# ADVANCED TOPICS

Krzysztof Kutt, PhD
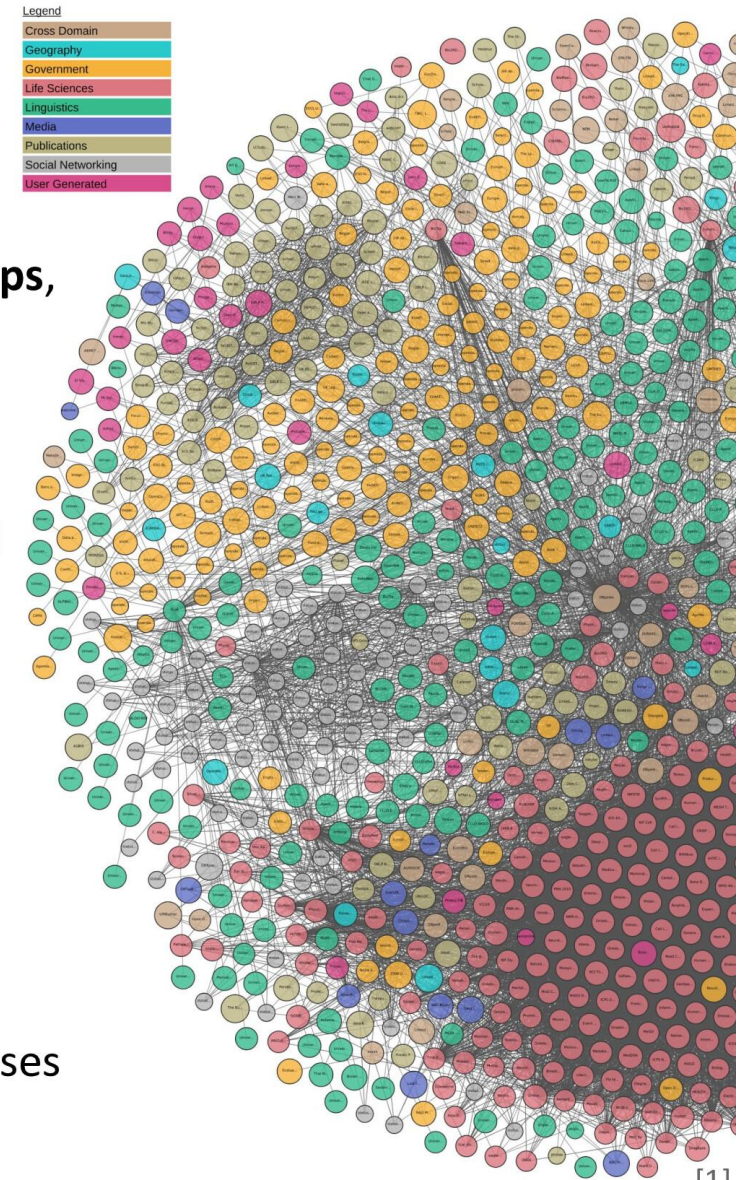Semantic Internet, WFAIS UJ

# OUTLINE

1. Knowledge graphs

2. Towards automated KG management

3. Semantic search and recommendations

4. Knowledge graph embeddings

5. Knowledge graph completion

6. ChatGPT is a bullshit. How can we fix it?

# THE GRAPHS

How to compare them?

# Knowledge Graph Recap

Legend
Cross Domain
Geography
Government
Life Sciences
Linguistics
Media
Publications
Social Networking
User Generated

- A **Graph** consisting of **concepts**, **classes**, **properties**, **relationships**, and **entity descriptions**
- Based on **formal knowledge representations** (RDF(S), OWL)
- Data can be **open** (e.g. DBpedia, WikiData), **private** (e.g. supply chain data), or **closed** (e.g. product models)
- Data can be **original**, **derived**, or **aggregated**
- We distinguish
  - **instance data** (ground truth),
  - **schema data** (vocabularies, ontologies)
  - **metadata** (e.g. provenance, versioning, licensing)
- **Taxonomies** are used to categorize entities
- **Links** exist between internal and external data
- Including **mappings** to data stored in other systems and databases
- *Fully compliant to **FAIR Data principles***

[1]

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Knowledge Base Definition

A Knowledge Graph is a **Knowledge Base** that is a Graph.

A **knowledge base** (KB) is a technology used to store complex structured and unstructured information used by a computer system. The initial use of the term was in connection with expert systems which were the first knowledge-based systems.

*Wikipedia*

### knowledge base

*Free Online Dictionary of Computing*

‹*artificial intelligence*›

A collection of knowledge expressed using some formal knowledge representation language. A knowledge base forms part of a knowledge-based system (KBS).
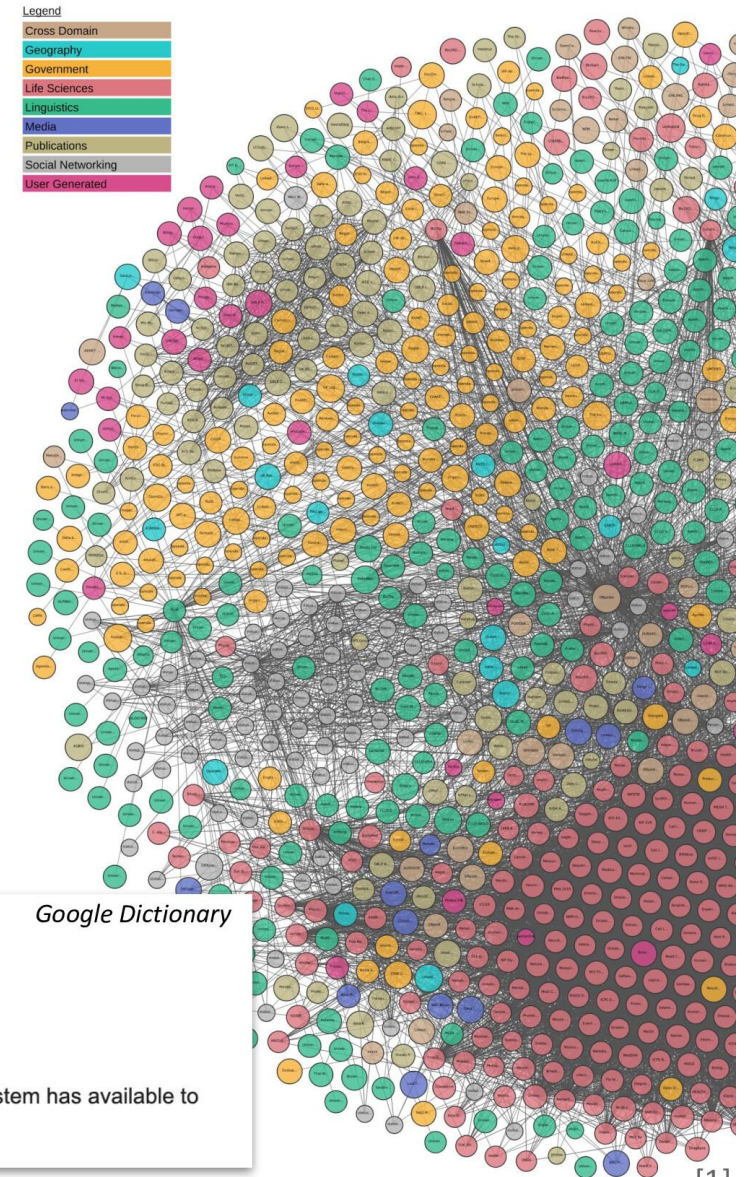
🔊 knowledge base

*Google Dictionary*

*noun*

1. a store of information or data that is available to draw on.
2. the underlying set of facts, assumptions, and rules which a computer system has available to solve a problem.



**Legend**
Cross Domain
Geography
Government
Life Sciences
Linguistics
Media
Publications
Social Networking
User Generated

The Linked Open Data Cloud from lod-cloud.net

[1]

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Graph Definition

A Knowledge Graph is a Knowledge Base that is a **Graph**.

---

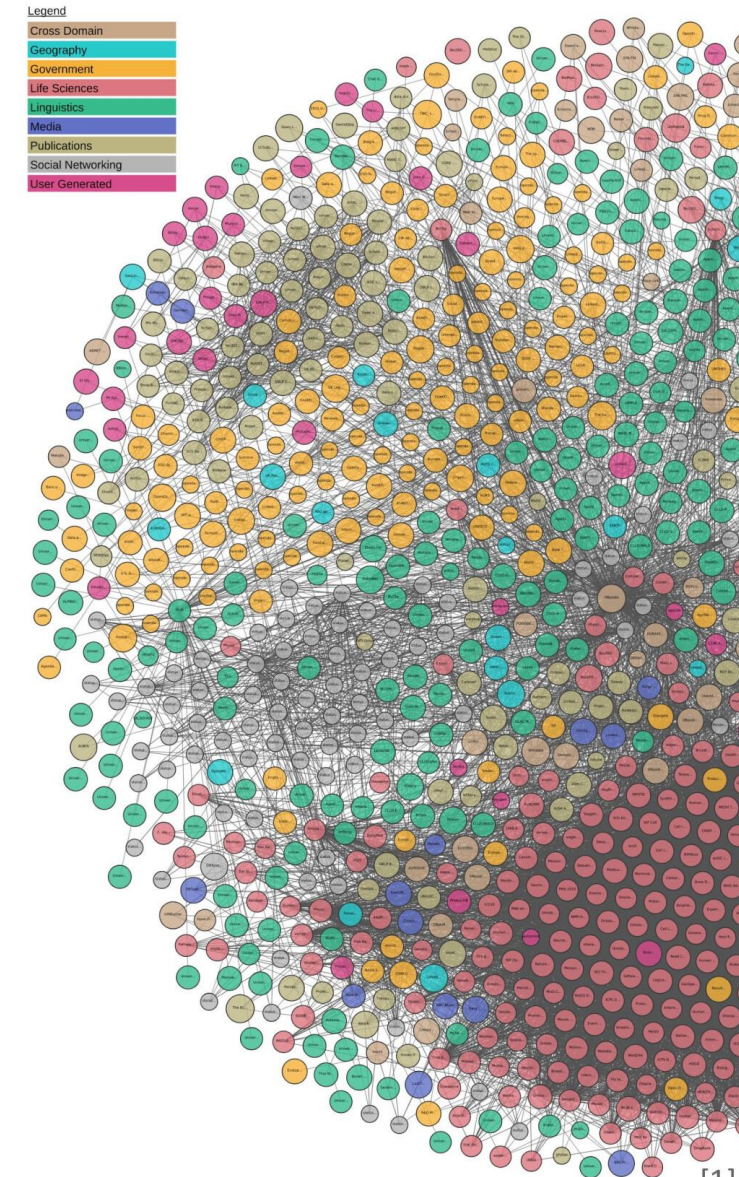**Definition**                                                                 **1.1**
A **simple directed graph** **G=(V,E)** consists of a set **V** of **vertices**, |V|=n, and a set **E** of **directed edges**, $E \subseteq V \times V$, where each edge $e_i=(v_k, v_l)$, $e_i \in E$
is an ordered pair of two vertices $(v_k, v_l)$ with $v_k, v_l \in V$.

---

**Definition 1.2**

- A **graph with self-loops** is a graph extended with the option of having edges that relate a vertex to itself.
- A **multi-graph** is a graph that may have multiple edges with the same vertices.
- An **edge-labelled graph** is a graph that has an additional **labelling function** $\lambda : E \rightarrow L$ that maps each edge in E to an element in a set of labels L (similarly for vertex-labelled graphs).
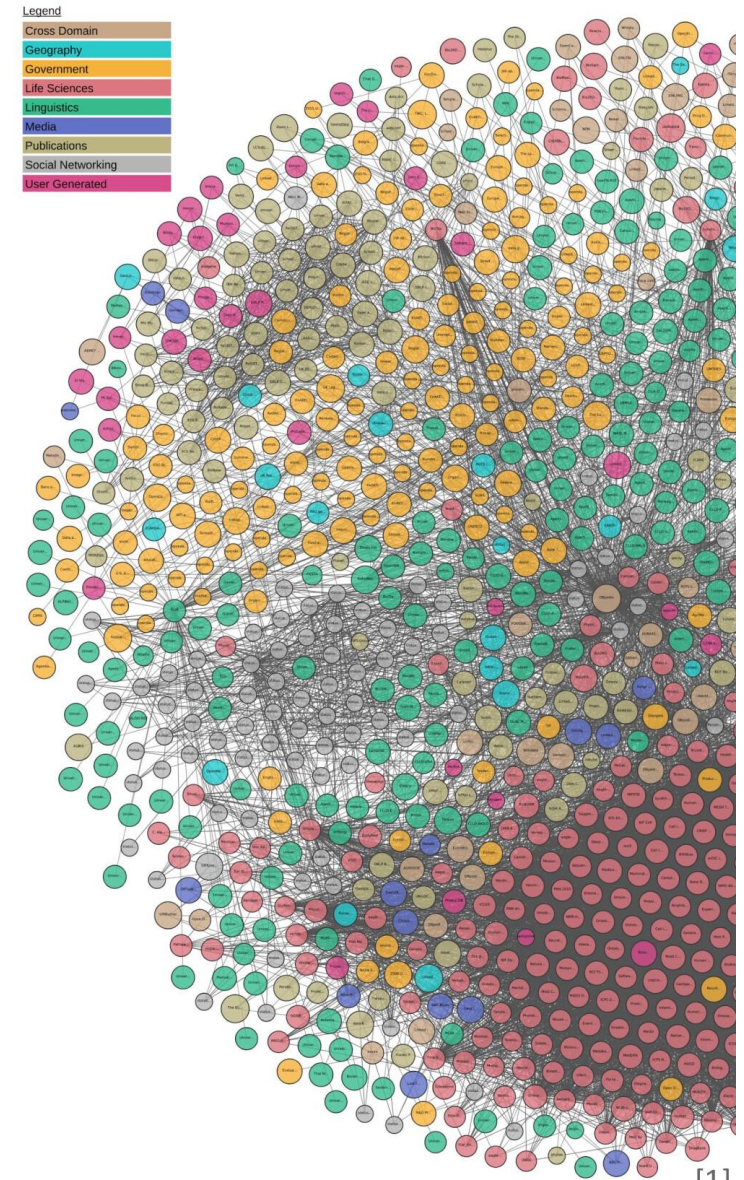
---

Legend
Cross Domain
Geography
Government
Life Sciences
Linguistics
Media
Publications
Social Networking
User Generated

The Linked Open Data Cloud from lod-cloud.net

[1]

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Graph Definition (cont.)

**Definition 1.3**

- An edge is said to be **incidental** to the vertices it connects.
- The **degree** of a vertex is the number of edges that are incidental to it.
- In a directed graph, the **in-degree** of a vertex is the number of edges pointing towards it; analogously for **out-degree**.

**Definition 1.4**

- A **directed path** in a directed graph is a sequence of consecutive edges $(e_1, e_2, \ldots, e_n)$ with $e_i=(v_l,v_k)$ and $e_{i+1}=(v_k,v_m)$.

- A directed graph is **strongly connected** if there is a directed path from any vertex to any other vertex.

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# How Can You Characterize a Knowledge Graph?

"Should I use Knowledge Graph A or Knowledge Graph B to solve my problem?"

- How to compare two Knowledge Graphs?
  - Size
  - Coverage
  - Completeness
  - Level of Detail
  - Accuracy
  - Reliability
  - etc.

- **Idea**: **Structural Comparison** by just comparing the Graphs

[1]

Slide from Knowledge Graphs course by prof. Harald Sack & Mewish Alam (openHPI, 2020).

# Graph Centrality Measures

- **Network analysis** has developed methods for **finding the most important vertices in a graph**.

- **Vertex importance** based on the structure of such graphs is called **centrality.**

- But, what makes a node important?

Knowledge Graphs 2020 , Prof. Dr. Harald Sack & Dr. Mehwish Alam, FIZ Karlsruhe - Leibniz Institute for Information Infrastructure & Karlsruhe Institute of Technology

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# What makes a Node important?

- Many networks can be considered to describe a **flow** of something (goods, information, etc.)
- A node might be **important**, if
  - a lot flows from it (in a supply chain),
  - to it (in a network of links), or
  - through it (in a communication network)
- Flow might be modelled by (weighted) paths, possibly factoring in their length and/or number
- Paths might be more important if they pass through important nodes
- In knowledge graphs, the importance of edges and nodes may also depend on more complex features (e.g., edge or vertex labels)

[1]

# What makes a Node important?

- **Wikidata Example:**
  - A Wikidata entity (node) might be important, if it is referenced by many Wikipedia pages
  - what are the most important Climatologists?

```
SELECT ?climatologistLabel (SUM(?link) AS ?importance)
WHERE {
    ?climatologist wdt:P106 wd:Q1113838 .
    ?climatologist wikibase:sitelinks ?link.
    ?climatologist rdfs:label ?climatologistLabel
      FILTER (lang(?climatologistLabel)="en")
} GROUP BY ?climatologistLabel
ORDER BY DESC(?importance)
```

Knowledge Graphs 2020 , Prof. Dr. Harald Sack & Dr. Mehwish Alam, FIZ Karlsruhe - Leibniz Institute for Information Infrastructure & Karlsruhe Institute of Technology

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

SPARQL query

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Degree Centrality

- A simple form of centrality restricts to incoming/outgoing paths of length one

**Definition 1.5**

- The **in-degree centrality** of a directed graph is given by the in-degree of each node.
- The **out-degree centrality** and the **degree centrality** (for undirected graphs) are defined analogously

- There are more sophisticated forms of centrality, as e.g.
  - Eigenvector centrality, Katz centrality, PageRank, etc.

Knowledge Graphs 2020 , Prof. Dr. Harald Sack & Dr. Mehwish Alam, FIZ Karlsruhe - Leibniz Institute for Information Infrastructure & Karlsruhe Institute of Technology

# Further Centrality Measures

- Further Measures to characterize a Knowledge Graph
  - Sizes
    - number of nodes
    - number of facts
    - avg number of facts per node
  - KG diameter

---

**Definition 1.6**

- The **eccentricity** of a node is the maximal distance between a certain node and any other node.
- The **diameter** of a graph is the maximum **eccentricity** of a graph, i.e. the greatest distance between any pair of nodes.
- To find the diameter of a graph, first find the **shortest path** between each pair of nodes. The greatest length of any of these paths is the **diameter of the graph**.

---

[1]

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Further Centrality Measures

- Further Measures to characterize a Knowledge Graph
  - Sizes
    - number of nodes
    - number of facts
    - avg number of facts per node
  - KG diameter
  - KG radius

**Definition 1.8**

- The **radius** of a graph is the minimum eccentricity of a graph, i.e. the shortest of the maximum distances between any pair of nodes.

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Further Centrality Measures

- Further (structural) measures to characterize a Knowledge Graph:
  - Sizes
    - number of nodes
    - number of facts
    - avg number of facts per node
  - KG diameter
  - KG radius
  - avg in/out degree
  - avg path length
  - and many more…

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Knowledge Graphs and Important Nodes

- In **Knowledge Graphs**, the importance of nodes might further be depending on
  - the properties (i.e. edge attributes)
  - the node labels (i.e. further attributes of nodes)
  - specific nodes or edges might be ignored, as e.g.
    - Basically for every entity in a (OWL encoded) knowledge graph the following fact holds:
      
      `:entity rdf:type owl:Thing`
    - Therefore, we might ignore this fact if we want to determine the importance of nodes

Knowledge Graphs 2020 , Prof. Dr. Harald Sack & Dr. Mehwish Alam, FIZ Karlsruhe - Leibniz Institute for Information Infrastructure & Karlsruhe Institute of Technology

[1]

# TOWARDS AUTOMATED KG MANAGEMENT

Mappings and alignment

# Knowledge Graph Challenges

- Building a small KG is easy but building a vast system like Google Knowledge Graph is a huge challenge

**Coverage**
Is the information complete?

**Freshness**
Is the information up to date?

**Correctness**
Is the information accurate?

Increase **Freshness** & **Coverage**
Harder to ensure **Correctness**

Increase **Correctness**
Harder to ensure **Freshness** & **Coverage**

**Correctness** is always hard – what is true and correct?

Knowledge Graphs 2020 , Prof. Dr. Harald Sack & Dr. Mehwish Alam, FIZ Karlsruhe - Leibniz Institute for Information Infrastructure & Karlsruhe Institute of Technology

# Towards Automated Knowledge Graph Management



- Unsupervised knowledge extraction from unstructured data in open domain

- Semantic embedding via Ontologies

- Ultra-scale knowledge representations

- Large scale entity linking and disambiguation

- Autonomous knowledge inference & verification

- Knowledge Graph versioning and archiving

- Knowledge Precision vs Comprehensiveness

# How to Automate Knowledge Graph Construction?

- Sound **Knowledge Graph Construction** relies on **Ontologies**

- Ontologies don't come for free, i.e. Ontology Design is very expensive wrt. time and resources

- Ontologies can be „learned" automatically

- **Ontology Learning** defines a set of methods and techniques
  - for **fundamental development** of new ontologies
  - for **extension or adaption** of already existing ontologies
- in a (partly) automated way from various resources.

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Fundamental Types of Ontology Learning

- Ontology Learning from Text
  - automatic or semi-automatic generation of lightweight ontologies by means of text mining and information extraction

- Linked Data Mining
  - detecting meaningful patterns in RDF graphs via statistical schema induction or statistical relational learning

- Concept Learning in Description Logics and OWL
  - learning schema axioms from existing ontologies and instance data mostly based on Inductive Logic Programming

- Crowdsourcing Ontologies
  - combines the speed of computers with the accuracy of humans, as e.g. taxonomy construction via Amazon Turk or games with a purpose

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Ontology Learning from Text

- **Ontology Learning from text** is the process of identifying terms, concepts, relations, and optionally axioms from textual information and using them to construct and maintain an ontology.

- Automatisation requires help from
  - Natural Language Processing (NLP)
  - Data Mining
  - Machine Learning techniques (ML)
  - Information Retrieval (IR)

# Ontology Learning from Text - Basic Approach



**document corpus**

**terminology**

**ontology**

**(1) term extraction**

<dog> <dogs>
<cat>
<siamese cat>

**(2) conceptualisation**

pet

dog        cat

siamese cat

**(3) evaluation & adaption**

**term extractions** requires linguistic processing (NLP) to identify important noun phrases and their internal semantic structure

**terms**: linguistic realisations of domain specific concepts

**Concepts**: clusters of semantically related terms

# The Ontology Learning Layer Cake

Country $\sqsubseteq$ ≤1 hasCapital.$\top$

General Axioms

River $\sqcap$ Mountain $\sqsubseteq$ $\bot$

Axiomatic Schemata

capitalOf $\sqsubseteq$ locatedIn

Relation Hierarchies

flowThrough(dom:River, range:GeoEntity)

Relations

Capital $\sqsubseteq$ City , City $\sqsubseteq$ InhabitedGeoEntity

Concept Hierarchies

c:=country:=<description(c), uri(c)>

Concept Description

{country, nation, land}

Multilingual Synonyms

river, country, nation, city, capital, ...

Terms

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Ontologies are NOT the Reality

- Ontologies are a **context-dependent projection (model) of the Reality**

- **Different ontologies** might represent the **same (or similar) knowledge**, as e.g. ontologies might
  - reflect different tasks and requirements for applications
  - follow different conventions or restrictions

[1]

# How Ontologies can Differ

- The **same term describes different concepts**
  - e.g. Author - *writer of a book vs. creator of a document*

- **Different terms describe the same concept**
  - e.g. Author vs. Writer

- **Different modeling conventions and paradigms**
  - e.g. intervals vs. points - *to describe temporal aspects*

- **Different level of granularity**
  - e.g. Fiction vs. PoliticalFiction, ScienceFiction, RomanticFiction, *etc. as literary Genres*

- **Different coverage or different point of view,** etc.

# Ontology Alignment

- **Ontology Alignment** or **Ontology Matching** is the process of determining *correspondences* between ontological concepts

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Correspondence or Mapping

- Given the ontologies $O_1$ and $O_2$, a **correspondence** or **mapping** among the entities $e_1$ and $e_2$ from $O_1$ and $O_2$ respectively, is defined as

$$\langle id, e_1, e_2, r, n \rangle$$

- with
  - **id** ... a unique **identifier** of the correspondence
  - **r** ... a **relation**, as e.g. equivalence (**=**), more general (**⊒,≥**), less general (**⊑,≤**), disjointness(**⊥**), part-of, etc...
  - **n** ... a **confidence measure** (typically in the range of **[0,1]**) holding for the correspondence between $e_1$ and $e_2$
- the correspondence$\langle id, e_1, e_2, r, n \rangle$asserts that the relation *r* holds between the entities $e_1$ and $e_2$ with confidence *n*

# Complexity of Correspondences

- Examples of **simple correspondences**:

    - `http://dbpedia.org/resource/Joseph_Fourier =`

        `https://www.wikidata.org/wiki/Q8772`

    - `Author = Writer`

    - `Gas ≥`$_{1.0}$ `GreenhouseGas`

    - `rdfs:label ≥`$_{0.9}$ `dc:title`

# Complexity of Correspondences

- Examples of **more complex correspondences**:

  ○         $\texttt{speed = velocity × 2.237}$
     $\texttt{0.477 × speed = velocity}$

  ○ $\texttt{Book(x) ∧ author(x,y) ∧ Writer(y) ⇒}_{.85}$
     $\texttt{writtenBy(x,concat(y.firstname, y.lastname))}$

# Alignment Example

Book =1.0 Volume

id ≥0.9 isbd

Person =0.9 Human

name ≥1.0 title

author =1.0 author

Science ≤0.9 Essay



*Jérôme Euzenat, Pavel Shvaiko: Ontology Matching, Springer, 2007, p.48*

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Ontology Matching Techniques

- **Element-level Ontology Matching Techniques** consider ontology entities or their instances in isolation from their relations with other entities or their instances

    - **String-Based** - *matching names or descriptions of entities*

    - **Linguistic-Based** - *use NLP, lexicons, or domain specific thesauri to match words based on linguistic relations (homonymy, synonymy, partonomy, etc.), or exploiting morphological properties*

    - **Constrained-Based** - *take into account internal constraints applied to the definitions of entities, as e.g. types, cardinality of properties, etc.*

    - **Extensional-Based** - *use individual representation of classes, i.e. classes are considered similar if they share many instances*

# Ontology Matching Techniques

- **Structure-level Ontology Matching Techniques** consider ontology entities or their instances to compare their relations with other entities or their instances

  - **Graph-Based** - *consider ontologies as labeled graphs, assumption: if nodes are similar, then also their neighbors must be similar*

  - **Taxonomy-Based** - *like graph-based algorithms, but consider only specialization/generalization relation*

  - **Method-Based** - *take into account semantic interpretation of the ontologies, assumption: if two entities are the same, then they share the same interpretation*

  - **Data Analysis and Statistics** - *take a large sample, try to find regularities, discrepancies, allows grouping or determining distance metrics, ...*

# Ontology Alignment



*Euzenat, Shvaiko: Ontology Matching, Springer 2007*

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# SEMANTIC SEARCH AND RECOMMENDATIONS

Would graphs help search engines?

# The Information Retrieval Dilemma



- Ambiguity of natural language (polysemy)
- Different words/expressions for the same concept (synonyms, metaphors, paraphrases,...)

# The Information Retrieval Process

Evaluation

User Interaction

Ranking

Document Store

Index

Retrieval Process

Indexing Process

Text Acquisition

Text Transformation

Index Creation

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Knowledge Graph Supported Retrieval Process

[1]

Slide from Knowledge Graphs course by prof. Harald Sack & Mewish Alam (openHPI, 2020).

# Knowledge Graph Supported Retrieval Process

- **Prerequisite:**
  **Document Annotation** with explicit semantics, e.g. semantic entities



Example for Linked Data Based Document Annotation

http://scihi.org/neil-armstrong/

- Enables **entity-based Information Retrieval**
  - Language independent

[1]

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Entity Based Search

**Query Processing:**

Armstrong on the Moon

**Named Entity Linking**

dbr:Neil Armstrong    dbr:Moon

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Indexing:**

**The first Man on the Moon**

....
On the Moon, the 38-year-old civilian commander, radioes to earth and the mission control room here: "Houston, Tranquility Base here, The Eagle has landed."
....

dbr:Neil Armstrong

dbr:Moon

**Named Entity Linking**

**Entity-Based Query Matching**

- **simple entity matching**
- similarity-based entity matching
- relationship-based entity matching
- ...

[1]

# Entity Based Search

**Query Processing:**

Armstrong on the Moon

**Named Entity Linking**

dbr:Neil Armstrong    dbr:Moon

---

**Indexing**

> The 2nd Man on the Moon
>
> ….
> Legendary astronaut Buzz Aldrin has revealed some captivating pieces of Apollo 11 memorabilia on social media in the last few days.
> ...

**dbr:Moon**

dbr:Buzz_Aldrin

↕ semantic similarity

**dbr:Neil_Armstrong**

**Named Entity Linking**

**Entity-Based Query Matching**

- simple entity matching
- **similarity-based entity matching**
- relationship-based entity matching
- ...

Two entities are considered **semantically similar**
- if they share property/value pairs
- if they share properties with similar values

# Entity Based Search

**Query Processing:**

Armstrong on the Moon

**Named Entity Linking**

dbr:Neil Armstrong          dbr:Moon

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Indexing**

The 2nd Man on the Moon

….
Legendary astronaut **Buzz Aldrin** has revealed some captivating pieces of Apollo 11 memorabilia on social media in the last few days.
...

**dbr:Moon**

dbo:Astronaut

dbr:Apollo_11

rdf:type
dbo:mission

**dbr:Neil_Armstrong**

**Named Entity Linking**

**Entity-Based Query Matching**

- simple entity matching
- similarity-based entity matching
- **relationship-based entity matching**
- ...

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

[1]

# Retrieval vs. Exploration

# The Retrieval Problem

- **Retrieval Problem:**
  - you are looking for **something specific**
    i.e. you know what you are looking for

- How to **specify your search request**?
  - e.g. for a (specific) book:
    *author name, title, etc.*

- Often you are using
  - (unique) identifier
  - descriptive metadata

[3]

*Author: Jules Verne*
*Title:    From the Earth to the Moon*

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# The Retrieval Problem



[2]

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Retrieval vs. Exploration

- *Find another („comparable") book, (that will be of interest for me...)*
- *Find books of the same or of related topics*
- *How did the author / the topic develop over time?*
- *What else would I like to read?*
- *...*

Exploratory Search

[2]

...and intelligent recommendations

Traditional Libraries enable Exploratory Search

[1]

# Exploratory Search

represents the activities carried out by searchers who are:

- unfamiliar with the domain of their goal (i.e. need to learn about the topic in order to understand how to achieve their goal),
- unsure about the ways to achieve their goals (either the technology or the process),
- or even unsure about their goals in the first place.

- ...**Browsing** instead of **Searching**
- ...to find something by chance, i.e. **Serendipity**
- ...to get an **overview**
- ...enable content based **navigation**

# Exploratory Search and Recommendation

# Exploratory Search via Knowledge Graphs



http://dbpedia.org/resource/From_the_Earth_to_the_Moon

# Exploratory Search via Knowledge Graphs

**:From_the_Earth_to_the_Moon**

**:Jules_Verne**

**dbo:Book** ← rdf:type

dbo:author

dct:subject

dbo:influenced

**:H._G._Wells**

category:1865_novels
category:Frence_science_fiction_novels
category:Novels_by_Jules_Verne
category:Moon_in_fiction
category:Fictional_rivalries
category:Novels_set_in_Florida
category:1860s_science_fiction_novels
...

dbo:previousWorkOf

**:In_Search_of_the_Castaways**

# Exploratory Search via Knowledge Graphs

**:From_the_Earth_to_the_Moon**



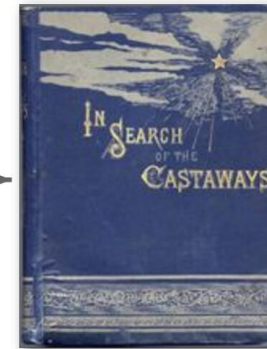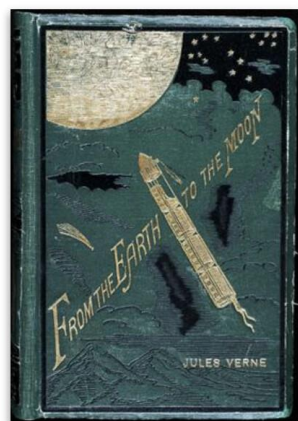rdf:type → **dbo:Book** ← rdf:type

dbo:subsequentWorkOf

rdf:type

**:A_Journey_to_the_Center_of_the_Earth**

dbo:previousWorkOf

**:In_Search_of_the_Castaways**

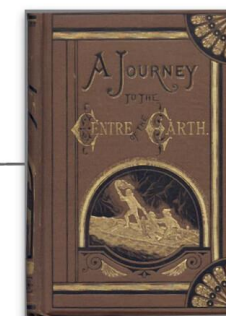# Exploratory Search via Knowledge Graphs
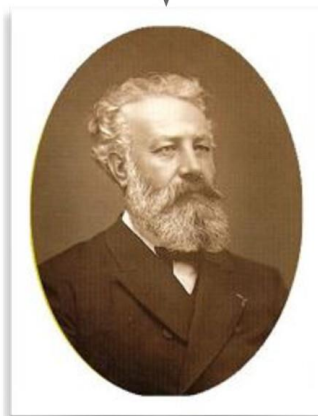


**:From_the_Earth_to_the_Moon**

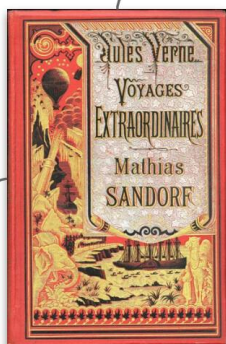**dbo:Book**

rdf:type

rdf:type

rdf:type

rdf:type

rdf:type

**:A_Journey_to_the_Center_of_the_Earth**

**:The_Mysterious_Island**

dbo:author

dbo:author

dbo:author

dbo:author

dbo:author

**:Matthias_Sandorf**

**:Jules_Verne**

**:Master_of_the_World_(novel)**

# Exploratory Search via Knowledge Graphs



:From_the_Earth_to_the_Moon

rdf:type → dbo:Book

rdf:type

rdf:type

rdf:type

:The_Invisible_Man

dbo:author

:The_Island_of_Doctor_Moreau

dbo:author

dbo:author → dbo:Writer

rdf:type

rdf:type

dbo:author

:Jules_Verne

dbo:influenced

:H._G._Wells

:The_War_of_the_Worlds

# Exploratory Search via Knowledge Graphs

- **Exploratory Search** represents the activities carried out by searchers who are either:
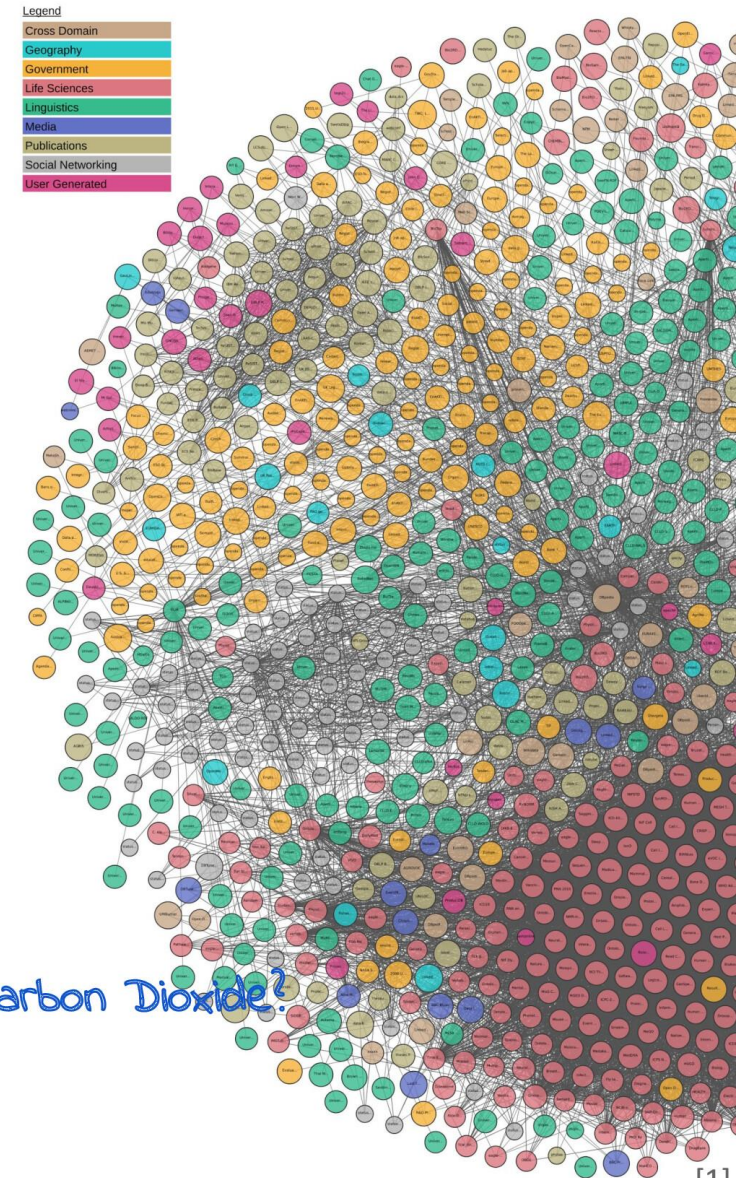  - **unfamiliar with the domain** of their goal (i.e. need to learn about the topic in order to understand how to achieve their goal),
  - **unsure about the ways** to achieve their goals (either the technology or the process)
  - or even **unsure about their goals** in the first place.

- **Recommender Systems** seek to predict the preference a user would give to an item.
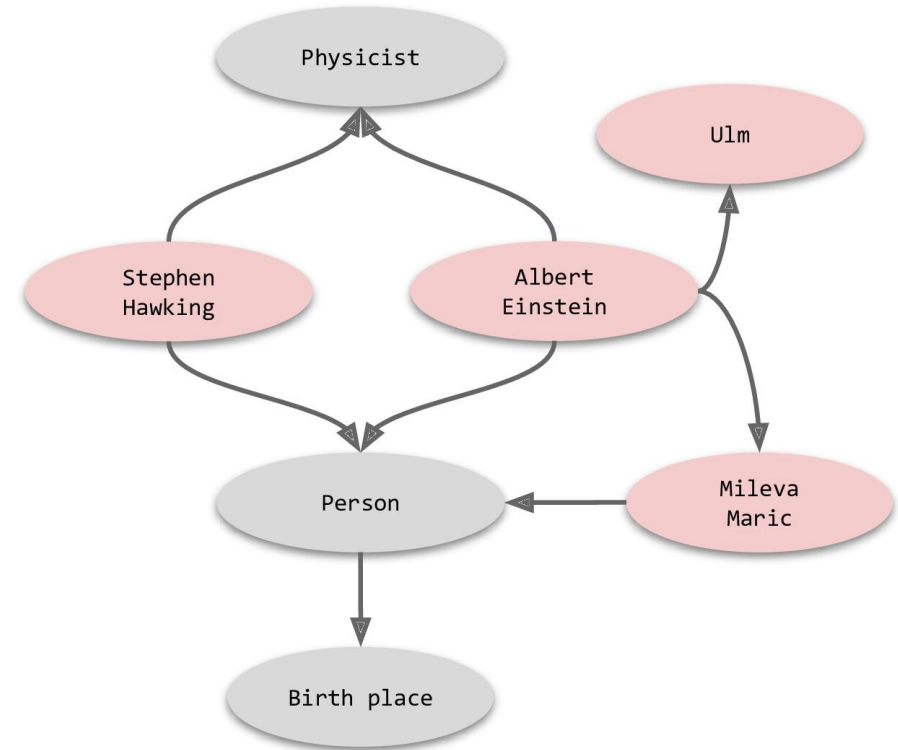
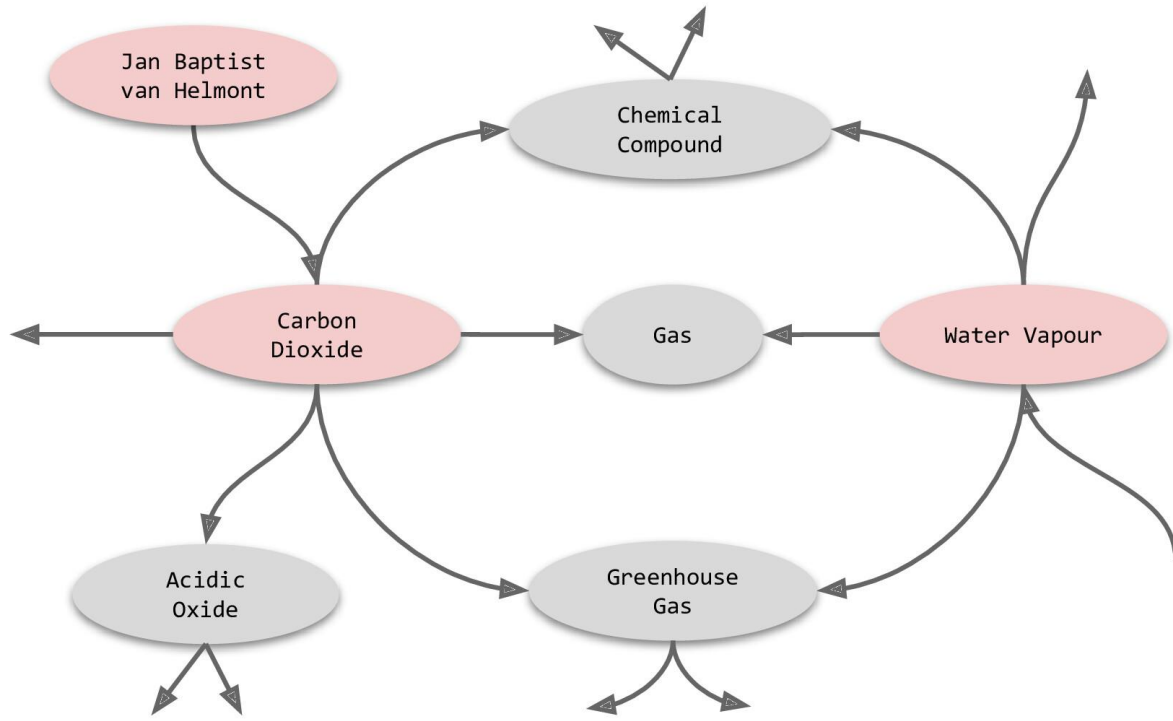# KNOWLEDGE GRAPH EMBEDDINGS

The graphs are vectors if you need it
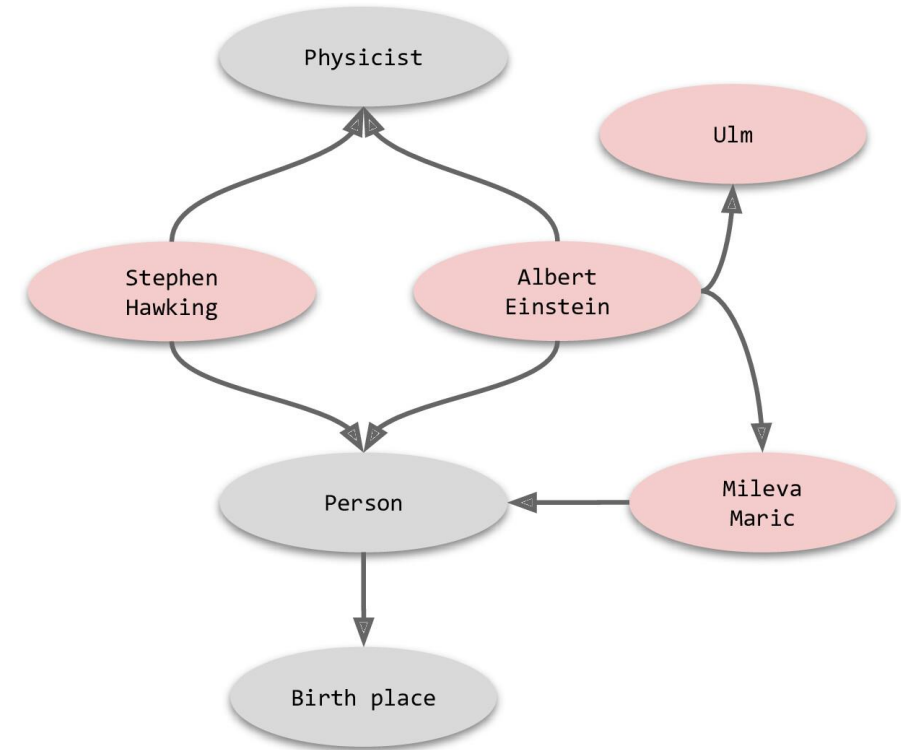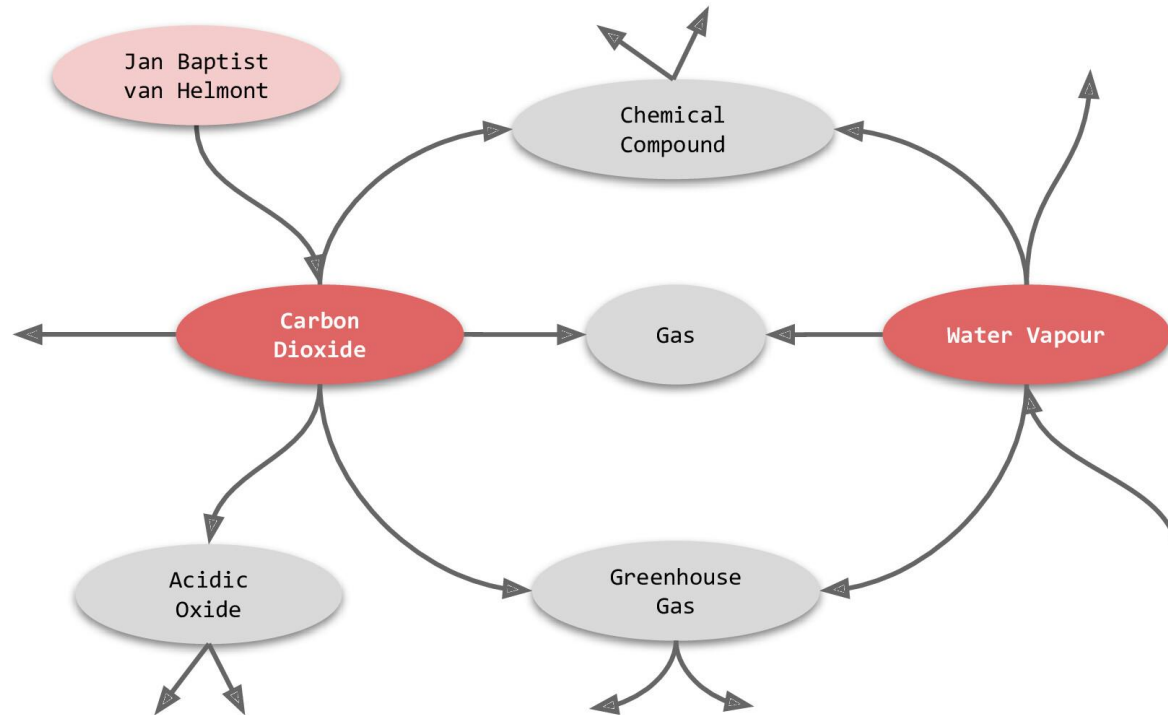
# Semantic Similarity

- For search and retrieval systems, **semantic similarity of entities** is an important feature, as e.g.
  - Given an entity find the most similar entities
  - Given an entity find the most similar documents
  - Given a document find the most similar documents, etc.
- **When are two entities (semantically) similar?**
  - If they can be described by the same/similar facts, as e.g.
  - Carbon Dioxide is a Greenhouse Gas and water vapour is a Greenhouse Gas
  - Albert Einstein is a Physicist and Stephen Hawking is a Physicist
  - Is Stephen Hawking more similar to Albert Einstein or to Carbon Dioxide?

Legend
Cross Domain
Geography
Government
Life Sciences
Linguistics
Media
Publications
Social Networking
User Generated

The Linked Open Data Cloud from lod-cloud.net

[1]

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Semantic Similarity

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Semantic Similarity
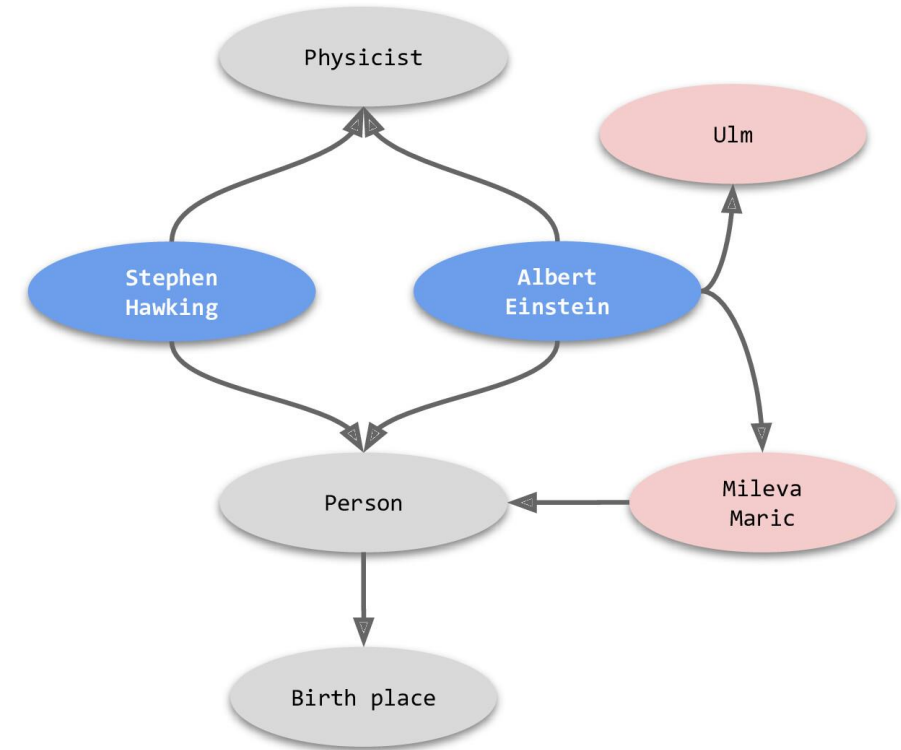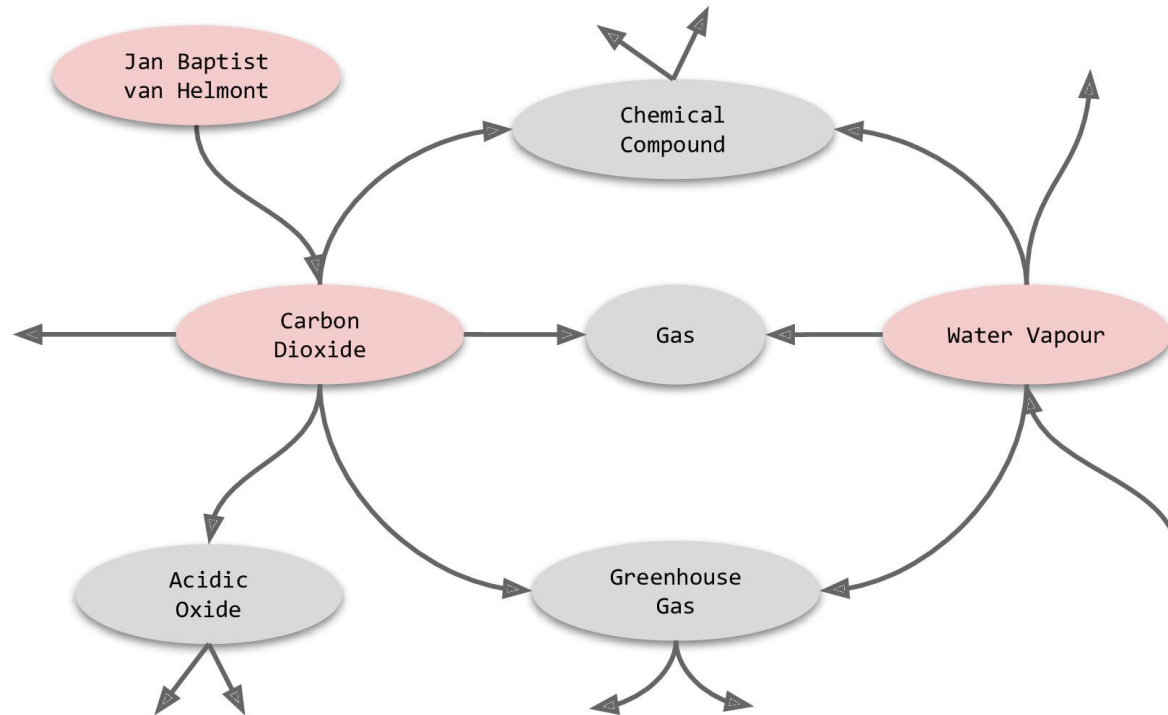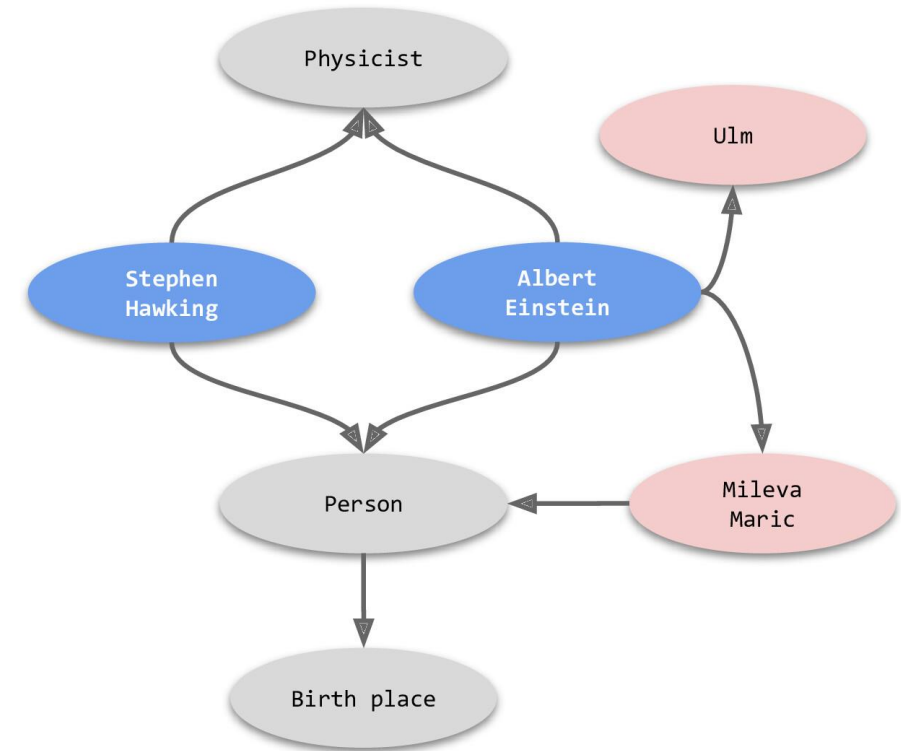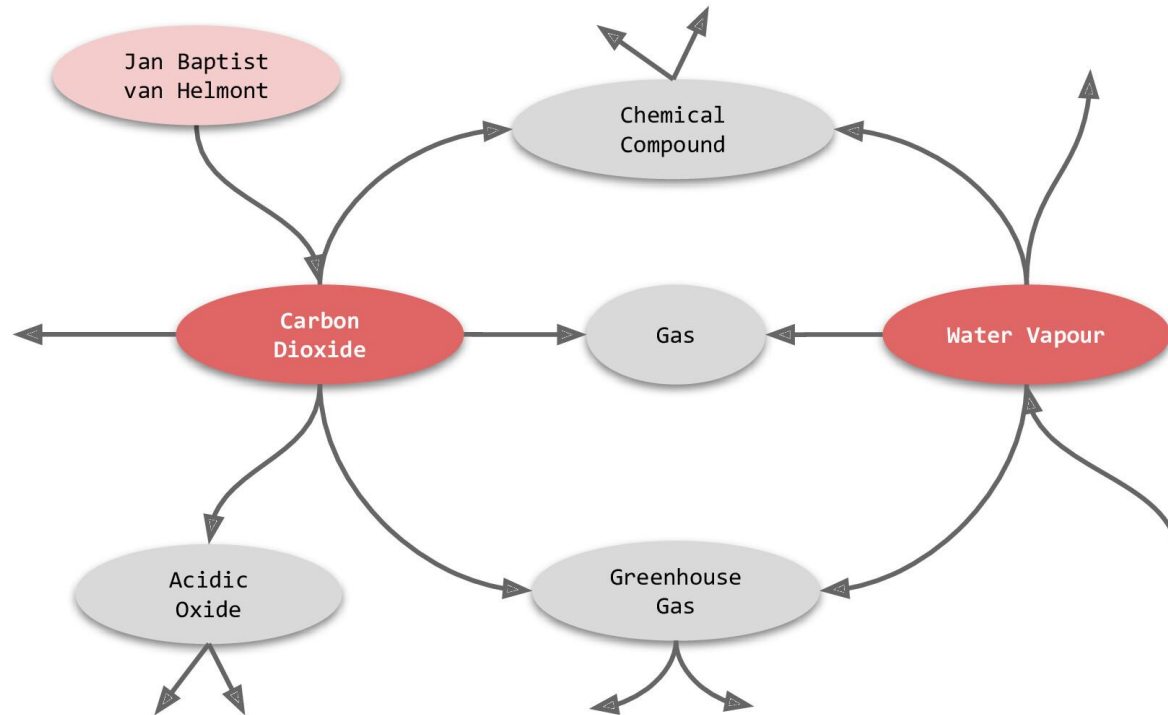


○ Carbon Dioxide and water vapour share similar (structural) context in the graph

# Semantic Similarity



○ Stephen Hawking and Albert Einstein share similar (structural) context in the graph

# Semantic Similarity



○ "You shall know a node by the company it keeps"

○ i.e. similar nodes can be identified by having the same/similar environment (context)

○ adjacency based similarity

# Semantic Similarity

- In a Knowledge Graph,
    - **similar entities** are represented by nodes that are connected to **similar/same facts**
    - i.e. that are connected to **similar graph structures**
    - To identify **similar entities**, we have to identify **similar graph structures**

- **Problem:**
    - Algorithms to determine semantic similarity in graphs are of high complexity, i.e. with large KGs, as e.g. Wikidata, they don't work efficiently.

- **Idea**:
    - Approximate the problem by transferring it from graph structures to vector spaces That are easier to handle.

From Nodes and Edges ...

# ... To Semantically Meaningful Vector Representations

# Excursion: Word Embeddings

- **Word Embeddings** map natural language words to a dense vector representation

- **Basic Assumption:** Similar words occur in similar contexts:

  (Carbon Dioxide, Water Vapour, Methane) is one of the driving agents of climate change.
  Climate change is caused by greenhouse gases like (Carbon Dioxide, Water Vapour, Methane)

- **Basic idea:** instead of counting co-occurrences of words, predict the likelihood of the appearance of words in the neighborhood of others

- Train a predictor (neural network) that can predict a word from its context (**CBOW**) or the context from a given word (**Skip Gram**)

Mikolov, Tomas; et al. (2013). "Efficient Estimation of Word Representations in Vector Space". arXiv:1301.3781

# Excursion: Word Embeddings

- **Skip Gram:**

  - Train a neural network with one hidden layer
  - Use output at hidden layer as vector representations

- **Observation:**

  - *Carbon Dioxide, Water Vapour, Methane* will activate similar context words
  - i.e. their output weights at the projection layer have to be similar



Mikolov, Tomas; et al. (2013). "Efficient Estimation of Word Representations in Vector Space". arXiv:1301.3781

# Word Embeddings



Male-Female     Verb tense     Country-Capital

- Semantics of words is preserved, i.e. it enables semantic arithmetic operations as e.g. analogies
  - "king" - "man" ≈ "queen" - "woman"
  - "king" - "man" + "woman" ≈ "queen"

Slide from Knowledge Graphs course by prof. Harald Sack & Mewish Alam (openHPI, 2020).

# Graph Embeddings

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Graph Embeddings - Encoder-Decoder Approach



- The goal is to encode the nodes of the graph in a way so that **similarity in the embedding space** (e.g., dot product) **approximates similarity in the original network**.

- $ENC: N \rightarrow \mathbb{R}^d$ , $u,v \in N$, $ENC(u) = z_u \in \mathbb{R}^d$, $ENC(v) = z_v \in \mathbb{R}^d$

- $DEC: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^+$ , $DEC(ENC(u), ENC(v)) = DEC(z_v , z_u) \approx$ similarity $(u,v)$

# Learning Graph Embeddings

1) Define an **encoder ENC** (i.e., a mapping from nodes to embeddings)

2) Define a **node similarity function** that specifies how relationships in vector space map to relationships in the original network.

3) Optimize the parameters of the encoder so that:

$$\text{similarity}(u, v) = z_v^T \, z_u$$

# Knowledge Graph Embeddings

Many ways to generate Knowledge Graph Embeddings:

- **Translational Methods**: TransE, TransH, TransR, TransEdge, …

- **Rotation Based**: RotatE

- **Graph Convolutional Networks**: R-GCN, TransGCN

- **Walk-Based Methods**: DeepWalk, RDF2Vec

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

## Translational Distance Models

- Exploit distance-based scoring functions

- Measure the **plausibility of a fact** as the **distance between two entities**

- A translation carried out by the relation.

- **Models**: TransE, TransH, TransR, TransD, TransSparse, TransM, TransEdge

Wang et al., Knowledge graph embedding: A survey of approaches and applications. IEEE Transactions on Knowledge and Data Engineering (TKDE), 2017.

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# TransE

- Entities and relations are embedded into **same vector space**.
- h = head, t = tail, r = relation
- Relation r is considered as translation from h to t
- Learning Assumption **h+r≈t**
- **Problem:** Symmetric functions, 1-N / N-1 / N-N functions



Entity and Relation Space

Bordes et al, Translating Embeddings for Modeling Multi-relational Data, NIPS 2013

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# TransH

- From original space to Hyperplane

- TransH enables **different roles of an entity in different relations**.

- Entities h and t are projected into specific **hyperplane of relation r**.

- Then predict new links based on translation on hyperplane.



Entity and Relation Space

Wang et al., Knowledge graph embedding by translating on hyperplanes. AAAI, 2014.

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Graph Convolutional Network

- **Graph Convolutional Networks (GCN)**
  - modeling structured neighborhood information of **unlabeled** and **undirected** graphs with **convolution operations**

- **Relational Graph Convolutional Network (R-GCN)**
  - Models Relational Data using GCN where Knowledge Graphs are considered as **directed labeled multigraphs**.
  - Models in RGCN
    - **Link Prediction:**
      - **an encoder:** an R-GCN producing latent feature representations of entities,
      - **a decoder:** a tensor factorization model exploiting these representations to predict labeled edges

# RDF2Vec

- Word2vec operates on sentences, i.e. sequences of words
- **RDF2Vec Basic Idea**:
  - Generate "sentences" from knowledge graph, i.e. sequences of interconnected RDF triples

```
:CarbonDioxide rdf:type :GreenhouseGas.
:GreenhouseGas, rdf:type, :Gas.
:Gas, rdf:type, :FundamentalStateOfMatter.
```

  - Selection strategies:
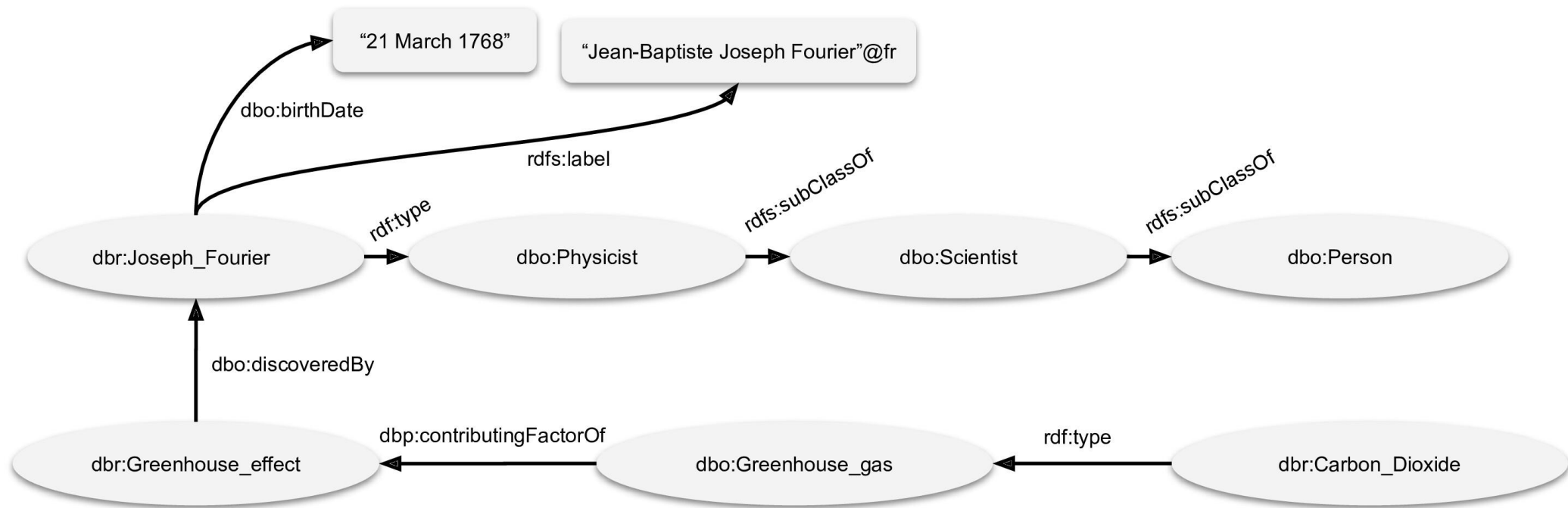    - Depth first search
    - Breadth first search
    - Random walk
    - RDF Graph Kernels

Petar Ristoski and Heiko Paulheim RDF2Vec: RDF graph embeddings for data mining, ISWC 2016

# Graph Walks RDF2Vec



**Generated Sequences of depth = 3:**

- dbr:Carbon_Dioxide→ rdf:type→dbo:Greenhouse_gas → dbp:contributingFactorOf → dbr:Greenhouse_effect
  → dbo:discoveredBy → dbr:Joseph_Fourier

# Libraries for KG Embedding



https://github.com/facebookresearch/PyTorch-BigGraph



https://github.com/Accenture/AmpliGraph



*PyKeen*

https://github.com/SmartDataAnalytics/PyKEEN

*OpenKE*

http://openke.thunlp.org/

Knowledge Graphs 2020 , Prof. Dr. Harald Sack & Dr. Mehwish Alam, FIZ Karlsruhe - Leibniz Institute for Information Infrastructure & Karlsruhe Institute of Technology

# KNOWLEDGE GRAPH COMPLETION

How to guess the missing triples?

# Knowledge Graph Refinement

- As a model of the real world or a part of it, **knowledge graphs cannot reasonably reach full coverage**, i.e., contain information about each and every entity in the universe.

- **It is unlikely,** in particular if heuristic methods are applied for knowledge graph construction, **that the knowledge graph is fully correct**.

- To address those shortcomings, various methods for **Knowledge Graph Refinement** have been proposed, as e.g.

  - Deduplicating entity nodes (entity resolution)
  - Collective reasoning (probabilistic soft logic)
  - **Link prediction** or **Knowledge Graph Completion**
  - Dealing with missing values
  - Anything that improves an existing knowledge graph

# Completion vs. Error Detection

- **Knowledge Graph Completion:**
  Adding missing knowledge to the Knowledge Graph

  E.g. adding a triple:
  *<JosephFourier, occupation, Physicist>*

- **Error Detection:**
  Identifying wrong information in the Knowledge Graph

  E.g. finding inconsistencies:
  *<JosephFourier, isA, Human>*
  *<JosephFourier, isA, FictionalCharacter>*

# Knowledge Graph Completion

- A promising approach for **Knowledge Graph Completion** is
  - to embed Knowledge Graphs into latent spaces (via Knowledge Graph Embeddings) and
  - make inferences by learning and operating on latent representations.

- Such embedding models, however, **do not make use of any rules** during inference and hence have limited accuracy.

- E.g. predict that in Wikidata the following fact may be complemented:

  *(AtsumoOmuhura occupation Climatologist)*
  `wd:Q462297 wdt:P106` `wd:Q1113838` .

  — Tail Prediction

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).

# Link Prediction

| | Task | Example | Result |
|---|---|---|---|
| **Link Prediction** | Triple Classification | (JosephFourier, occupation, physicist)? | (yes, 95%) |
| | Tail Prediction | (JosephFourier, occupation, ?) | (1, physicist, 0.95), (2, chemist, 0.93) … |
| | Head Prediction | (?, occupation, physicist) | (1, AlbertEinstein, 0.91) (2, StephenHawking, 0.90) |
| | Relation Prediction | (JosephFourier, ?, physicist) | (1, occupation, 0.95) |
| | Entity Classification (Type Prediction) | (JosephFourier, isA, ?) | (1, Person, 0.99) (2, Human, 0.99),… |

# Type Prediction

- **Predicting a type or class** for an entity given some characteristics of the entity is a very common problem in machine learning, known as **classification**.
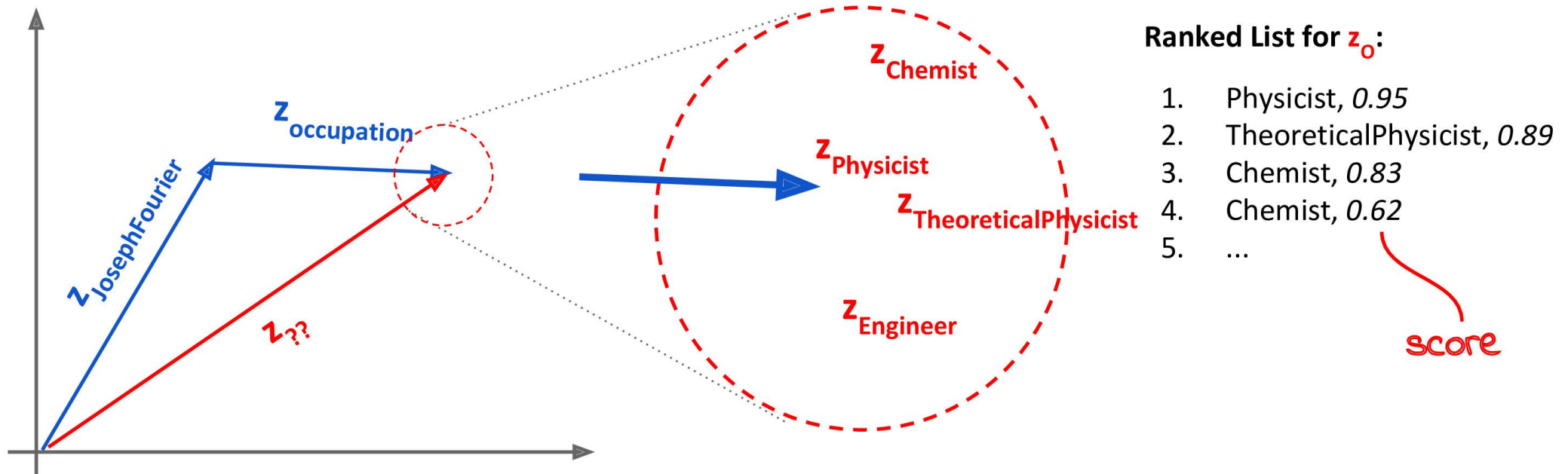
<p align="center"><code>&lt;JosephFourier, isA, ?&gt;</code></p>

- **Supervised Learning Approach**:
  - Type Prediction can be addressed via a **classification model** based on **labeled training data**,
  - typically the set of entities in a Knowledge Graph which have types attached.

# Type Prediction

- **Multi-Class Prediction:**
  - In Knowledge Graphs usually there are more than two types/classes of entities to distinguish
    E.g. Classes Physicists, Chemists, Climatologists, etc.

- **Single-Label Classification:**
  - Only one type/class can be assigned per entity
    E.g.:   `<JosephFourier, isA, Person>`

- **Multi-Label Classification:**
  - In Knowledge Graphs some entities might allow for the assignment of more than one type
    E.g.:       `<electron, isA, Particle>` and
               `<electron, isA, Wave>`

# Methods for Knowledge Graph Link Prediction

- Use **Translational Embeddings**
  - **Unsupervised** methods, e.g. **TransE**, use $z_s + z_p$ to predict $z_o$
  - **Supervised** Methods for prediction based on embedding vectors



Ranked List for $z_o$:

1. Physicist, *0.95*
2. TheoreticalPhysicist, *0.89*
3. Chemist, *0.83*
4. Chemist, *0.62*
5. ...

Slide from Knowledge Graphs course by prof. Harald Sack & Mehwish Alam (openHPI, 2020).
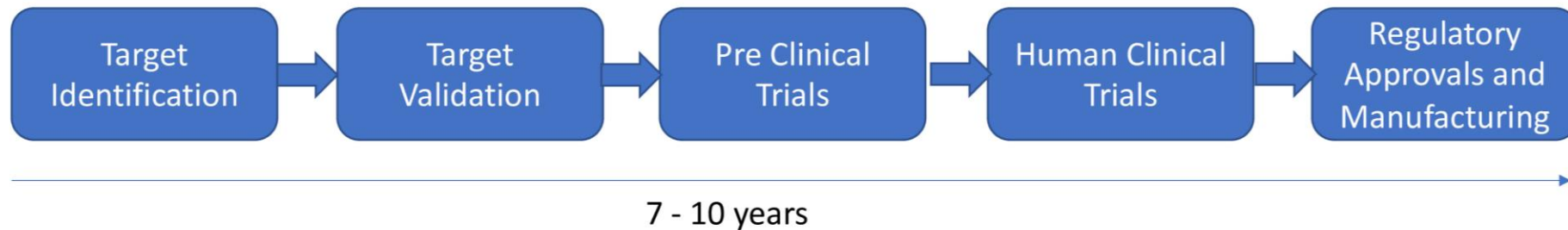
Slide from Knowledge Graph Embeddings Tutorial by L. Costabello et al. (ECAI 2020).
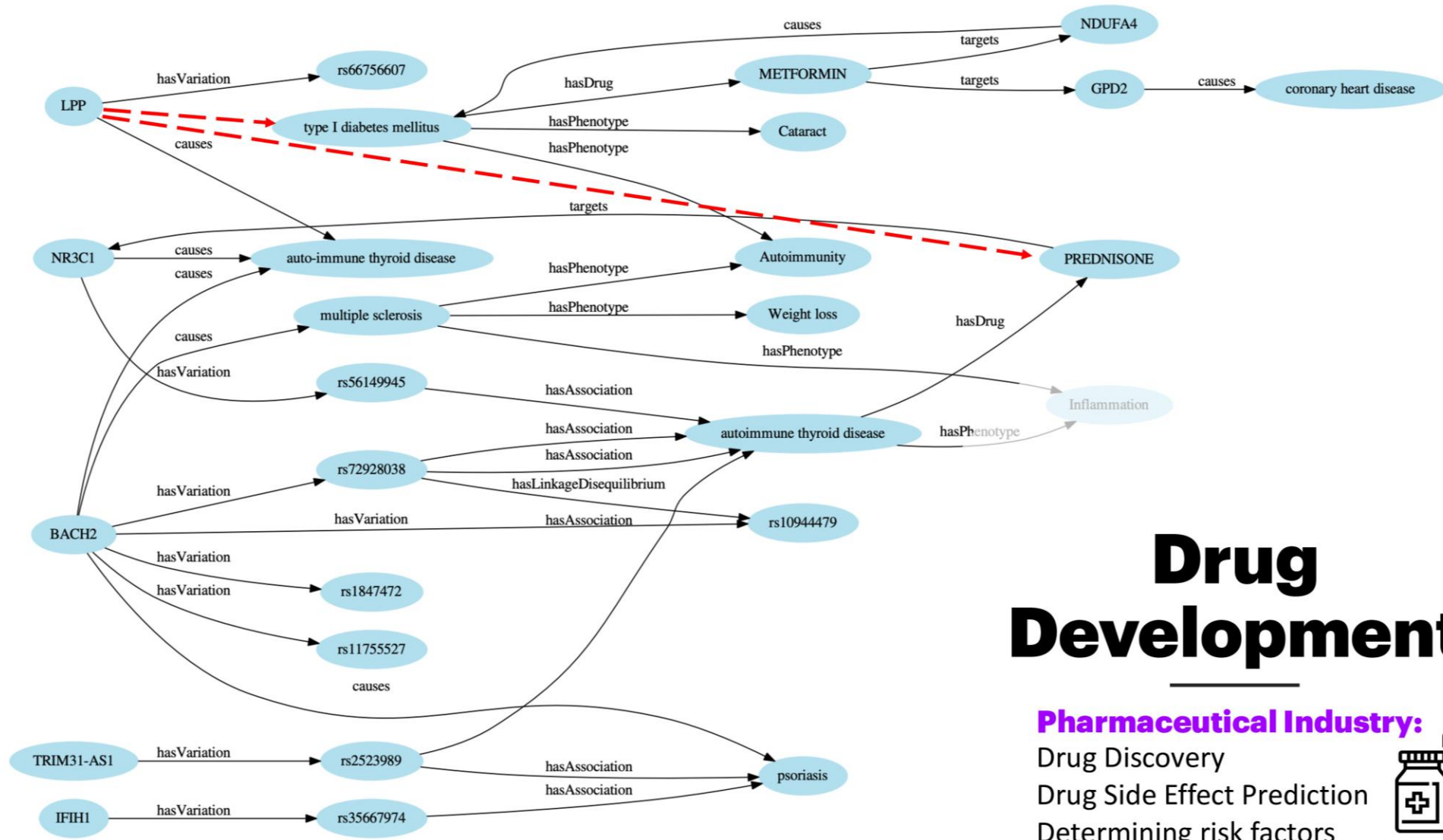
# Drug Development

- Drug Development is a time consuming and expensive process which ranges from gene identification, identifying a compound to target the gene, and finally experimentation on animals and humans.



7 - 10 years

- The initial step of identification of gene/drug takes several years and if not identified correctly may result in loss of time and money.
- "Drug Developers" identify the genes/drugs by reading the latest research before proceeding with experimentation. But it is highly dependent on the experience of the person.

Slide from Knowledge Graph Embeddings Tutorial by L. Costabello et al. (ECAI 2020).
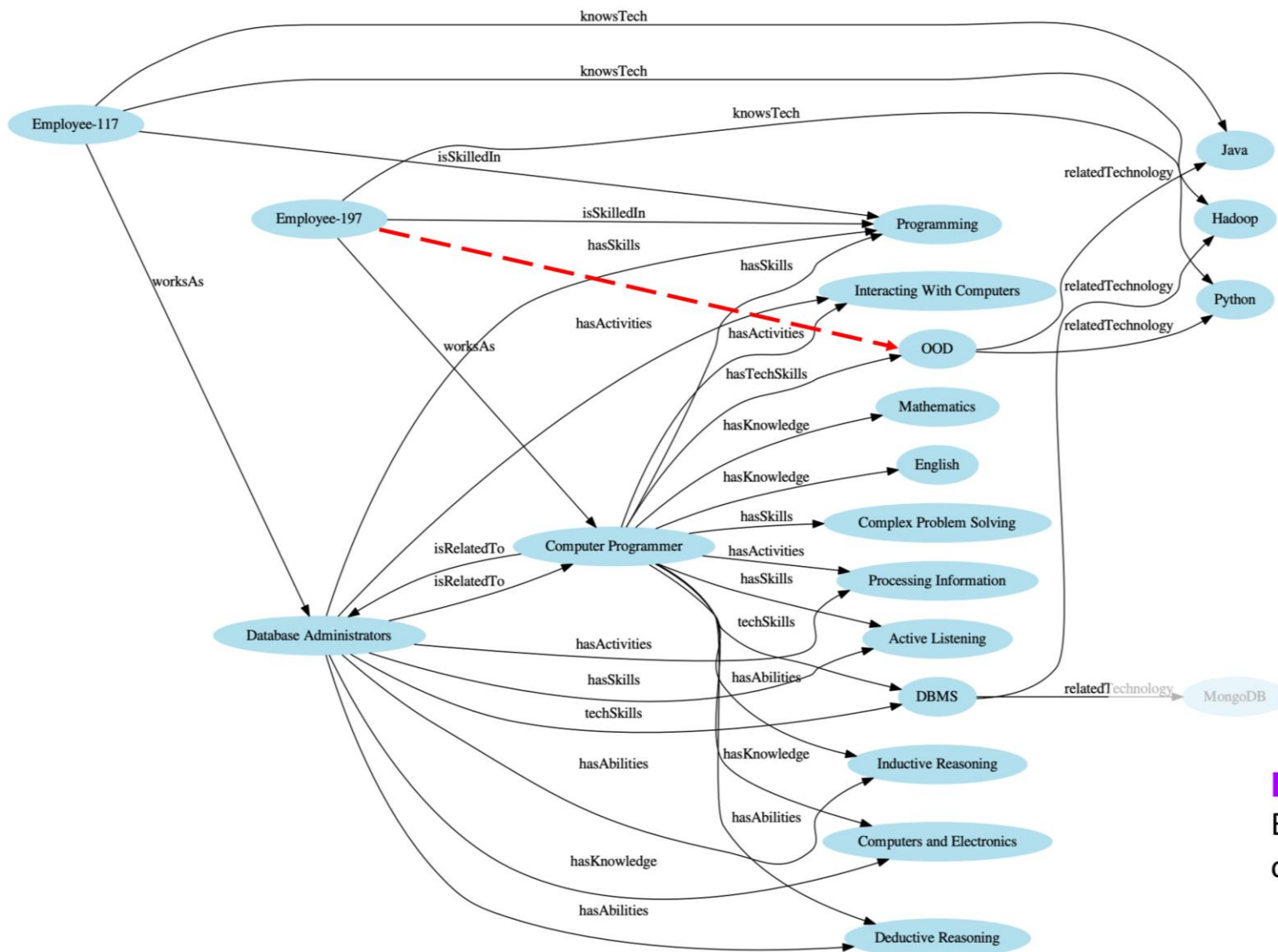
# Human Resource

- Technology is evolving at an extremely fast pace. People need to learn new skills to be relevant in the market.

- Due to automation, a lot of roles are becoming obsolete and companies are forced to lay off people.

KGEs can be used for following tasks:

- Suggest new technology/tasks for career progression.

- Recommend similar roles within the organization when existing role becomes obsolete.

# Human Resource

**Human Resources:**
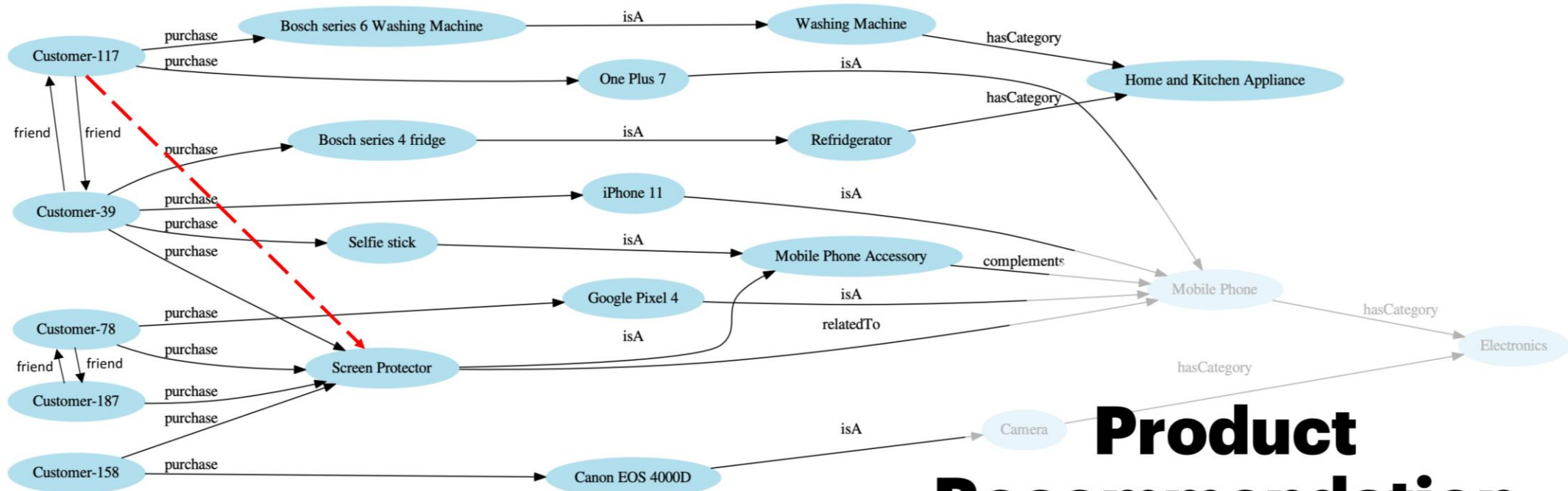Employee Career Progression or Transition

# Product Recommendation

KGEs can leverage relation between customers and products.

KGEs can be used for following tasks:
- Recommend new products to customers
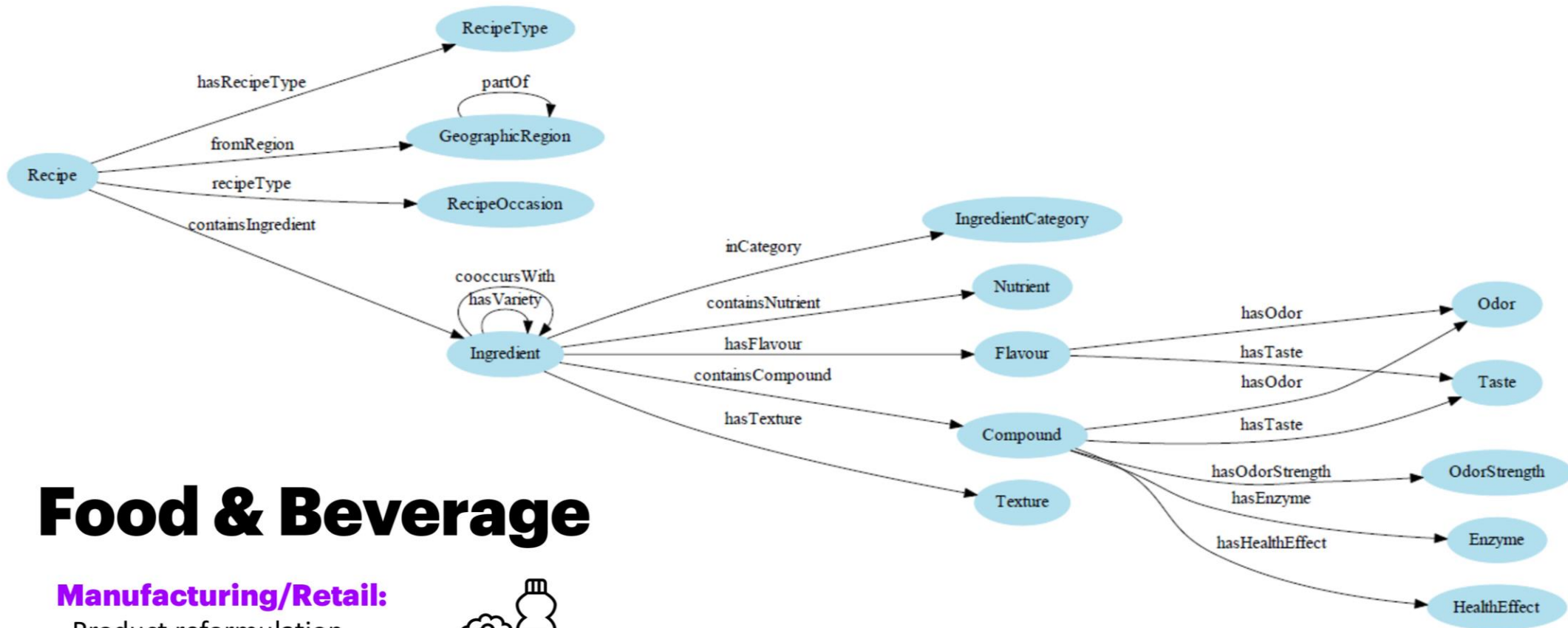- Group customers based on their purchase history

Slide from [Knowledge Graph Embeddings Tutorial](#) by L. Costabello et al. (ECAI 2020).

Slide from Knowledge Graph Embeddings Tutorial by L. Costabello et al. (ECAI 2020).

# Food & Beverage

**Manufacturing/Retail:**
- Product reformulation
- Adapting to consumer trends

Slide from Knowledge Graph Embeddings Tutorial by L. Costabello et al. (ECAI 2020).

**Product Reformulation:**
- Item substitution based on embedding proximity

Slide from Knowledge Graph Embeddings Tutorial by L. Costabello et al. (ECAI 2020).

# Item Recommendation

- Use vector algebra to find latent region that satisfy input criteria

- Example:
  - "I want *Indian recipes* that contain garlic and tomato"

  - $nearest(\ avg\ (avg(GARLIC, TOMATO) - containsIngredient, India - recipeOrigin))$

- *Note: the above is pseudo-code, actual solutions will depend on model, data, fine-tuning, etc*
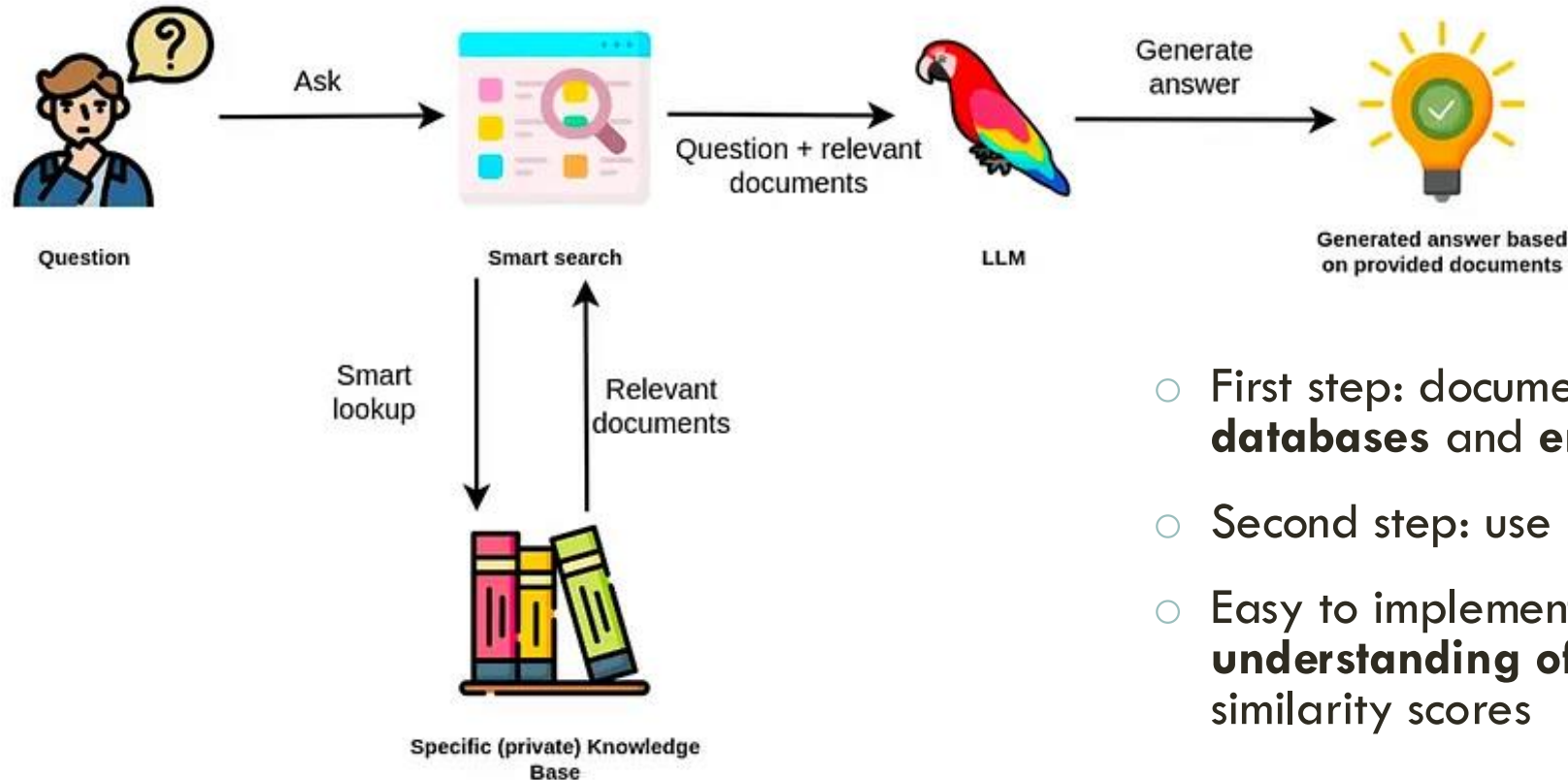
- Alternatively use Bayesian optimization ..

# CHATGPT IS A BULLSHIT
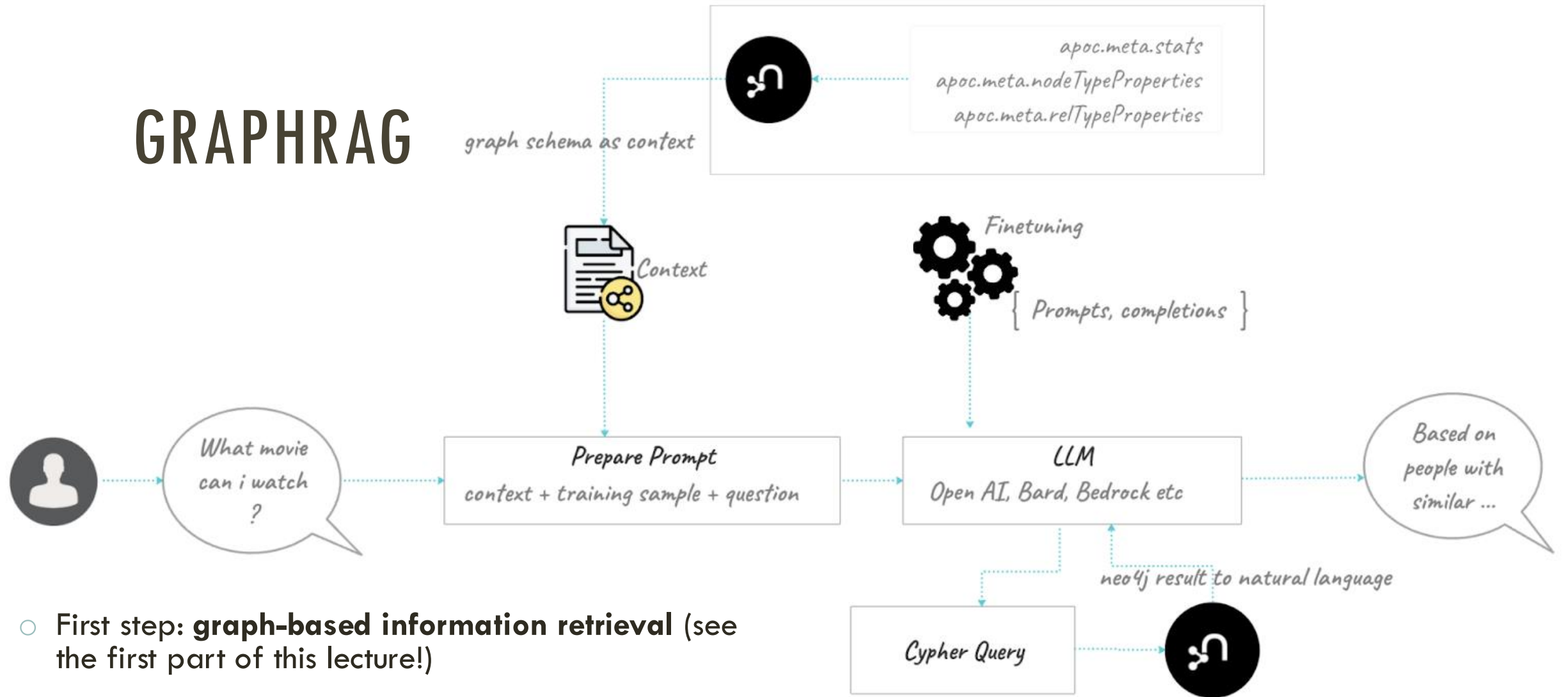
How can we fix it?

# IT'S NOT ABOUT HALLUCINATIONS...

We argue against the view that when ChatGPT and the like produce false claims they are lying or even hallucinating, and in favour of the position that the activity they are engaged in is bullshitting, in the Frankfurtian sense (Frankfurt, 2002, 2005). Because **these programs cannot themselves be concerned with truth**, and because they are designed to produce text that looks truth-apt **without any actual concern for truth**, it seems appropriate to call their outputs **bullshit.**

Source: M.T. Hicks, J. Humphries & J. Slater (2024), ChatGPT is bullshit, Ethics Inf Technol, 26, pp. 38.

# RETRIEVAL-AUGMENTED GENERATION (RAG)



- First step: documents retrieval (based on **vector databases** and **embeddings**)

- Second step: use LLM to generate output for user

- Easy to implement, but **lacks a comprehensive understanding of data**, relying primarily on similarity scores

Sources: (1) M. Gupta (2024), GraphRAG vs RAG: Which is Better?
(2) Z. Blumenfeld & E. Htet (2024), What Is Retrieval-Augmented Generation (RAG)?

# GRAPHRAG



- First step: **graph-based information retrieval** (see the first part of this lecture!)

- Second step: use LLM to generate output for user

- More complicated, but **offers enhanced data understanding by capturing the context** (associated information and related entities)

Sources: (1) M. Gupta (2024), GraphRAG vs RAG: Which is Better?
(2) M. Hunger (2024), Get Started With GraphRAG: Neo4j's Ecosystem Tools

# OUTLINE

1. Knowledge graphs

2. Towards automated KG management

3. Semantic search and recommendations

4. Knowledge graph embeddings

5. Knowledge graph completion

6. ChatGPT is a bullshit. How can we fix it?

# THANK YOU FOR YOUR ATTENTION!

GEIST Research Group: https://geist.re/

Krzysztof Kutt: https://krzysztof.kutt.pl/

KEEP CALM AND ASK QUESTIONS!

keep-calm.net