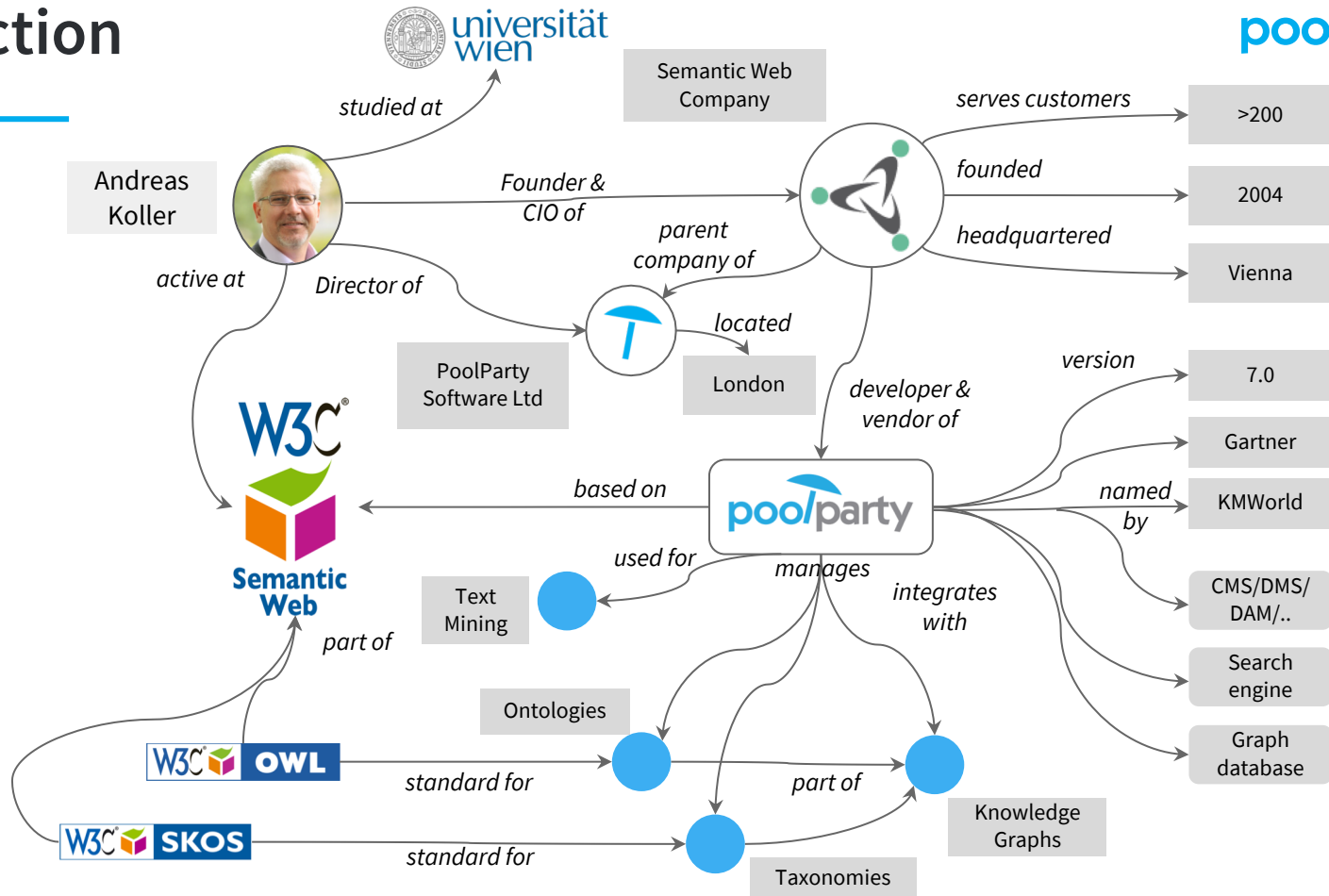# Deep Text Analytics

## 2019 NKOS Workshop

How to extract hidden information and 'aboutness' from text using SKOS, ontologies, corpus analysis and linked data

**Andreas Koller**

Co-founder and CIO
Semantic Web Company

# Introduction

# At a glance

## PoolParty Semantic Suite

▶ PoolParty is the most complete and most secure **semantic middleware** on the global market.

▶ PoolParty integrates with leading **graph database** technologies, NoSQL stores, and search engines.

▶ PoolParty was launched in **2009**.

## Semantic Web Company

▶ SWC is **developer, licensor and vendor** of PoolParty Semantic Suite.

▶ Established in 2004, SWC has pioneered the areas of **Enterprise Semantics & Semantic AI**.

▶ SWC is dedicated to use **W3C standards** like SKOS, OWL, SPARQL, etc.

# Why Knowledge Graphs?

# WHAT'S THE PROBLEM?

Many Challenges in (Enterprise) Information Management & AI

# "Search" is still about documents only

… and in enterprises it's painful

poolparty®

wind farms OR wind parks OR wind power plants OR wind power stations

# Missing Context

# Context-free data models and User-agnostic

Do machines understand user intent? Do they have enough context?

French SUV

Which **Sport-utility vehicle** from **France** provides **enough space** for my family with 3 kids?

1. **Intent recognition**
2. **Entity linking**
3. **Background knowledge**

The Peugeot 5008 breaks new ground as a large SUV with many features.

| Car | Loadspace | Max no. of seats |
|-----|-----------|------------------|
| KIA Sorento | 550 litres | 7 |
| Peugeot 5008 | 823 litres | 7 |
| BMW X3 | 550 litres | 5 |

# Background knowledge (context) is key

poolparty®

## Support complex Q&A:

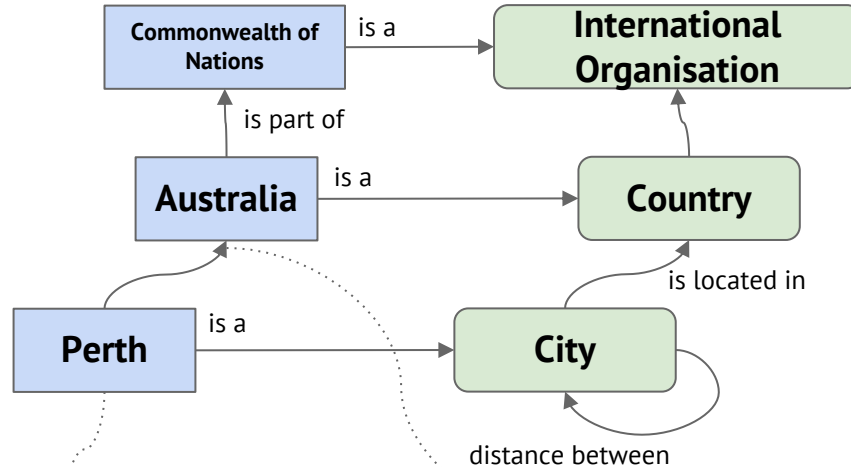Which cities located in the Commonwealth of Nations have a population of more than 2 mio. people?

## Avoid illogical answers:

How far am I away from Perth, Australia?

WESTERN AUSTRALIA

Au

1,195 miles

Distance between Perth and Australia

Commonwealth of Nations — is a → International Organisation

Commonwealth of Nations ← is part of — Australia

Australia — is a → Country

Country ← is located in — City

Perth — is a → City

City — distance between

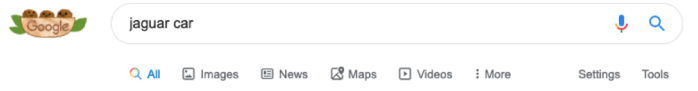Perth is one of the most isolated major cities in the world, with a population of 2,022,044 living in Greater Perth.

Australia is a member of the OECD, United Nations, G20, ANZUS, and the World Trade Organisation.

# KNOWLEDGE GRAPHS

Bring your metadata to the next level!

# Google Knowledge Graph

# Example **Knowledge Graph**

- ▶ URIs
- ▶ Concepts
- ▶ Labels
- ▶ Relations



Cat

Felid

prefLabel

altLabel

http://localhost/Jaguars/4

Keystone species

prefLabel

http://localhost/Jaguars/7

hasBroader

hasBroader

Jaguar

prefLabel

http://localhost/Jaguars/5

altLabel

Panthera onca

Car

prefLabel

http://localhost/Jaguars/1

hasBroader

Disambiguation

Jaguar

prefLabel

http://localhost/Jaguars/6

Demo

# 'Things' but no strings:
## Using a 'Semantic Knowledge Graph'



Retina

*prefLabel*

http://www.my.com/taxonomy/62346723

*image*

http://www.my.com/images/90546089

Home → Medical Encyclopedia → Retina

### Retina

The retina is the light–sensitive layer of tissue at the back of the eyeball. Images that come through the eye's lens are focused on the retina. The retina then converts these images to electric signals and sends them along the optic nerve to the brain.

The retina usually looks red or orange because there are many blood vessels right behind it. An ophthalmoscope allows a health care provider to see through your pupil and lens to the retina. Sometimes photos or special scans of the retina can show things that the provider cannot see just by looking at the retina through the ophthalmoscope. If other eye problems block the provider's view of the retina, ultrasound can be used.

Anyone who experiences these vision problems should get a retinal examination:

- Changes in sharpness of vision
- Loss of color perception
- Flashes of light or floaters
- Distorted vision (straight lines look wavy)

Watch this video about:
Retina

http://www.my.com/taxonomy/97345854

*prefLabel*

Funduscope

Ophthalmoscope

*altLabel*

*has broader*

http://www.mycom.com/taxonomy/4543567

*prefLabel*

Diagnostic Equipment

# Semantic Enrichment

**Taxonomy & Ontology Server**

- ▸ Graph-based annotation →
  **Entity linking**

- ▸ Machine-learning-based annotation →
  **Named entity recognition**

- ▸ Machine-learning based classification →
  **Document Classification**

- ▸ Annotation based on
  **Regular expressions**

Taxonomy tree:
- ▼ Sporting Goods (12)
  - ▸ Air Sports (3)
  - ▸ Combat Sports (7)
  - ▸ Dancing (1)
  - ▸ Exercise & Fitness (22)
  - ▸ Gymnastics (8)
  - ▸ Indoor Games (8)
  - ▸ Jumping (5)
  - ▸ Outdoor Recreation (25)
  - ▸ Racquet Sports (7)
  - ▸ Team Sports (17)
  - ▸ Water Sports (14)
  - ▼ Winter Sports (8)
    - Bobsledding (0)
    - Luge (0)
    - ▸ Skiing (8)
    - ▸ Sledding (3)
    - ▸ Snow Boots (1)
    - ▸ Snowboarding (5)
    - Snowmobiling (0)
    - ▸ Snowshoeing (2)
- ▸ Toys & Games (5)
- ▸ Vehicles & Parts (1)

**Entity Extraction & Text Mining**

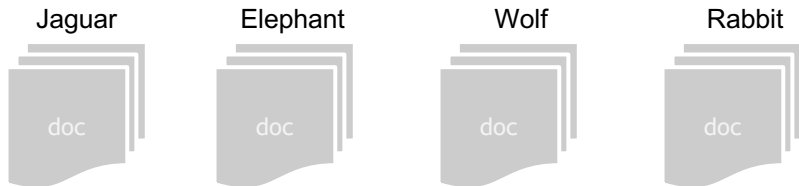Men's SuperBoots are warm, waterproof slip-on Snow boots. They are perfect for quick trips out to grab the paper or all-day escapades in the snow. They're built with all-over insulation and rich, durable, waterproof leather, with a seam-sealed construction for bonus protection from the wet and rugged traction to keep you stable on wet, slippery surfaces.

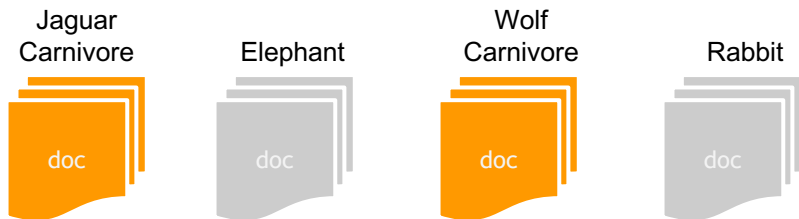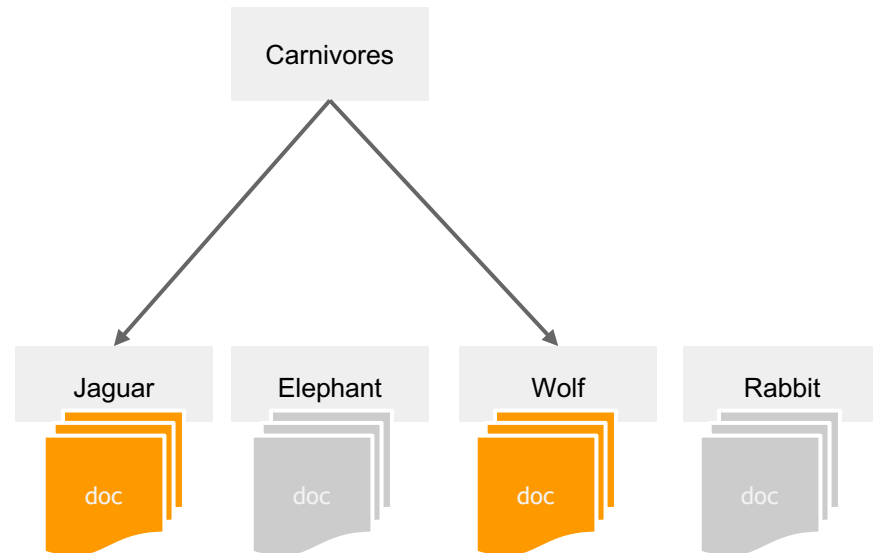# Traditional vs. Graph-based Metadata Management

**Traditional approach**

Show me all documents about *Carnivores*

Graph-based approach

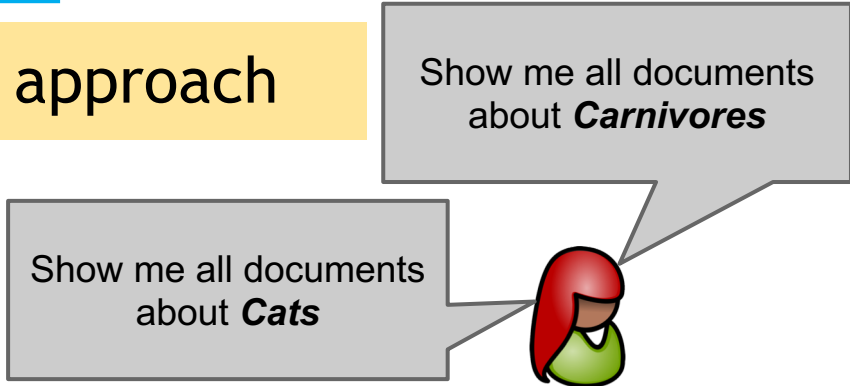| Jaguar | Elephant | Wolf | Rabbit |
|--------|----------|------|--------|
| doc | doc | doc | doc |

| Jaguar | Elephant | Wolf | Rabbit |
|--------|----------|------|--------|
| doc | doc | doc | doc |

# Traditional vs. Graph-based Metadata Management

# Traditional vs. Graph-based Metadata Management

# Traditional vs. Graph-based Metadata Management

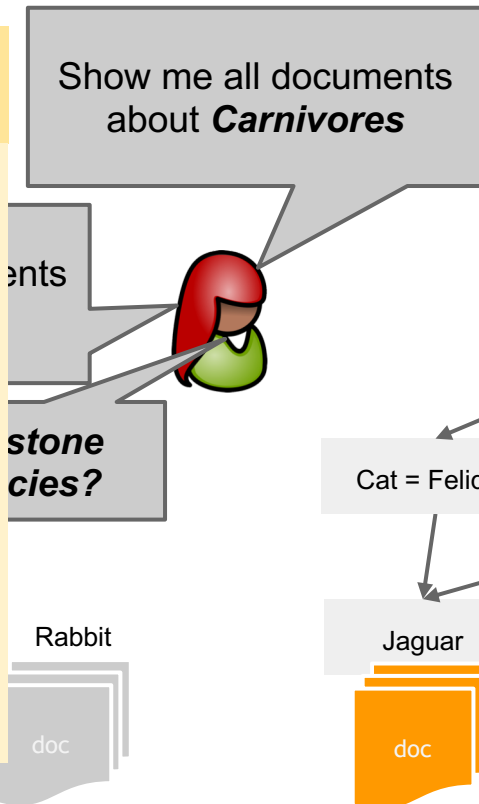# Traditional vs. Graph-based Metadata Management

# Traditional vs. Graph-based Metadata Management

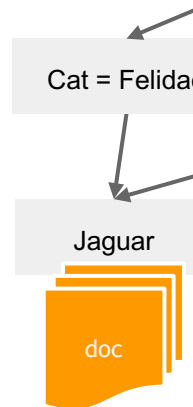

## Traditional approach

### Metadata per document

1. No or little network effects
2. No reuse of metadata
3. Metadata resides in silos
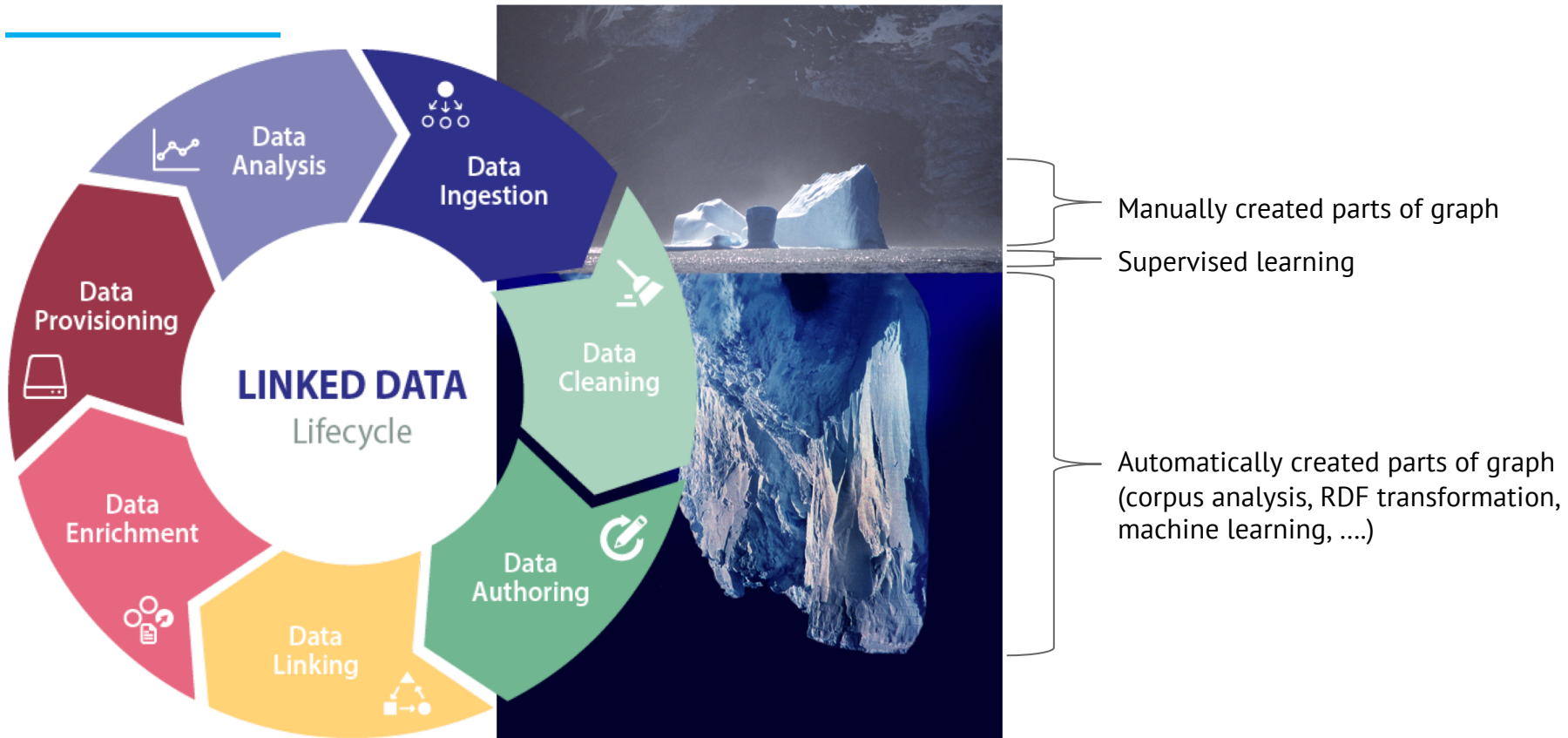4. Data quality hard to measure
5. Not machine-readable

## Graph-based approach

### Knowledge about metadata

1. Explicit knowledge models
2. Reusable and measurable
3. Metadata is machine-processable
4. Standards-based metadata
5. Linkable metadata opens silos

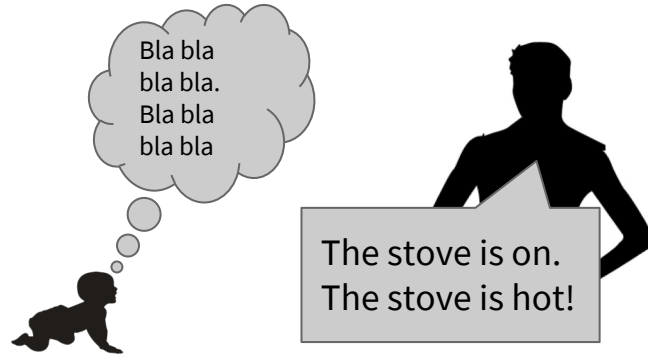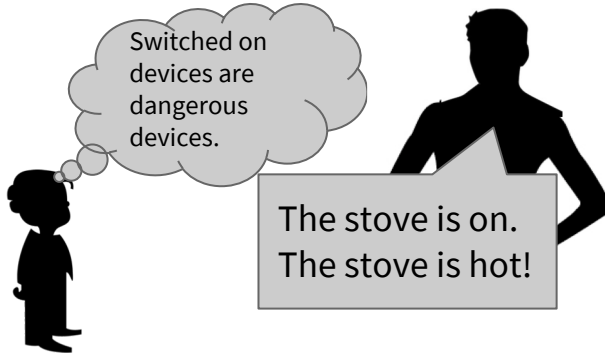# Knowledge Graphs as a result of human-machine cooperation



Manually created parts of graph

Supervised learning

Automatically created parts of graph (corpus analysis, RDF transformation, machine learning, ....)

# How does it work?

# CORPUS ANALYSIS

Semantic Information Retrieval

# Bionics

## How do we learn from a lot of text?

# Extraction of the non-obvious

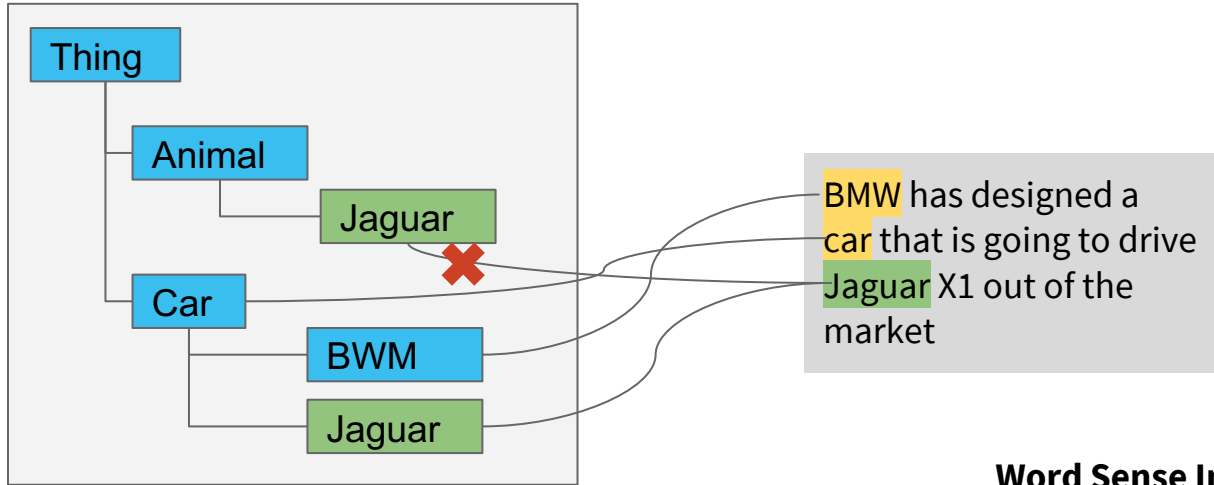| Mechanism | How does it work? |
|---|---|
| Entity extraction based on concepts instead of simple term-based extraction | Make use of synonyms in a taxonomy; disambiguation when necessary |
| Extraction of super classes and broader concepts | Make use of hierarchical structures in a knowledge graph |
| Extraction of related terms and concepts | Co-occurrence models and word embeddings |
| Extraction of 'Shadow Concepts' | Combining co-occurrence models and knowledge graphs |
| Semantic Classifier | Enrichment of training documents with metadata from a knowledge graph |
| Deep Text Analytics | Derivation of new metadata and classifications based on mechanisms as described above combined by rules |

Complexity

# Entity extraction based on concepts instead of simple term-based extraction



Vienna
http://localhost/DemoGeo/3

+ Add to Collection  ⊘ Add to Blacklist  ⊘ Add to ExactMatch  🗑 Delete Concept

**Details** | Notes | Documents | Linked Data | Triples | Visualization | Quality Management

History

**SKOS**  👤 +

**Broader Concepts**
Austria

**Narrower Concepts**
Carinthia
Tyrol
Upper Austria

**Related Concepts**

**Top Concept of Concept Schemes**

**Preferred Label**
⊘ Vienna                    `en`
⊘ Wien                      `de`

**Alternative Labels**
⊘ Austrias Capital          `en`
⊘ Capital City of Austria
⊘ Capital of Austria
⊘ Vedunia

**Hidden Labels**
⊘ Wean                      `en`

Schaut man vom Kahlenberg auf die Donau hinunter, kann man Wien mit allen Sinnen spüren. Weinberge sind da zu sehen, dahinter glänzt das bauliche Erbe der mitteleuropäischen Metropole. Ein halbes Jahrtausend wurde hier Weltgeschichte geschrieben. Kunstgeschichte sowieso.
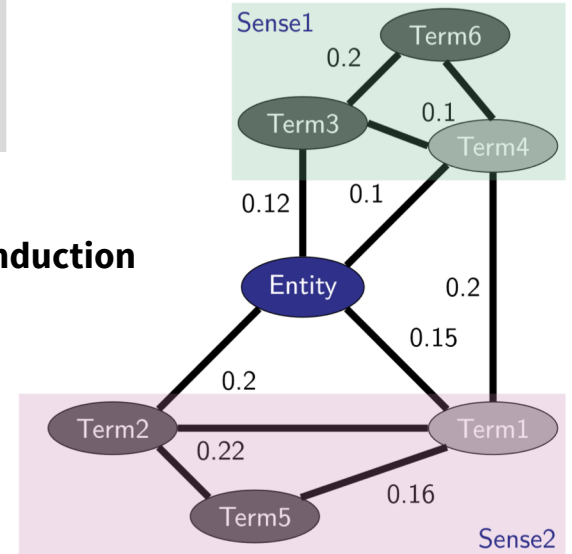
Austria's capital, lies in the country's east on the Danube River. Its artistic and intellectual legacy was shaped by residents including Mozart, Beethoven and Sigmund Freud. The city is also known for its Imperial palaces, including Schönbrunn, the Habsburgs' summer residence. In the MuseumsQuartier district, historic and contemporary buildings display works by Egon Schiele, Gustav Klimt and other artists.

# Extraction of super classes and broader concepts

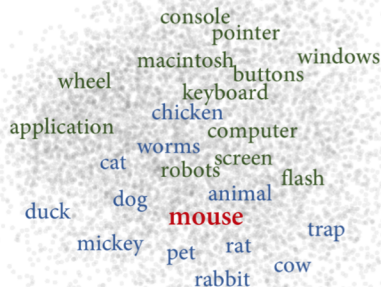Make use of hierarchical structures in a knowledge graph

Thing

Animal

Jaguar

❌

Car

BWM

Jaguar

BMW has designed a car that is going to drive Jaguar X1 out of the market

**Word Sense Induction**

Sense1

Term6

0.2

Term3

0.1

Term4

0.12

0.1

Entity

0.2

0.15

0.2

Term2

0.22

Term1

Term5

0.16

Sense2

# From Word to Sense Embeddings

## Unsupervised sense representation (Word sense induction)

... number of **cells** in plants and animals varies ... officers wait with prisoners in **cell** ... equilibrium is reached, the **cell** cannot provide further voltage ... outer membrane of the **cell** ... new lithium ion **cell** in the Model S Tesla ... carried out a pioneering human embryonic stem **cell** operation ... **cell** towers are usually interconnected ...

(1) Get occurrences of a word from text corpora

... number of **cells** in plants and animals varies ... officers wait with prisoners in **cell** ... equilibrium is reached, the **cell** cannot provide further voltage ... outer membrane of the **cell** ... new lithium ion **cell** in the Model S Tesla ... carried out a pioneering human embryonic stem **cell** operation ... **cell** towers are usually interconnected ...
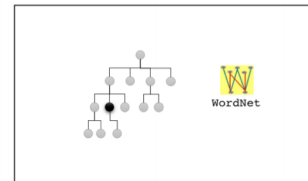
(2) Analyze contexts and induce senses of the word

cell#1
cell#2
cell#3
cell#4

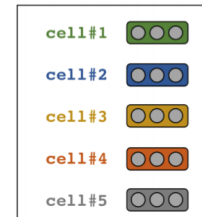(3) Compute sense representation

## Knowledge-based sense representation

1. **cell#1** (jail_cell, prison_cell): a room where a prisoner is kept.
2. **cell#2** the basic structural and functional unit of all organisms.
3. **cell#3** (cellphone, mobile_phone): a hand-held mobile radiotelephone.
4. **cell#4** (electric_cell): a device that delivers an electric current.
5. **cell#5** (cubicle): small room in which a monk or nun lives.

(1) Get senses as defined by a sense inventory (e.g., WordNet)

WordNet

(2) Gather information for each sense (e.g., by exploiting the structural properties of sense inventory's semantic network, and (optionally) then from text corpora)
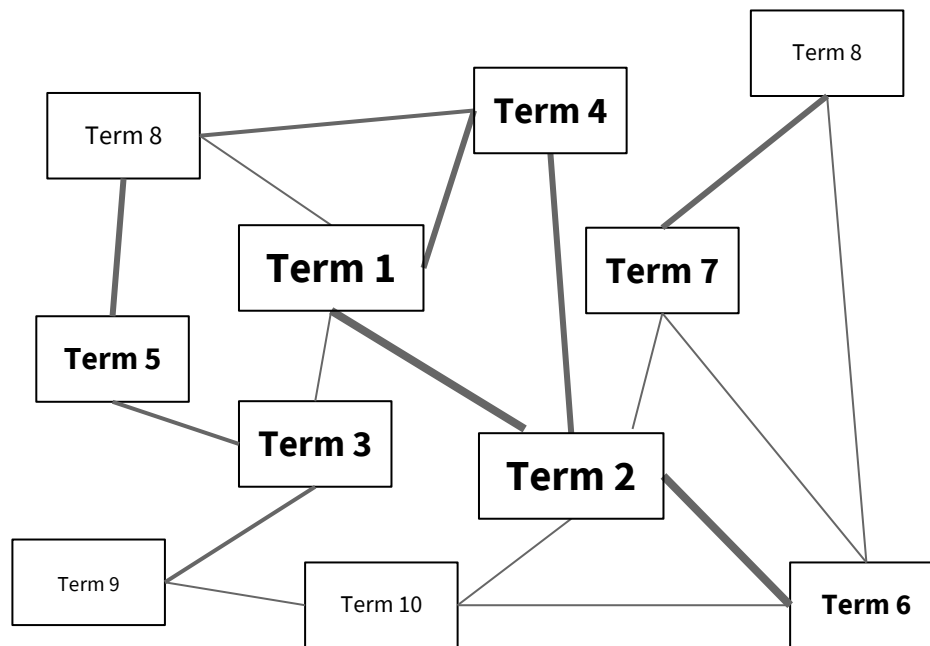
cell#1
cell#2
cell#3
cell#4
cell#5

(3) Compute sense representation

console pointer
macintosh windows
buttons
wheel keyboard
chicken
application computer
worms screen
cat robots flash
dog animal
duck mouse
mickey rat trap
pet cow
rabbit

### The main problem:
Meaning conflation deficiency
(Word sense ambiguity)

# Extraction of related terms and concepts



Document
Corpus
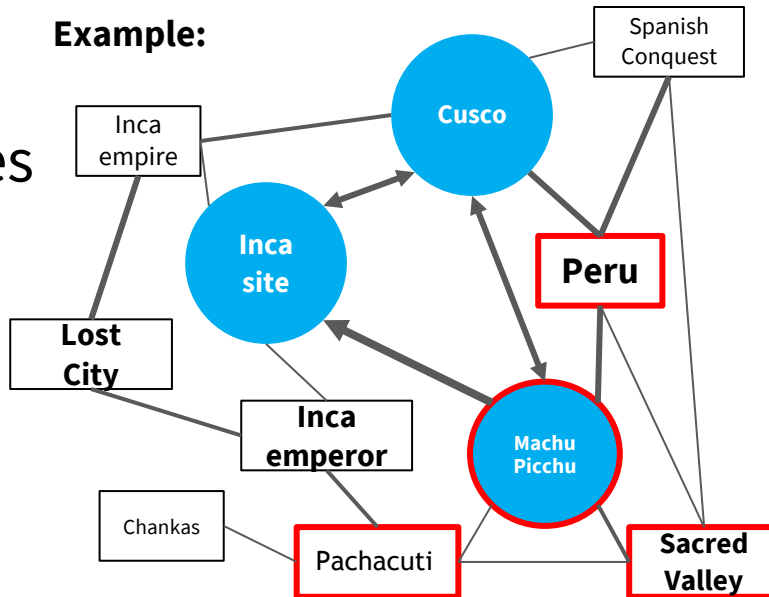
▸ Websites

▸ PDF, Word, …

▸ Abstracts from
  DBpedia

▸ RSS Feeds

▸ Relevant terms and phrases

▸ Relevancy of terms

▸ co-occurrence between terms and terms

# Extraction of 'Shadow Concepts'

**Example:**

Use co-occurrences between concepts and terms to extract 'shadow concepts'



This site is a 15th-century Inca site located 2,430 metres above sea level. It is located in Cusco, **Peru**.

It is situated on a mountain ridge above the **Sacred Valley** through which the Urubamba River flows. Most archaeologists believe that it was built as an estate for the Inca emperor **Pachacuti**. Often mistakenly referred to as the "Lost City of the Incas", it is the most familiar icon of Inca civilization. The Incas built the estate around 1450, but abandoned it a century later at the time of the Spanish Conquest.

In addition to explicitly used concepts and terms, ***Machu Picchu*** is extracted from the article as a *Shadow Concept*. As a prerequisite, one has to provide and analyze a representative text corpus first.

# Semantic Classifier



# Semantic Classifier

ML based on semantically enriched training data

**PoolParty Semantic Classifier** combines machine learning algorithms (SVM, Deep Learning, Naive Bayes, etc.) with Semantic Knowledge Graphs.

# Deep Text Analysis
Annotation, Extraction, Classification, Linking

- ▸ Corpus statistics / Word embeddings
  → **Keyphrase extraction**

- ▸ Graph-based annotation
  → **Entity/Concept linking**

- ▸ Corpus Statistics embedded in graphs
  → **Shadow Concepts**

- ▸ Machine-learning-based annotation
  → **Named entity recognition (NER)**

- ▸ Machine-learning based classification
  → **Document Classification**

- ▸ Annotation based on rules
  → **Regular expressions**

ML-based
Entity Extraction

Private Equity
▸ Funds (1)
▸ Locations (3)
  ▸ Americas (3)
  ▸ Asia (1)
  ▸ Europe (7)
▸ Markets (2)
  ▸ Industry (22)
  ▸ Region (1)
▸ Organizations (2)
  ▸ Company (47)
  ▸ Private equity firm (26)
▸ Persons (2)
  ▸ Chairman (6)
  ▸ Founder (13)
▸ Strategies & Products (18)

Graph-based
Entity
Linking

Bain Capital is a venture capital company based in Boston, MA.
Since inception it has invested in hundreds of companies. In 2018, Bain had $75b AUM.

Regular
Expressions-based
Annotation

Semantic
Rules Engine

Give me all paragraphs in documents about **"US based Private Equity firms with AUM higher than $20B"**

poolparty®

# USE CASES

See how it can be used in your environment!

# Extract concepts from text even if not used explicitly

Some domains use text that doesn't always call a spade a spade. With 'shadow concept extraction' those 'masked' concepts still can be surfaced.

Since these technologies would have become conventional technologies that are made into products and introduced into market at the time of their introduction, it would be difficult to differentiate them as innovative environmental and energy technologies from other global warming prevention technologies that have already been put to practical use in the industrial, commercial, residential, and energy conversion sectors.
- The Innovative Global Warming Prevention Technology Working Group under the Research and Development Subcommittee
- Council assessed that innovative global warming prevention technologies would bring about a reduction effect of 7.49 million t-CO2 case of average emissions factor for all power sources of carbon dioxide in 2010. In view of the difficulty in putting innovative carbon dioxide sequestration technology into practical use by 2010, the Working Group reassigned it as an issue of global warming prevention technology to be tackled by 2030.
The Central Environment Council, however, has not had the opportunity to examine the contents of these technologies in detail. (Promotion of climate change prevention activities by every social actor)
- The Programme encourages every social actor to take actions to prevent global warming. The actions include measures undertaken by the public sector.

**Climate Change**

Since these technologies would have become conventional technologies that are made into products and introduced into market at the time of their introduction, it would be difficult to differentiate them as innovative environmental and energy technologies from other ▮▮▮▮▮▮▮ prevention technologies that have already been put to practical use in the industrial, commercial, residential, and energy conversion sectors.
- The Innovative ▮▮▮▮▮▮▮ Prevention Technology Working Group under the Research and Development Subcommittee
- Council assessed that innovative ▮▮▮▮▮▮▮ prevention technologies would bring about a reduction effect of 7.49 million t-CO2 case of average emissions factor for all power sources of carbon dioxide in 2010. In view of the difficulty in putting innovative carbon dioxide sequestration technology into practical use by 2010, the Working Group reassigned it as an issue of ▮▮▮▮▮▮▮ prevention technology to be tackled by 2030.
The Central Environment Council, however, has not had the opportunity to examine the contents of these technologies in detail. (Promotion of ▮▮▮▮▮▮▮ prevention activities by every social actor)
- The Programme encourages every social actor to take actions to prevent ▮▮▮▮▮▮▮. The actions include measures undertaken by the public sector.

**Climate Change**

**Mini Countryman**

And it's probably more of a crossover than ever, with the design to match, Being a Mini, the Countryman is clearly meant to be the driver's car among small crossovers. The suspension is sophisticated, and there are lots of chassis options (a stiffer sports setup, variable damping, the electronically-controlled ALL4 all-whee...

But it's also the crossover for people who've bag... luxury.

There's been a lot of effort on ramping up the ca... was a sad let-down in that department.

On the outside, plastic wheel-arch extensions, w... as well as roof bars and sill protectors all add to ... only Mini with angular rather than oval headlamp... on in the lower face.

There are eight versions at launch, and they're e... S, each fuelled by petrol or diesel, each of them ... auto, too, if you count that as a separate choice. The Cooper petrol is a three-cylinder, the rest fours.

You get extra kit as standard versus the old car, ... and park sensors. Upgrades include a bigger to... various posher seats, a HUD, and driver aids. O... boot so you can sit on the rear bumper without g...

In June 2017 a Cooper E will launch, which has ... front wheels, and an electric motor for the rears, ... gentle all-electric running. So it has the performa... advantages of a plug-in hybrid. And you wouldn... distance.

The platform is BMW's contemporary transverse... sizes. That means it shares a lot with the BMW X1. The 4WD system is more sophisticated than the previous Countryman's. The proportion of drive to the rear is computed by a controller that takes into account parameters including grip, steering angle and throttle position, as well as whether you've got the sports mode and sports traction systems selected.

### Concepts

Front-wheel drive

**Four-wheel drive**

Throttles

Bodies

# Gasoline

Hardware        Head Rests

### Shadow Concepts

Diesel Engines

United States

# Engine

Spoiler

V6 engine        Sedan

Mid-size car        Interiors

Four-cylinder engine

# Use Case **Recommender System**

Connecting

▶ content to content

▶ people to content

▶ people to people

Demo

I am interested in ...

Architecture

Mediterranean Sea

Wine Tasting

Occitanie, vineyard, sea, church, village

**Languedoc-Roussillon**

Languedoc is a significant producer of wine, and a major contributor to the surplus known as the "wine lake". Today it produces more than a third of the grapes in France, and is a focus for outside investors.

The region contains the historic cities of Carcassonne, Toulouse, Montpellier, countless Roman monuments, medieval abbeys, Romanesque churches, and old castles.

# Example **EIP on Water**

At the center of the marketplace is the matchmaking function



https://www.eip-water.eu/

# Contract Intelligence

Semantic analysis and pre-selection of relevant clauses in contracts

# Connect

## Andreas Koller

CIO & Managing Partner, Semantic Web Company

- ▸ andreas.koller@semantic-web.com
- ▸ https://www.linkedin.com/in/ankoller/
- ▸ https://twitter.com/semwebcompany

Q & A