

Module IV: Knowledge Graph Fundamentals and Construction

1:00 pm - 2:20 pm

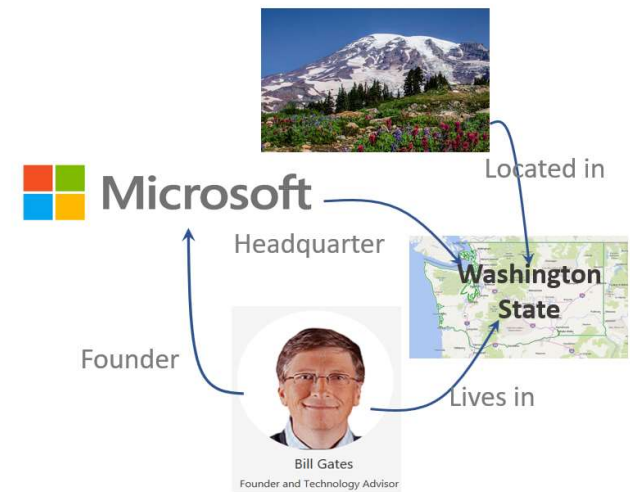
A brief history of Knowledge Graph (KG)

KG representation and examples

KG Construction

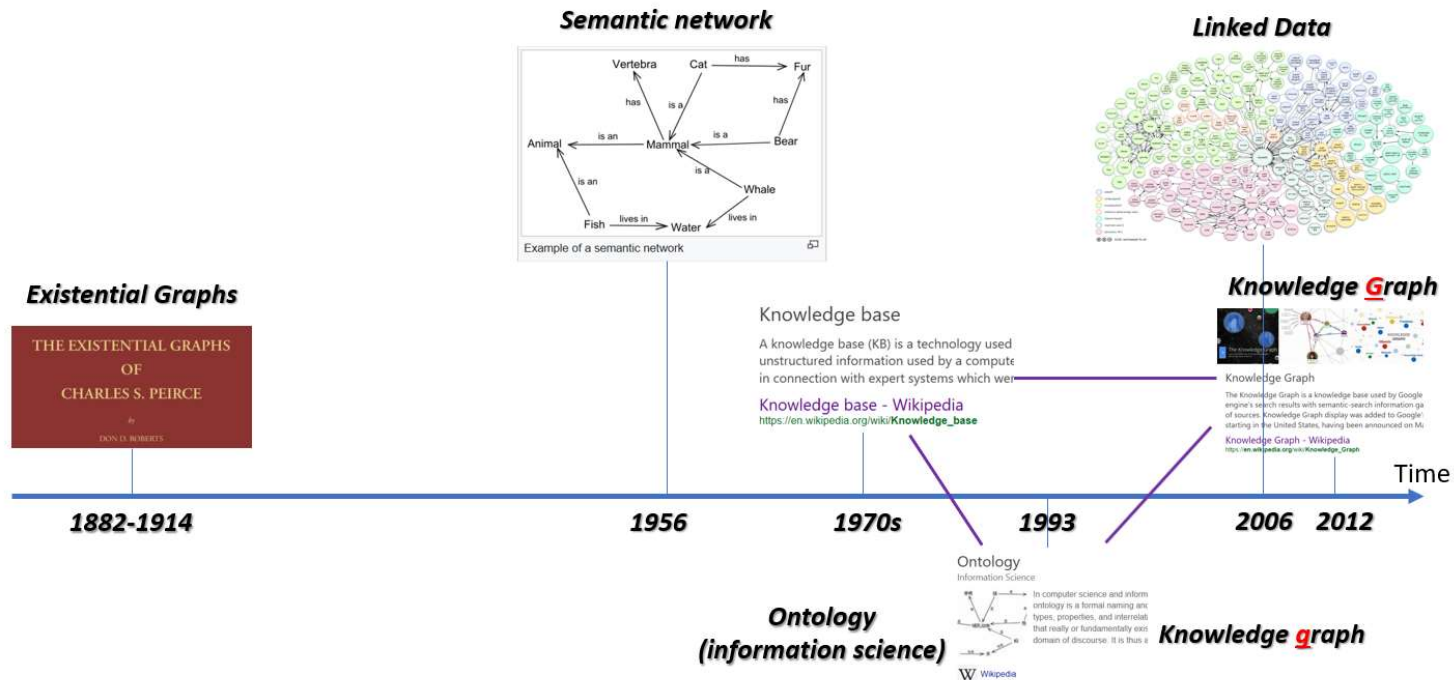
- General overview
- NLP techniques
- MAG case study

Lab 4 – Enrich concepts and taxonomy



Module 4 Overview

KG fundamentals and construction



A brief history of Knowledge graph

► Graph

► Node

- Attribute1
- Attribute2
- Attribute3
- ...

► Edge

(Node1, Node2, weights)

Homogeneous	vs	Heterogeneous
-------------	----	---------------

Node: 1 table Edge: 1 table		
--------------------------------	--	--

		Node: multiple tables Edge: multiple tables
--	--	--

• (Subject, Predicate, Object)

- Each **node** has a universal id (S)

- Its **attribute** is represented as:

(S, attributeName (P), attributeValue(O))

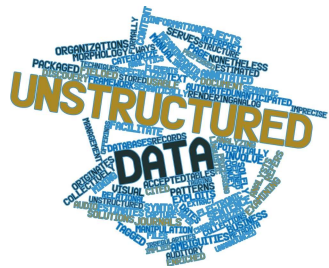
- An **edge** connected two nodes (e.g. S1, S2):

(S1, relationName (P), S2(O))

Homogeneous	vs	Heterogeneous
-------------	----	---------------

Node + Edge:		SINGLE table
--------------	--	--------------

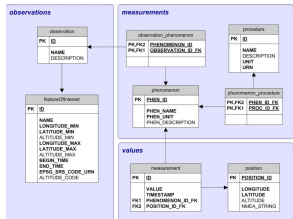
Knowledge graph Representation



Unstructured Documents



Natural language processing (NLP)



Existing Relational Databases

Data pipeline processing

Common sense
is NOT so
COMMON.

Human
Common
Sense

Manual efforts



Knowledge in the **graph** form

Knowledge Graph Construction

► Sentence Level

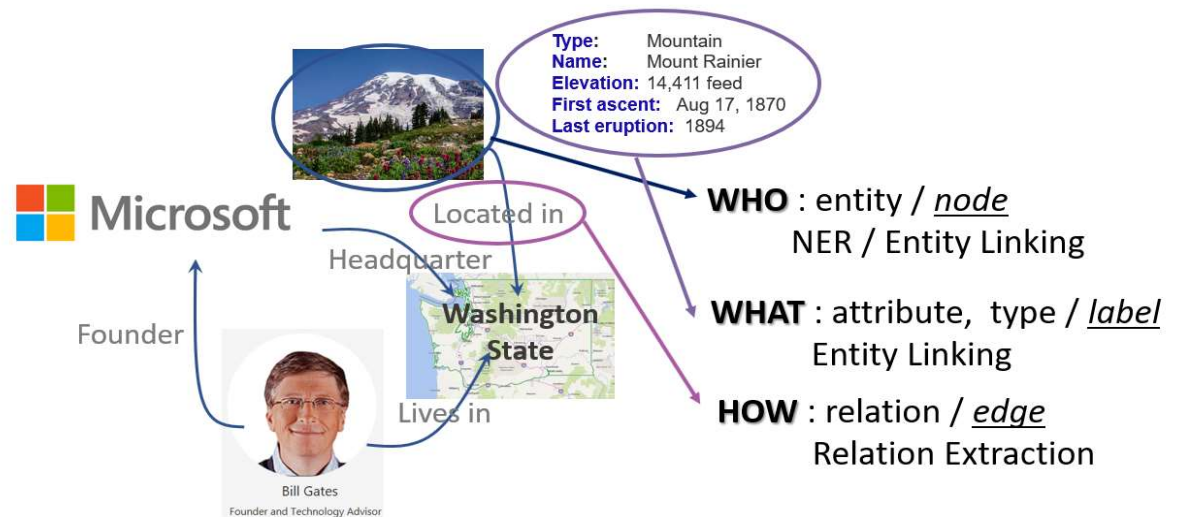
- Part-of-speech (PoS) Tagging
- Named entity recognition (NER)
- Dependency Parsing

► Document Level

- Coreference resolution
- Topic model
- Classification

► Information Extraction

- Entity resolution
- Entity linking
- Relation extraction



NLP Techniques

for Knowledge Graph Construction – At a glance

► Challenges:

- Incomplete
- Inconsistent
- Ambiguous

*Precision
Slow*



Head

Supervised



Torso

Semi-supervised
(Distantly-supervised)



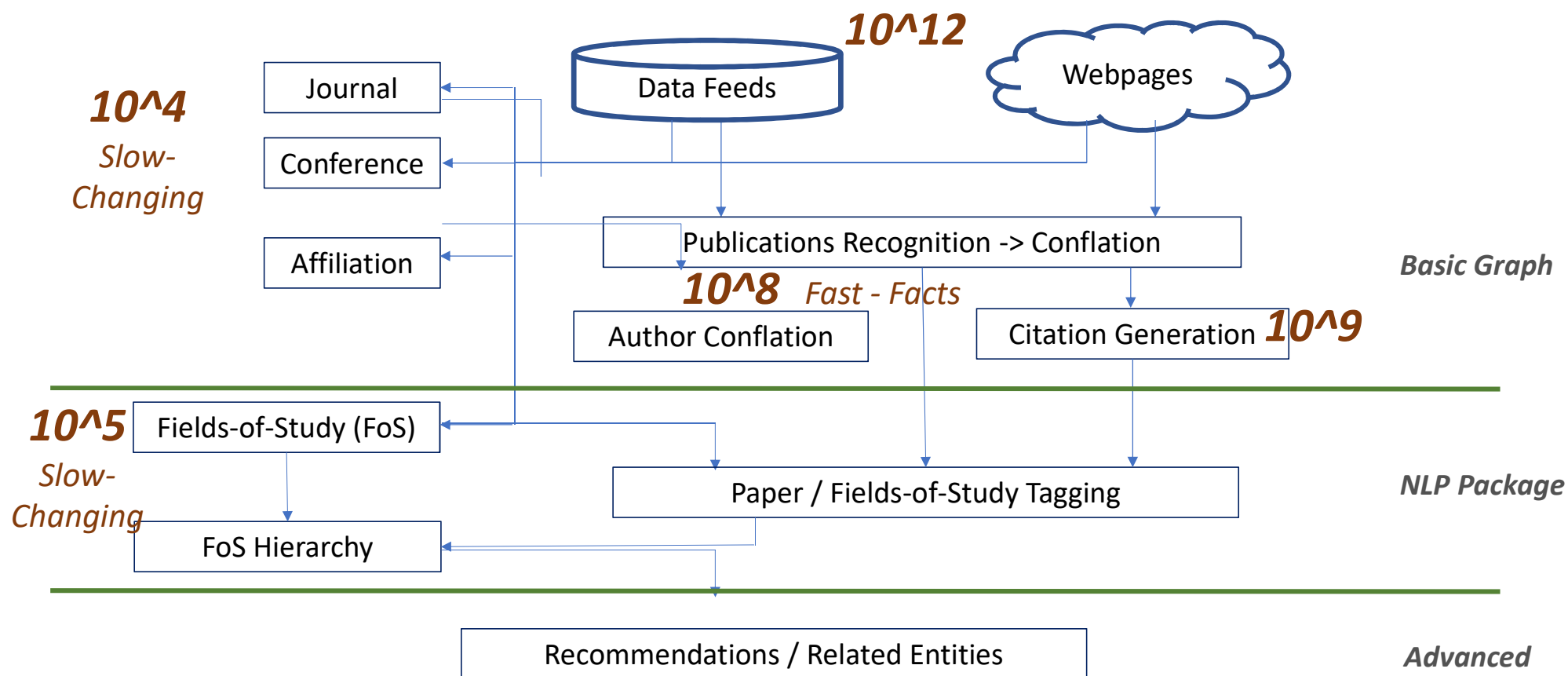
Long Tail

Unsupervised



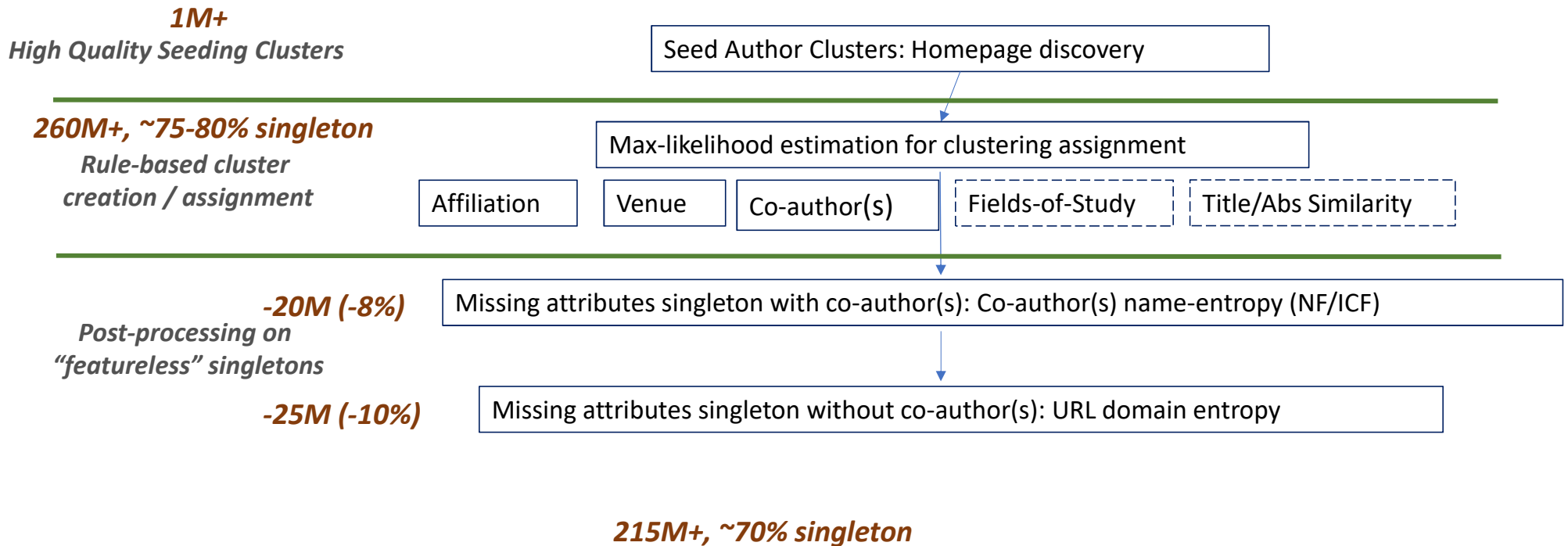
*Recall
Fast*

KNOWLEDGE GRAPH CONSTRUCTION



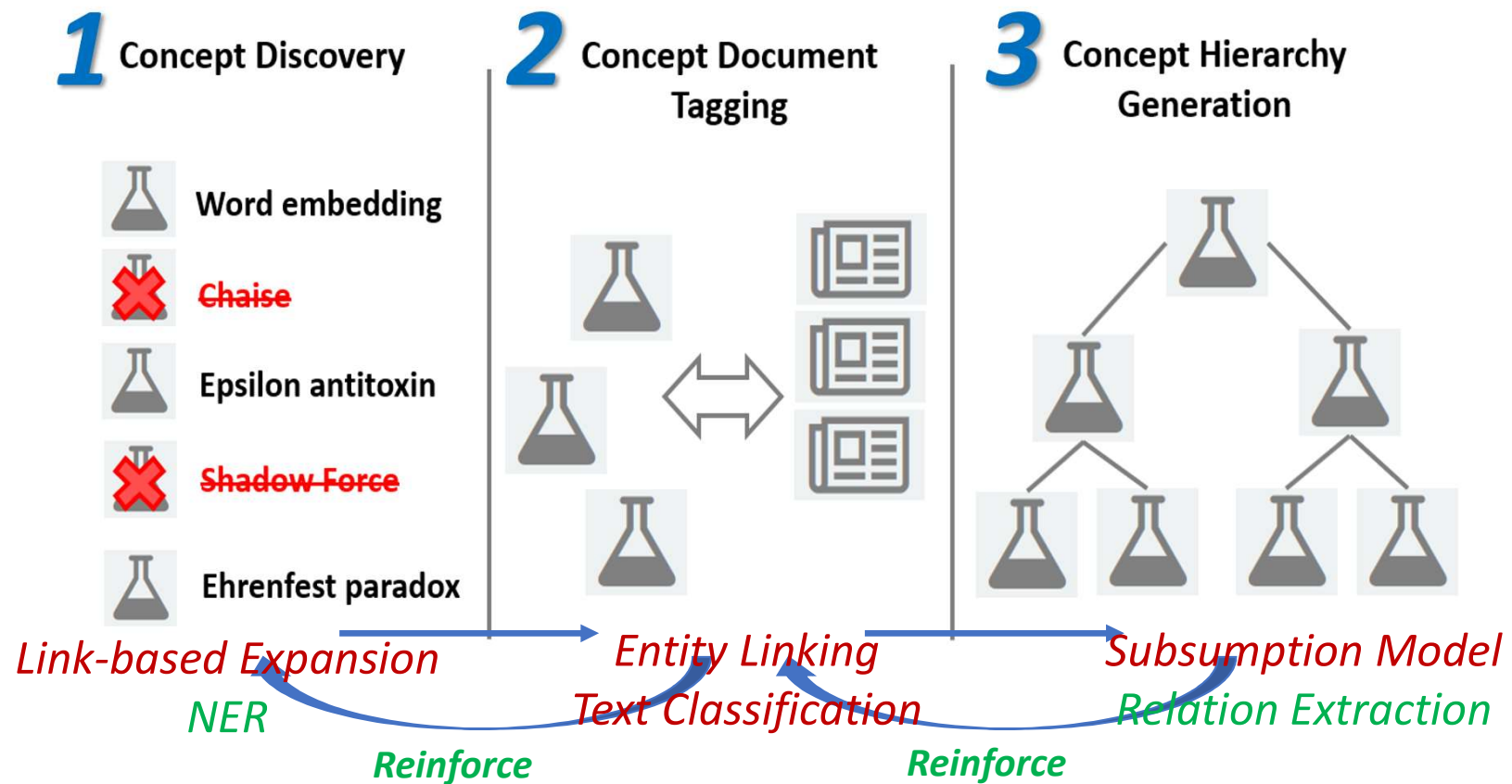
Microsoft Academic Graph (MAG) Construction

Guiding principles: prefer under-conflation than over-conflation (<0.1%)



MAG Construction

Author Recognition/conflation



MAG Construction

Concept Recognition/Tagging/Hierarchy

		Word	Concept	Publication
<i>Language Features</i>	Discrete Space	1-hot vector	Bag-of-Words (Desc.)	Bag-of-Words (Title + Abstract)
	Continuous Space	Word Vector *	Concept Embeddings **	Bag-of-Concepts (Title + Abstract)
<i>Structural Features</i>	Discrete Space		Related Concepts	Citation / Venue [/Author]
	Continuous Space		Heterogenous Graph Embedding	Homogenous Graph Embedding

Mixture Model to represent Concepts & Publications (or any documents)

*Word vector are train on Skip-gram model with 13B tokens of 130M English publications, vocab size: ~2M

Concept embeddings - average of the word vector in concept description. (Description** is important!)

MAG Construction

Concept / Publication representation

L5	L4	L3	L2	L1	L0
Convolutional Deep Belief Networks	Deep belief network	Deep learning	Artificial neural network	Machine learning	Computer Science
(Methionine synthase) reductase	Methionine synthase	Methionine	Amino acid	Biochemistry / Molecular biology	Chemistry / Biology
(glycogen-synthase-D) phosphatase	Phosphorylase kinase	Glycogen synthase	Glycogen	Biochemistry	Chemistry
	Fréchet distribution	Generalized extreme value distribution	Extreme value theory	Statistics	Mathematics
Hermite's problem	Hermite spline	Spline interpolation	Interpolation	Mathematical analysis	Mathematics

Completely Data-Driven (L2-L5)

6-Level Hierarchy with 660K+ Concepts

MAG Construction

Step	Accuracy
1 Concept Discovery	94.75%
2 Concept Tagging	81.20%
3 Hierarchy Building	78.00%

- Randomly sampled
- 500 data points divide to five groups per step
- 1 judge per group on step 1,2
- 3 judges per group on step 3

Concept Hierarchy Results

- Demo: new concept discovery
- Task : Generate concept hierarchy using subsumption

Lab 4: Enrich concepts and taxonomy