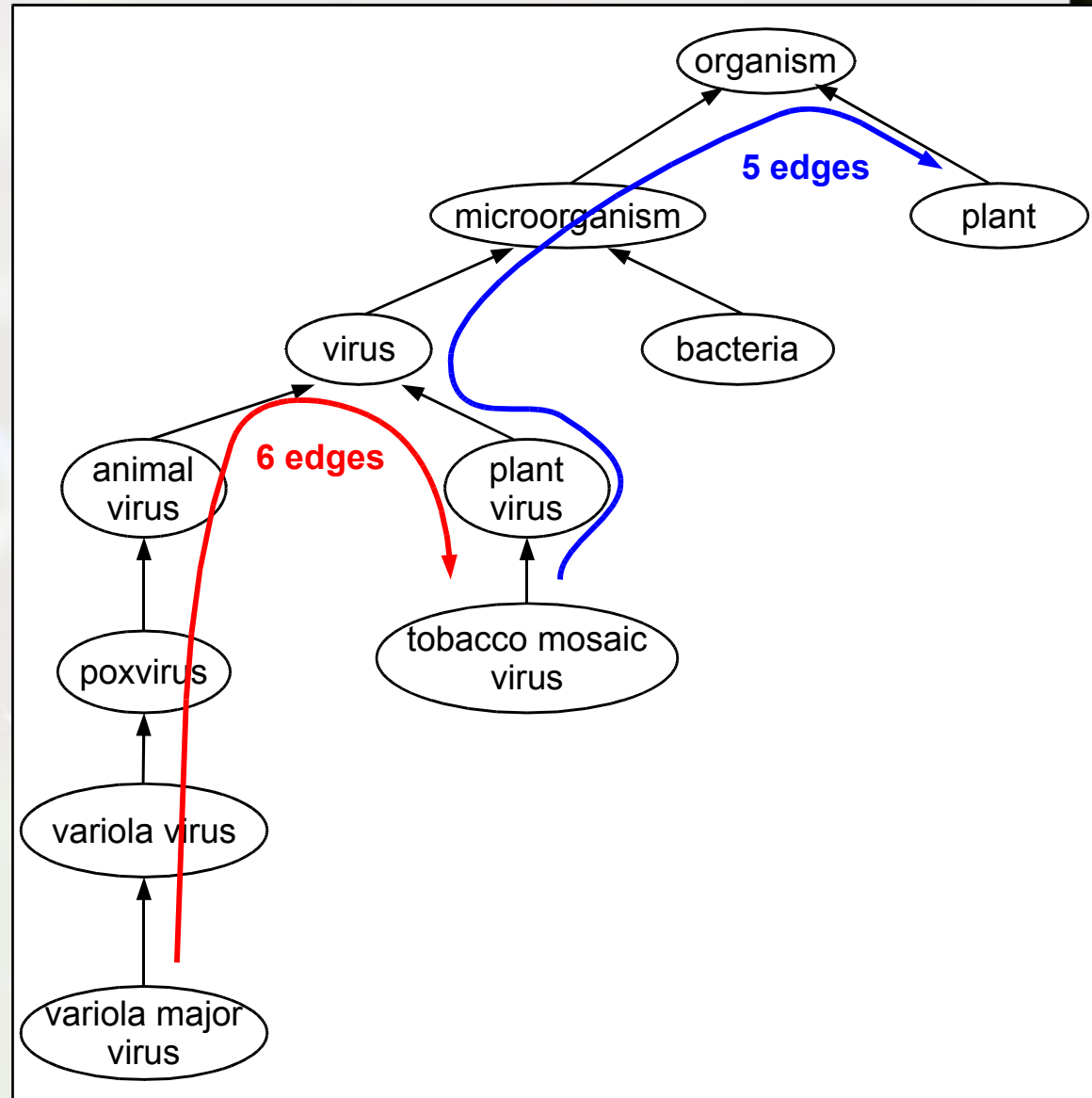
  a) variola major virus

  b) plant

# Similarity
## Minimum path length

- What is more similar of a tobacco mosaic virus?
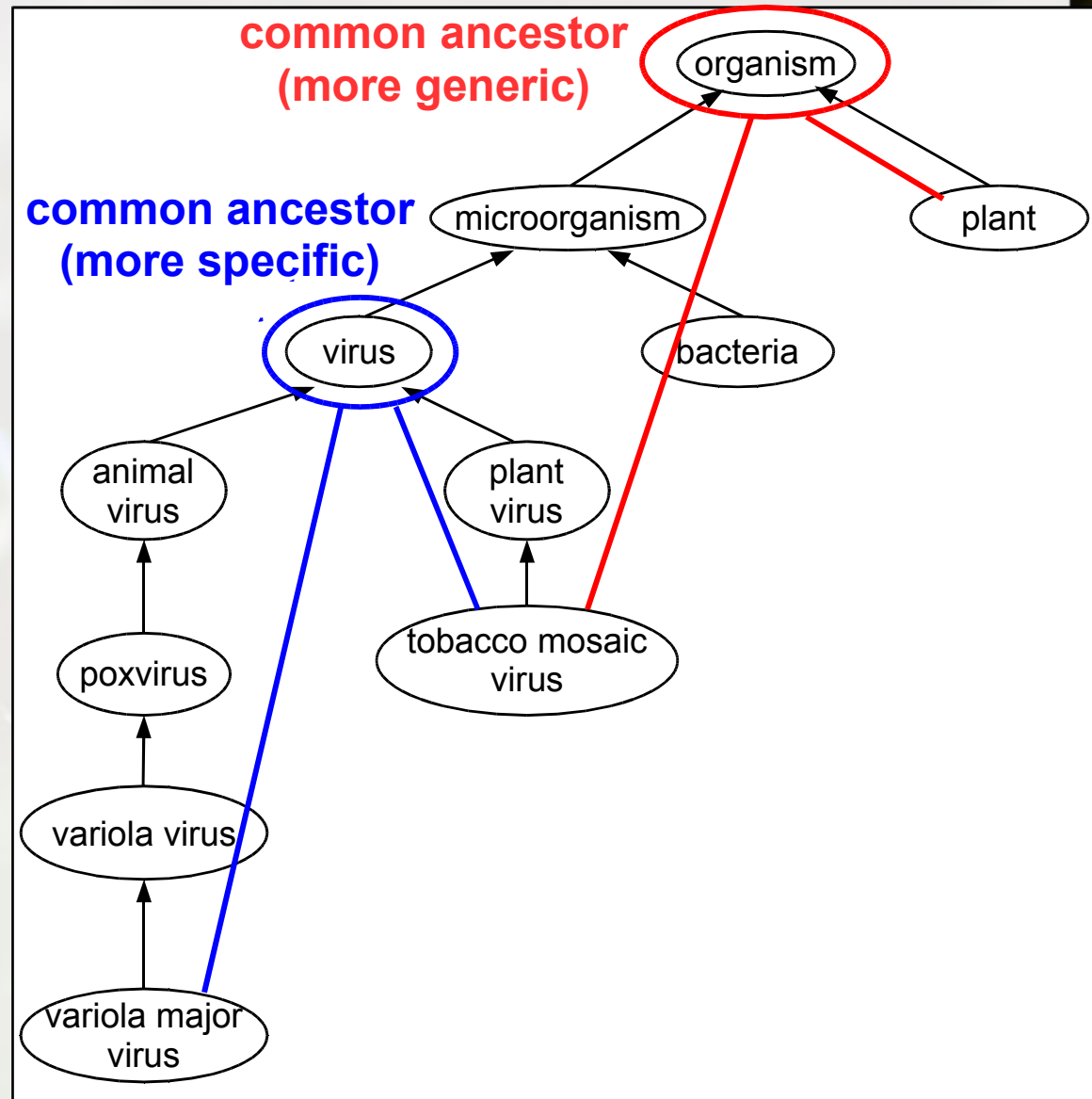
   a) variola major virus

   b) plant (?)

# Similarity
## More specialized common ancestor

- What is more similar of a tobacco mosaic virus?
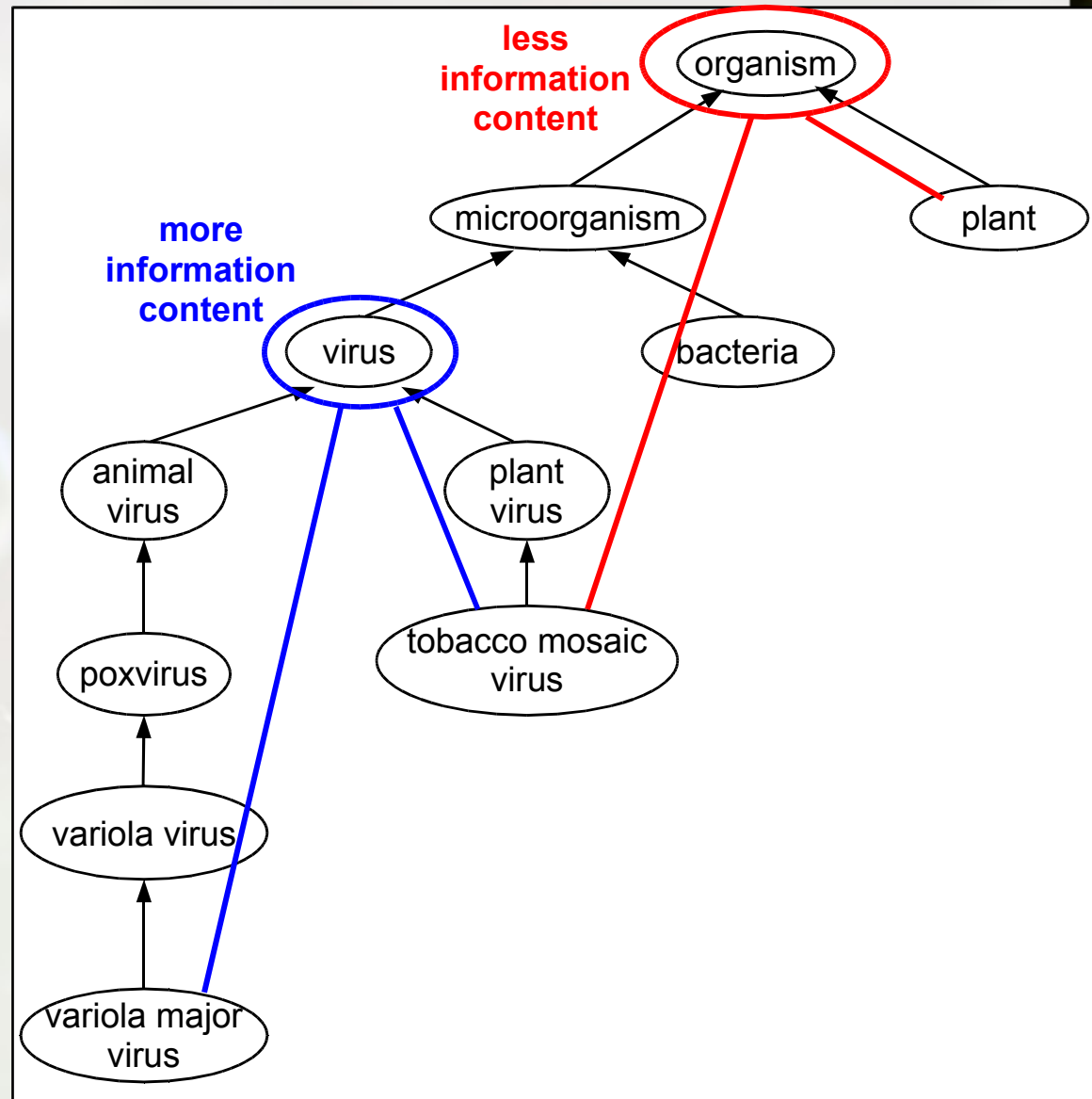
  a) variola major virus

  b) plant

# Similarity
## Maximum information content

- What is more similar of a tobacco mosaic virus?
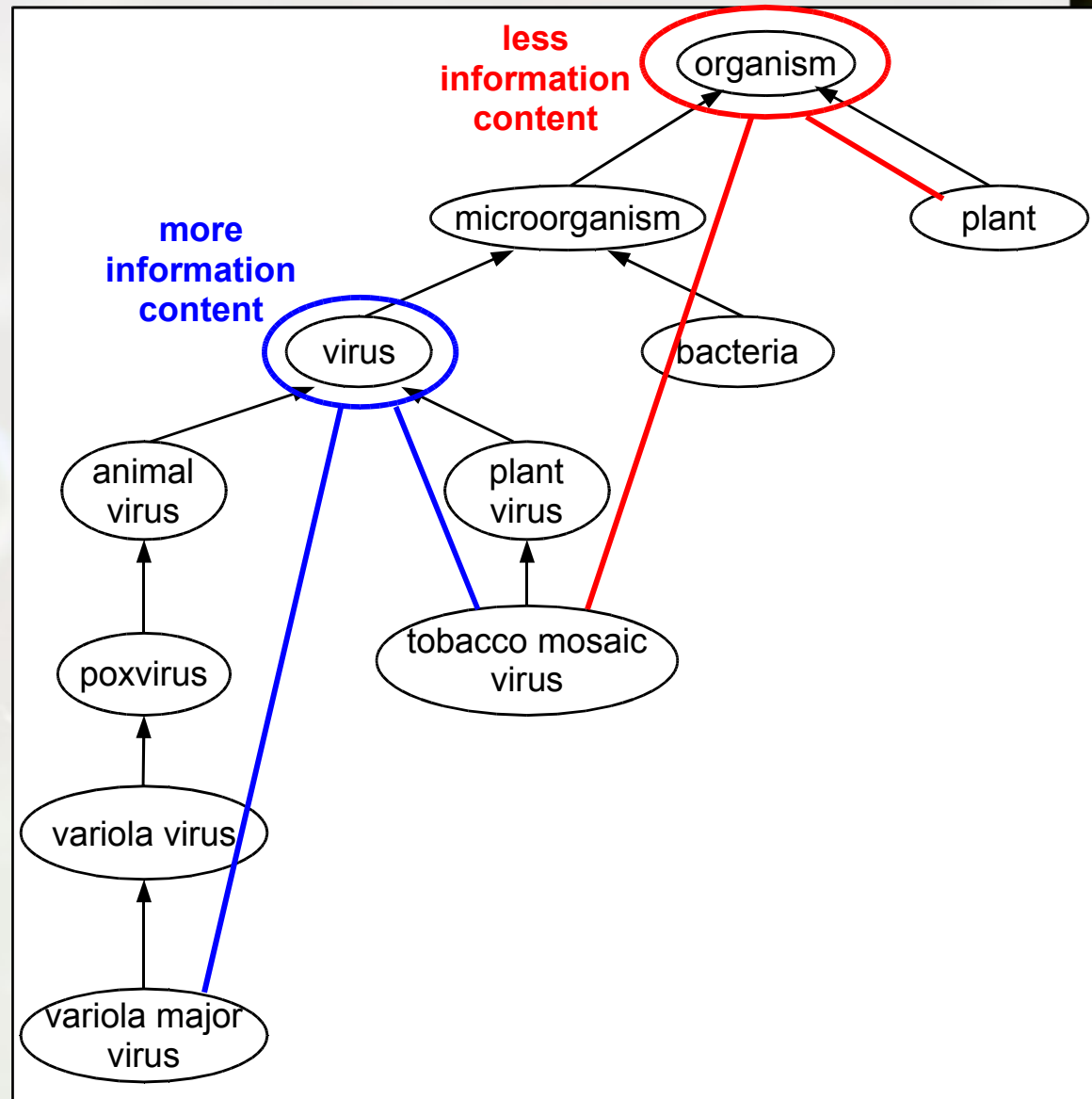
  a) variola major virus

  b) plant

# Similarity
## Maximum information content

- What is more similar of a tobacco mosaic virus?

  a) variola major virus

  b) plant

# Graph Databases

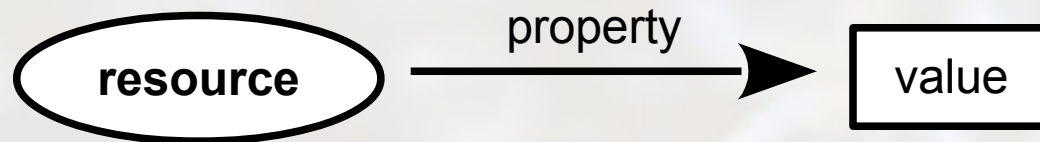# Graph Model

Graph Model
# RDF Graph

# RDF Graph
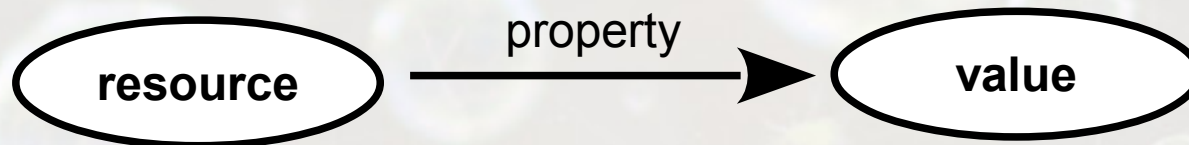
- Triple model: resource, property, value

# RDF Graph

- Triple model: resource, property, value



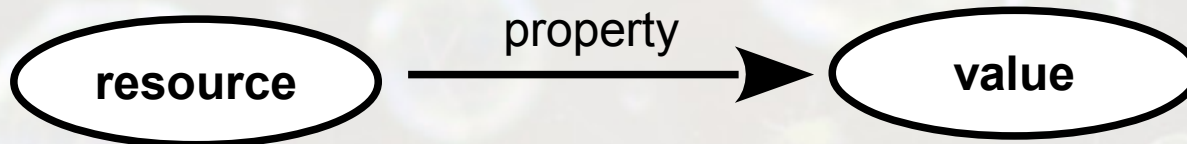- The value can be another resource

# RDF Graph

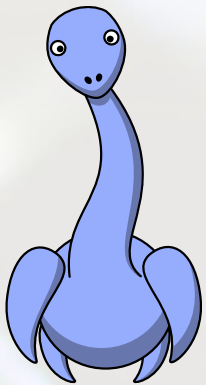- Triple model: resource, property, value



- The value can be another resource



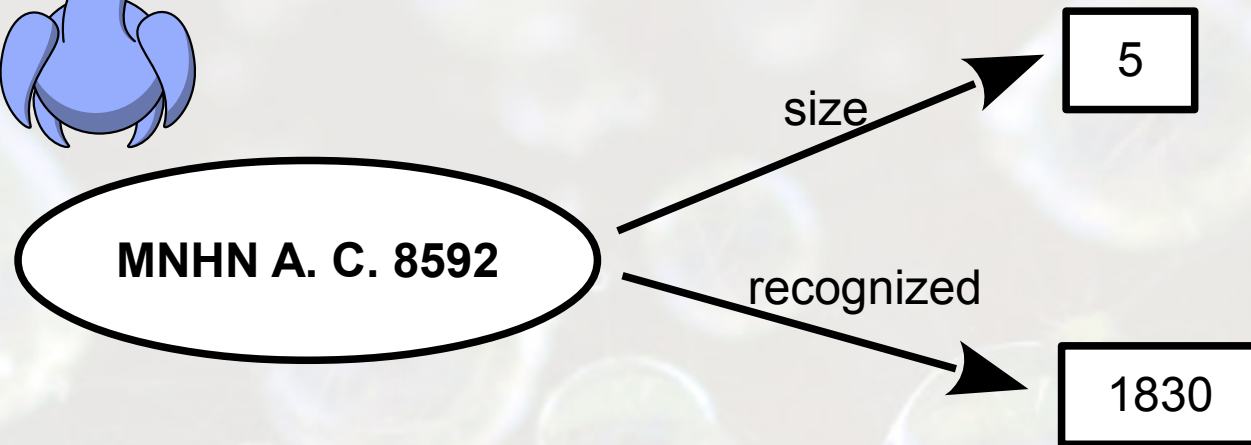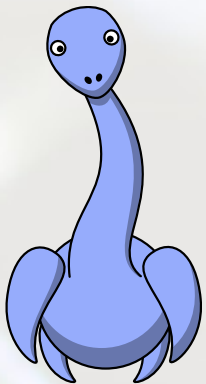- ... described by a property, value

# Plesiosaurus dolichodeirus in RDF
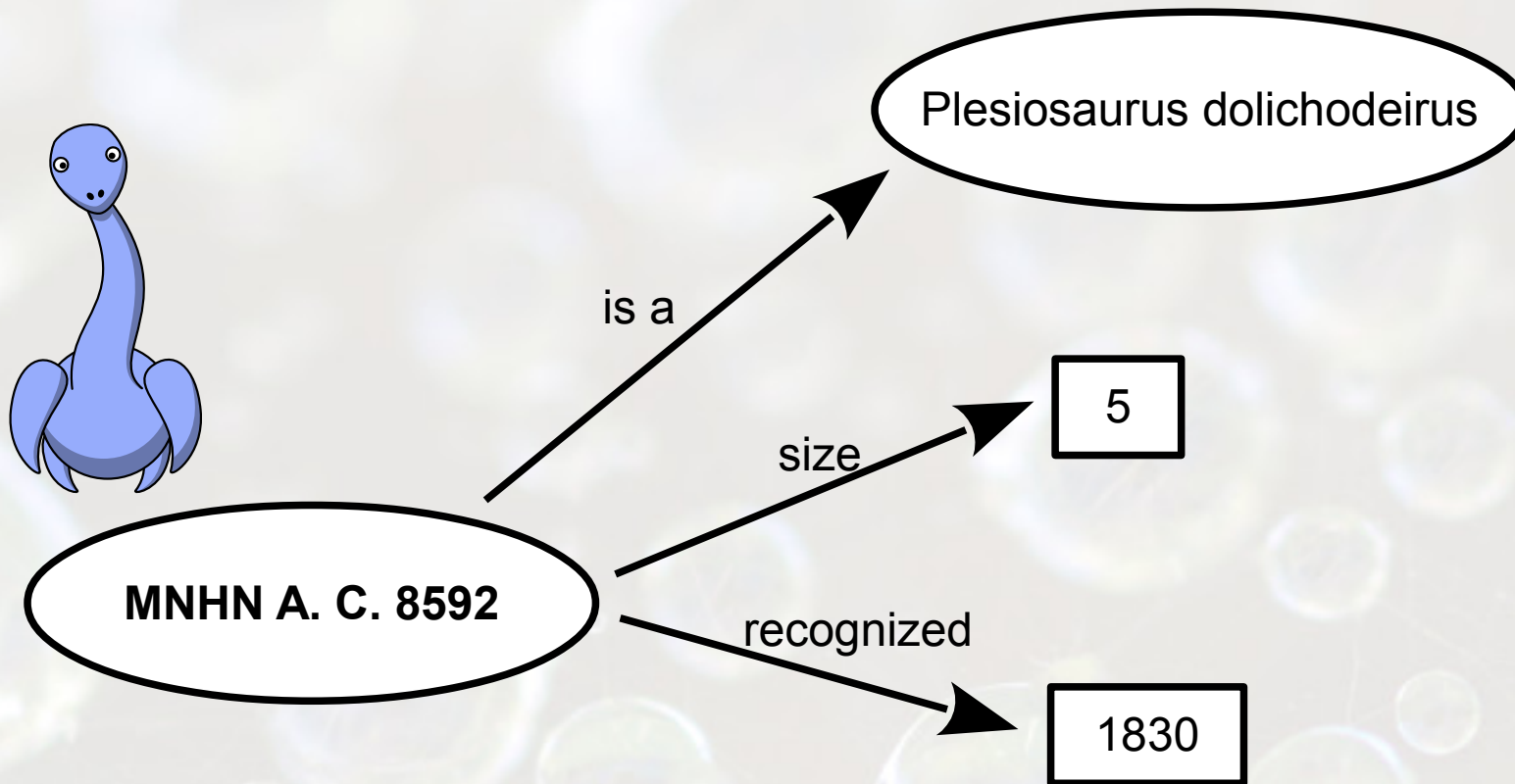


- Triple: <MNHN A. C. 8592>, size, 5

# Plesiosaurus dolichodeirus in RDF



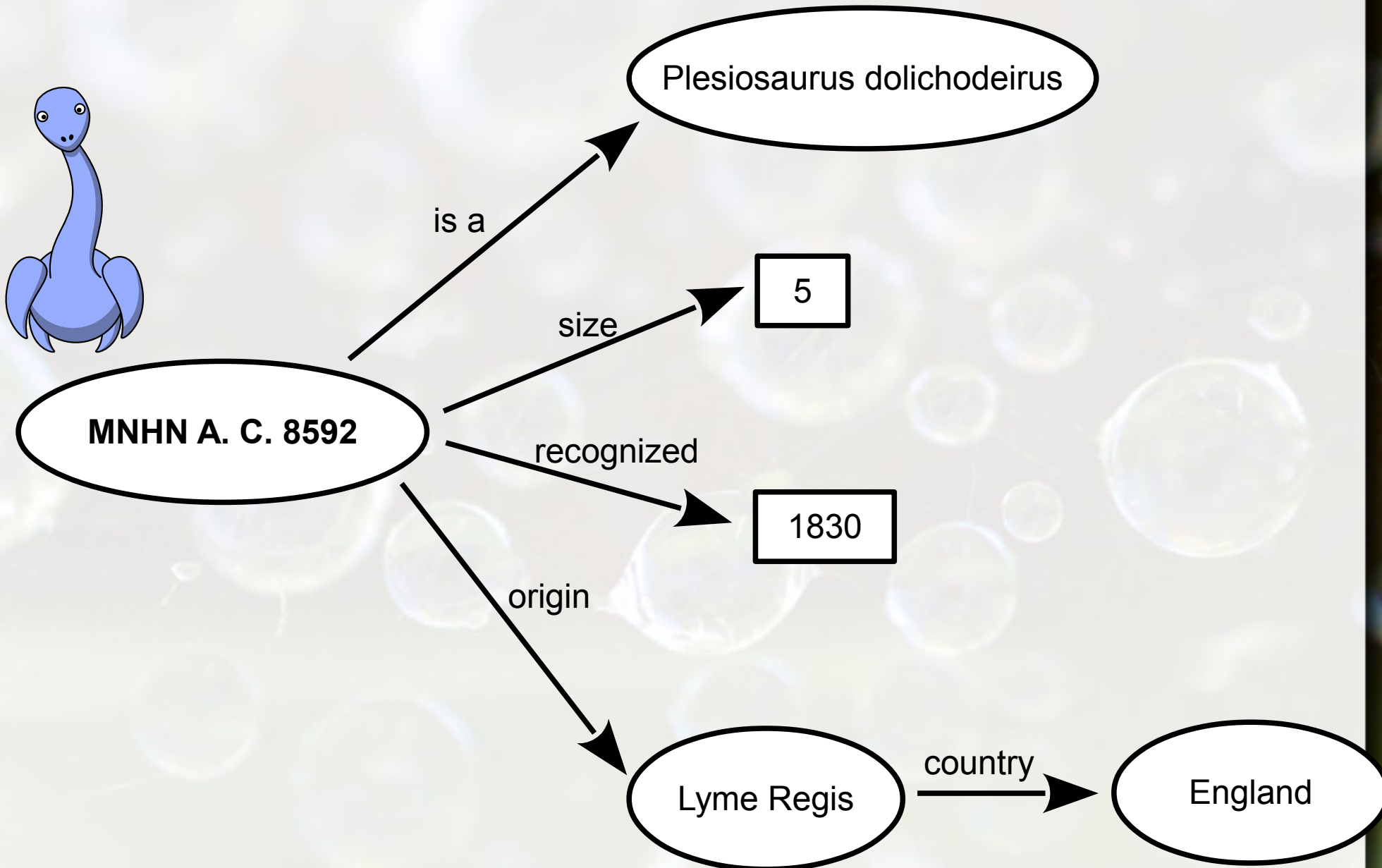- Triple: <MNHN A. C. 8592>, size, 5
- Triple: <MNHN A. C. 8592>, recognized, 1830
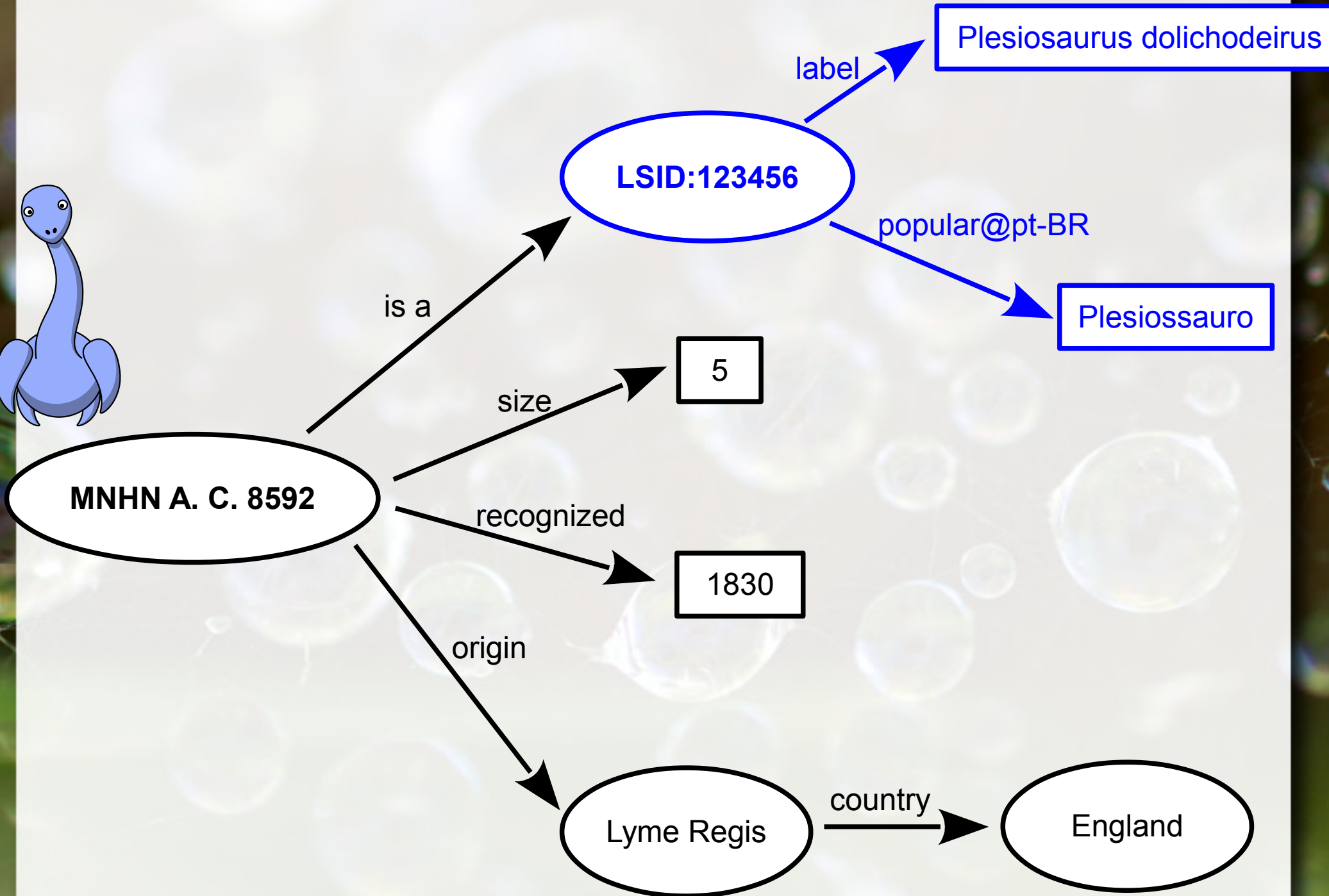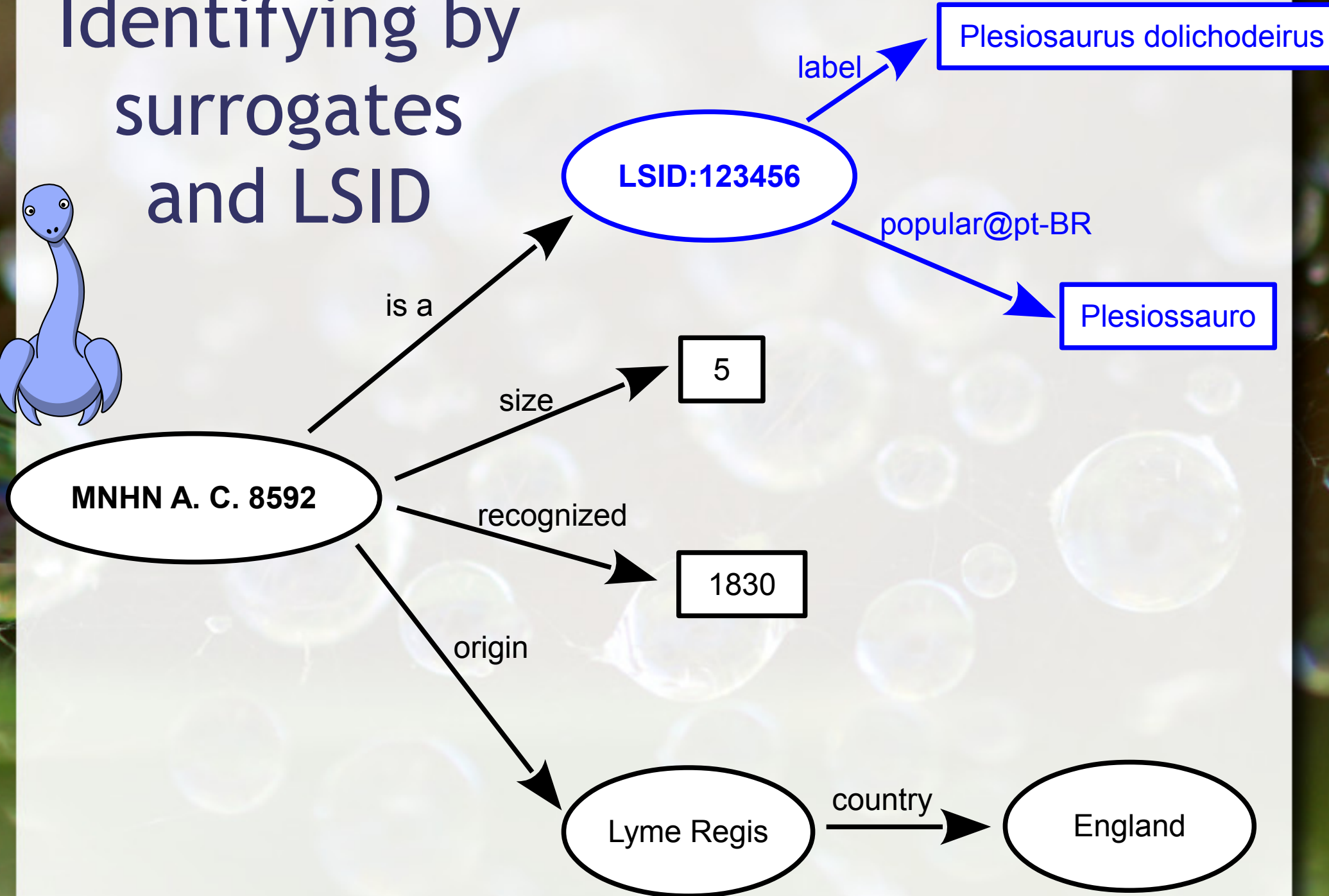
# Plesiosaurus dolichodeirus in RDF



- <MNHN A. C. 8592>, is_a, <Plesiosaurus dolichodeirus>
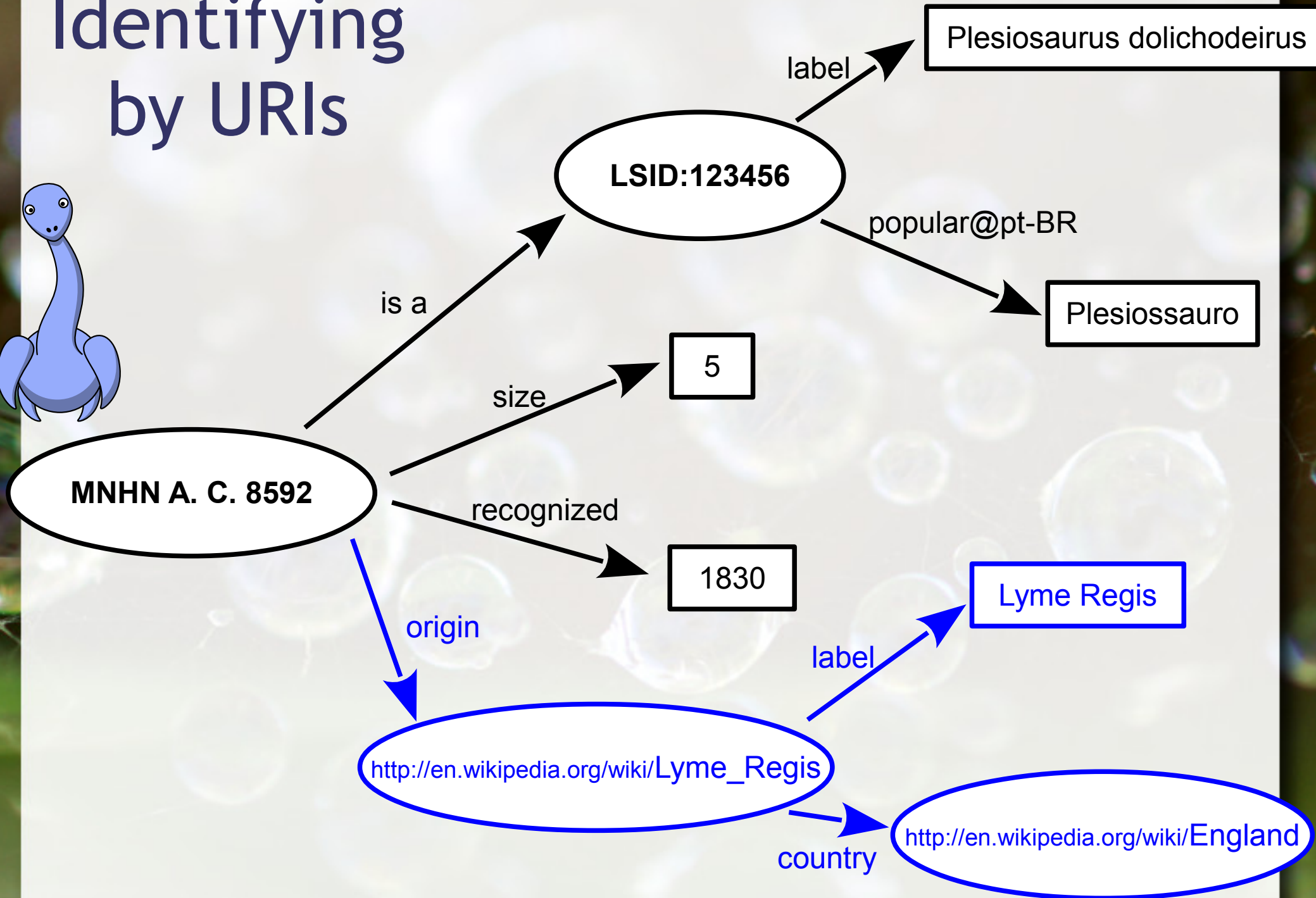
# Plesiosaurus dolichodeirus in RDF

# Identifying by surrogates and LSID

LSID:123456

label → Plesiosaurus dolichodeirus

popular@pt-BR → Plesiossauro

MNHN A. C. 8592

is a → LSID:123456

size → 5

recognized → 1830

origin → Lyme Regis

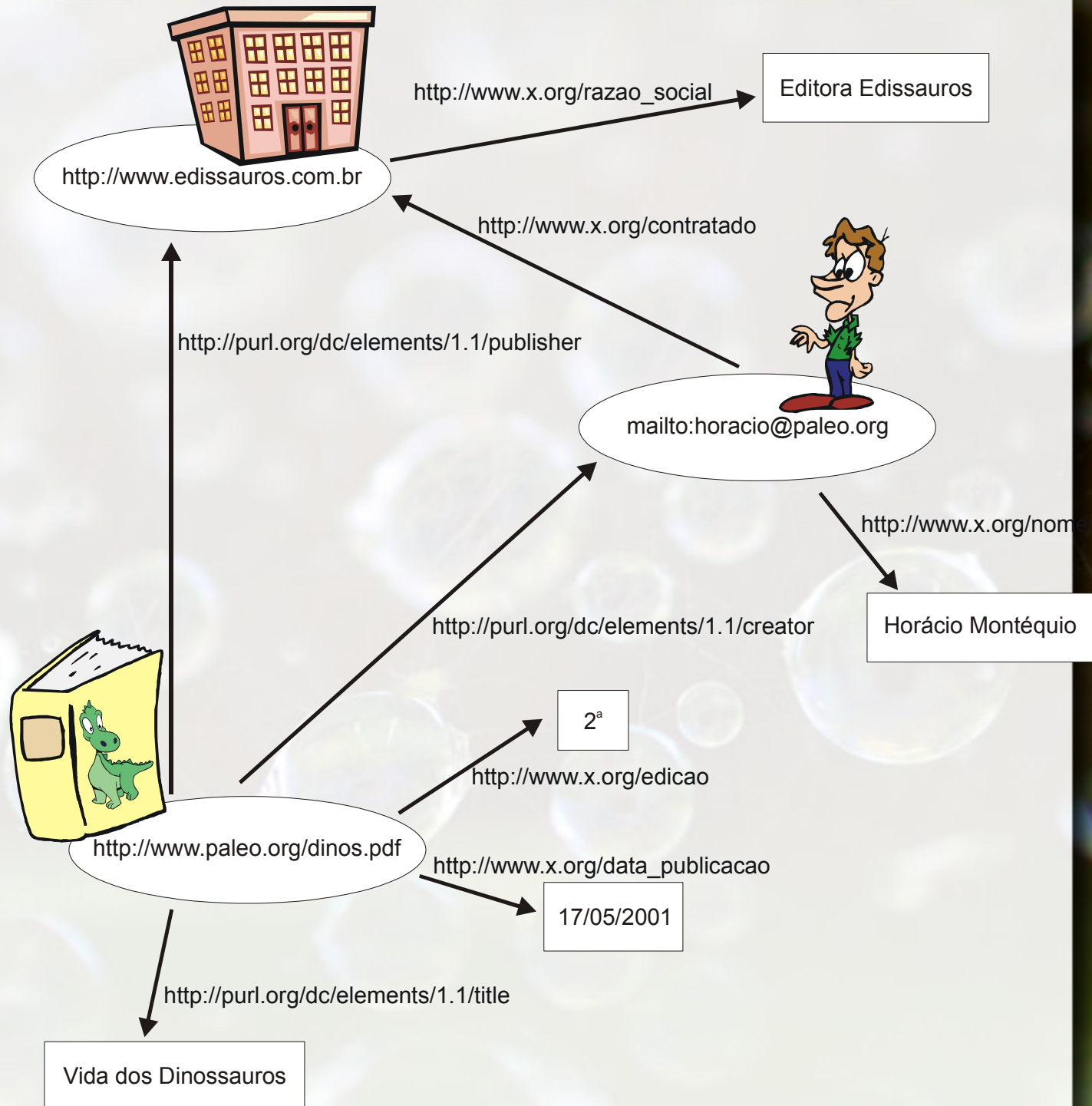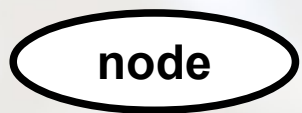country → England

Identifying by URIs

# Short URIs (namespaces)

RDF Graph
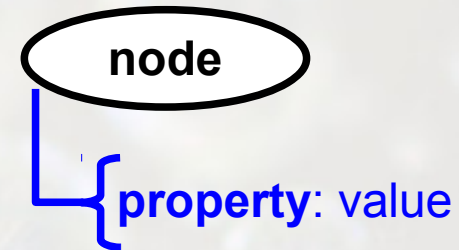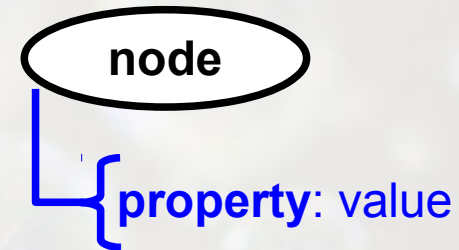
Graph Model
# Property Graph

# Property Graph

- (node)    ( **node** )

# Property Graph

- (node)

  - Nodes can have properties
    (node { property: value })
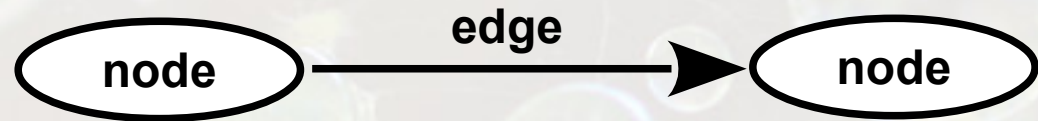
node

node

property: value

# Property Graph

- **(node)**

- Nodes can have properties
  (node { property: value })

- **(node)-[edge]->(node)**

# Property Graph

- **(node)**     node

  - Nodes can have properties
    (node { property: value })

    node
    property: value

- **(node)-[edge]->(node)**

    node —— edge ——> node

  - Edges can have properties
    (node)-[edge {property:value}]->(node)

    node —— edge ——> node
    property: value

# Plesiosaurus in a Property Graph



MNHN A. C. 8592

# Plesiosaurus in a Property Graph

# Plesiosaurus in a Property Graph

# LSID Surrogate



**LSID:123456**

**label**: Plesiosaurus dolichodeirus
**popular@pt-BR**: Plesiossauro

is a

**MNHN A. C. 8592**

**size**: 5
**recognized**: 1830

origin

Lyme Regis

country

England

# URI Identifiers

# Property in the `origin` edge

LSID:123456
{ **label**: Plesiosaurus dolichodeirus
**popular@pt-BR**: Plesiossauro

is a

MNHN A. C. 8592

{ **size**: 5
**recognized**: 1830

**discovered**: 1824

origin

wiki:Lyme_Regis — country → wiki:England

{ **label**: Lyme Regis

{ **label**: England

# Back to the Property Graph

# Two properties for `origin`

MNHN A. C. 8592
- is a → LSID:123456
  - label → Plesiosaurus dolichodeirus
  - popular: pt-BR → Plesiossauro
- size → 5
- recognized → 1830
- origin → ( )
  - place → wiki:Lyme_Regis
  - discovered → 1824

# References

- GOMES-JR, L. C.; JENSEN, R.; SANTANCHÈ, A. (2013) Towards Query Model Integration: Topology-aware, IR-inspired Metrics for Declarative Graph Querying, II Workshop on Querying Graph Structured Data 2013, Genoa.
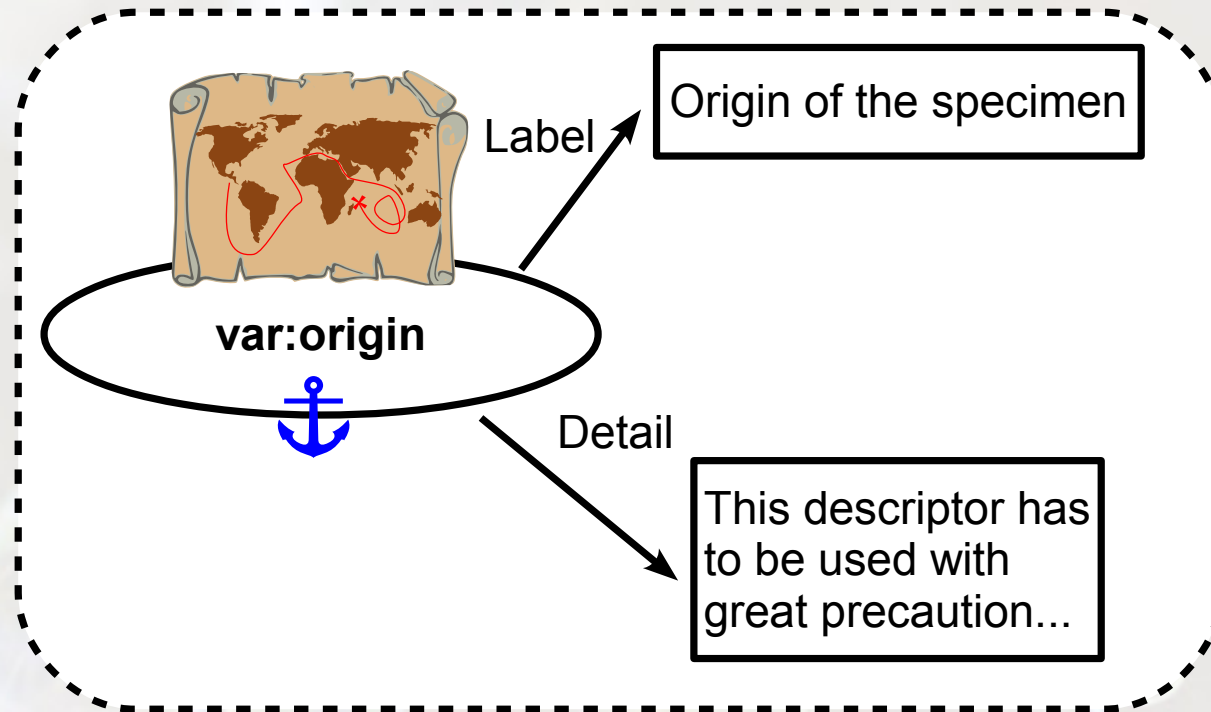
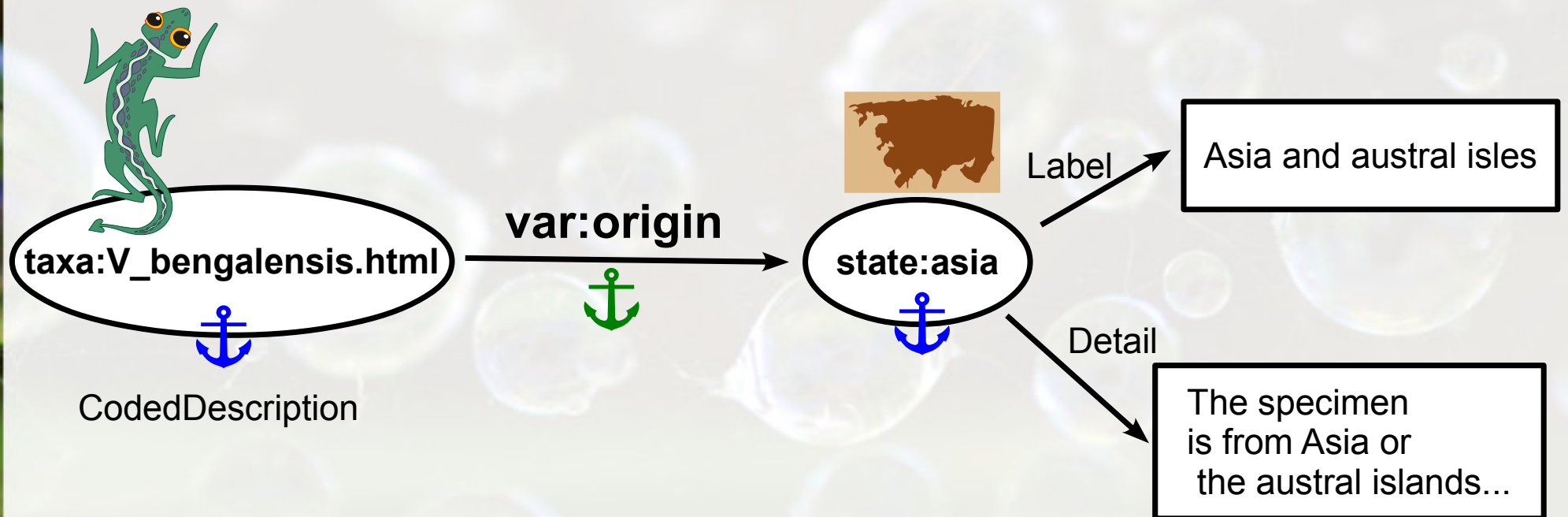# André Santanchè

http://www.ic.unicamp.br/~santanche
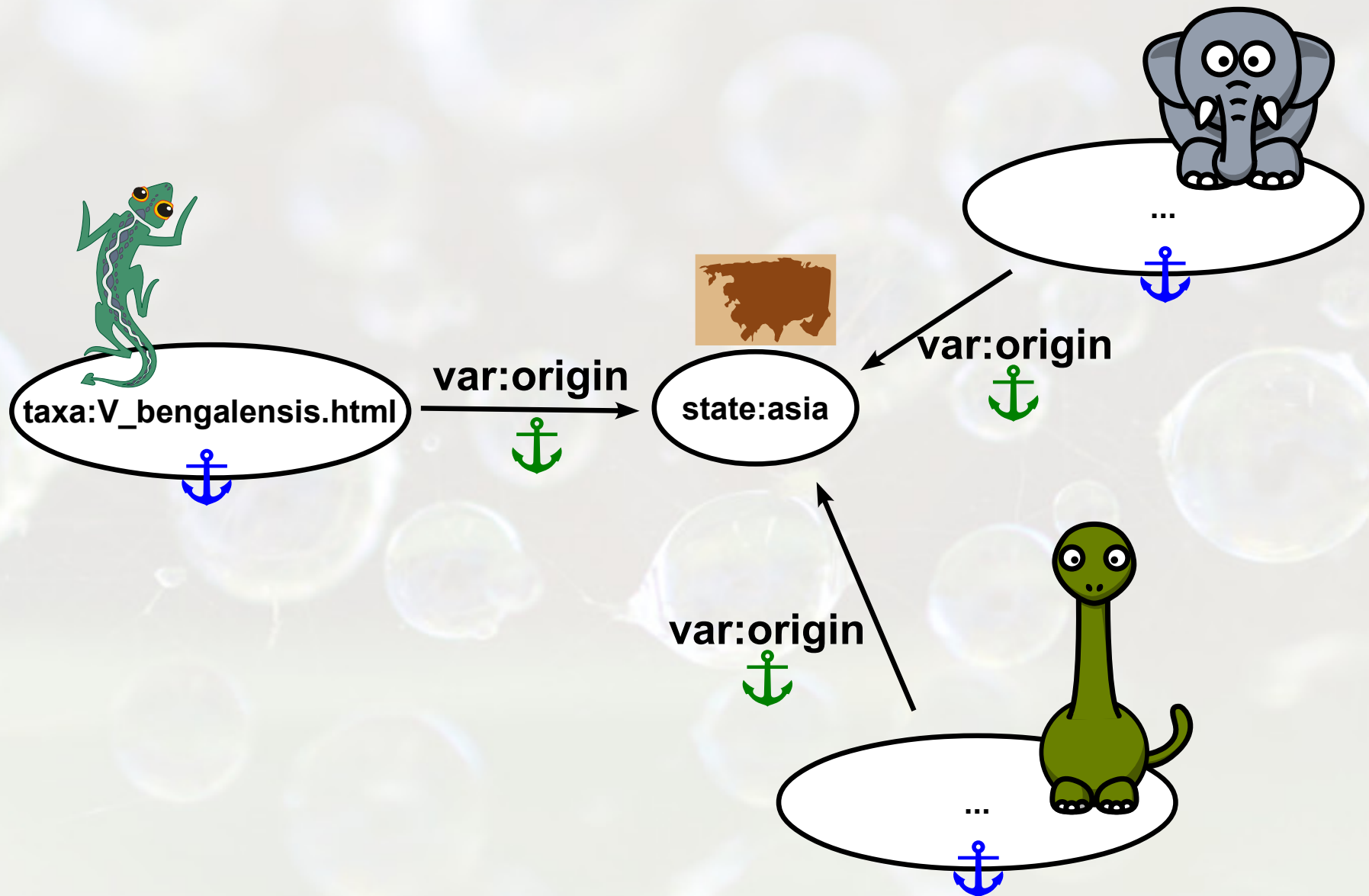
# Property

# Using the Property



taxa:V_bengalensis.html — var:origin → state:asia

CodedDescription

Label → Asia and austral isles

Detail → The specimen is from Asia or the austral islands...

# Connected Graph

# Origin in GeoNames

# Geo Tree

Asia

http://www.geonames.org/6255147/

gn:parentFeature

India

http://www.geonames.org/1269750/

# Hypergraph



(Jaudete Daltio, 2013)

# Licença
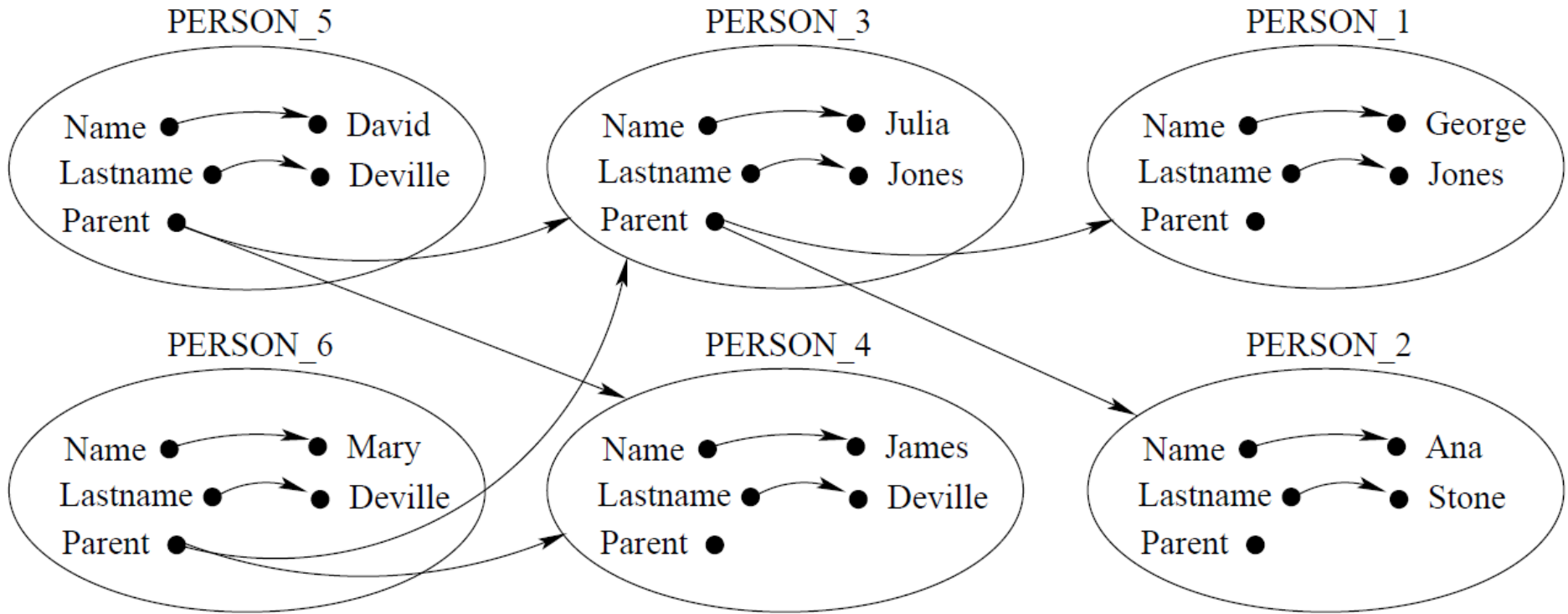
- Estes slides são concedidos sob uma Licença Creative Commons. Sob as seguintes condições: Atribuição, Uso Não-Comercial e Compartilhamento pela mesma Licença.

- Mais detalhes sobre a referida licença Creative Commons veja no link: http://creativecommons.org/licenses/by-nc-sa/3.0/

- Fotografia de capa e fundos: web-drops por Jeremy Hiebert [http://www.flickr.com/photos/jeremyhiebert/] dispinível em http://www.flickr.com/photos/jeremyhiebert/6081389428/