

Aprendizado de Máquina aplicado a Grafos de Conhecimento

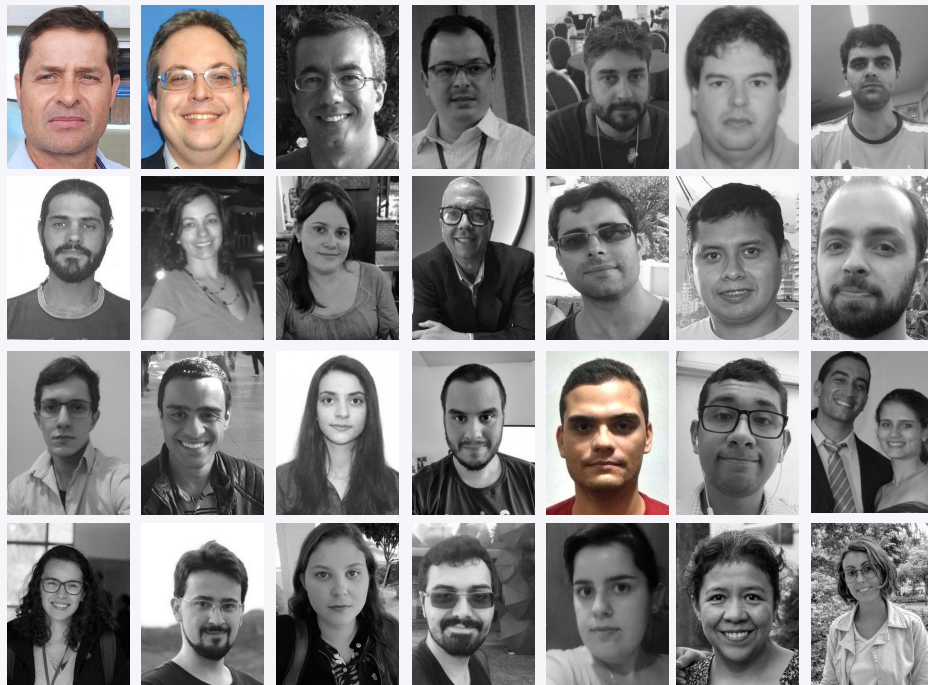
Daniel N. R. da Silva, Artur Ziviani e Fabio Porto
dramos, ziviani, fporto@Incc.br

Machine Learning for Knowledge Graphs

<http://dexl.lncc.br/>

- Big Data Management
- Complex Networks
- Computational Reproducibility
- Knowledge Bases
- Machine Learning
- Scientific Workflows

**LNCC: master and doctoral degrees on
Computational Modeling (a.k.a. Scientific
Computing) - CAPES 6.**



Acknowledgements:



Helpul Links:

<https://www.lncc.br/~ziviani/papers/Texto-MC1-SBBD2019.pdf>

<https://github.com/dnasc/knowledge-graphs>

Outline

Introduction and Motivation	04 to 20 (25 min)
Data Models and Systems	21 to 25 (10 min)
KG Tasks Intro	26 to 29 (10 min)
Break	15 min
KG Construction	30 to 38 (20 min)
KG Completion	39 to 65 (30 min)
Final Remarks	66 to 69 (10 min)

The third (current) rise of AI

Board Games

CRM

Drug Discovery

Image
Captioning

Medical Image
Analysis

Medical
Informatics

NLP

Online
Advertisement

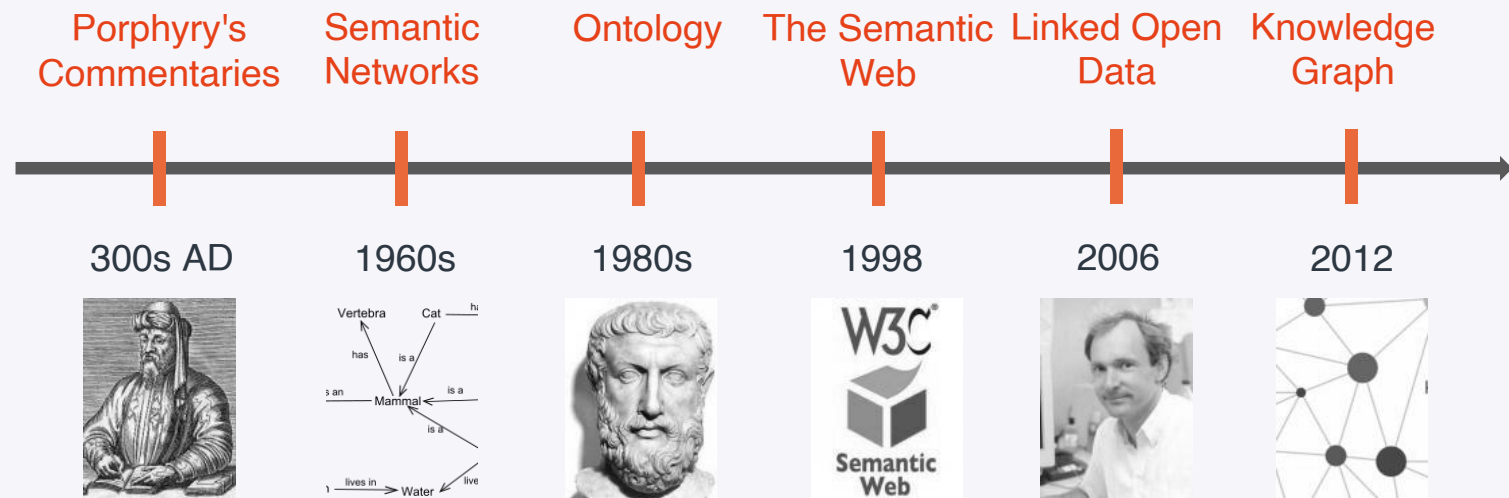
Recommendation
Systems

Robotics

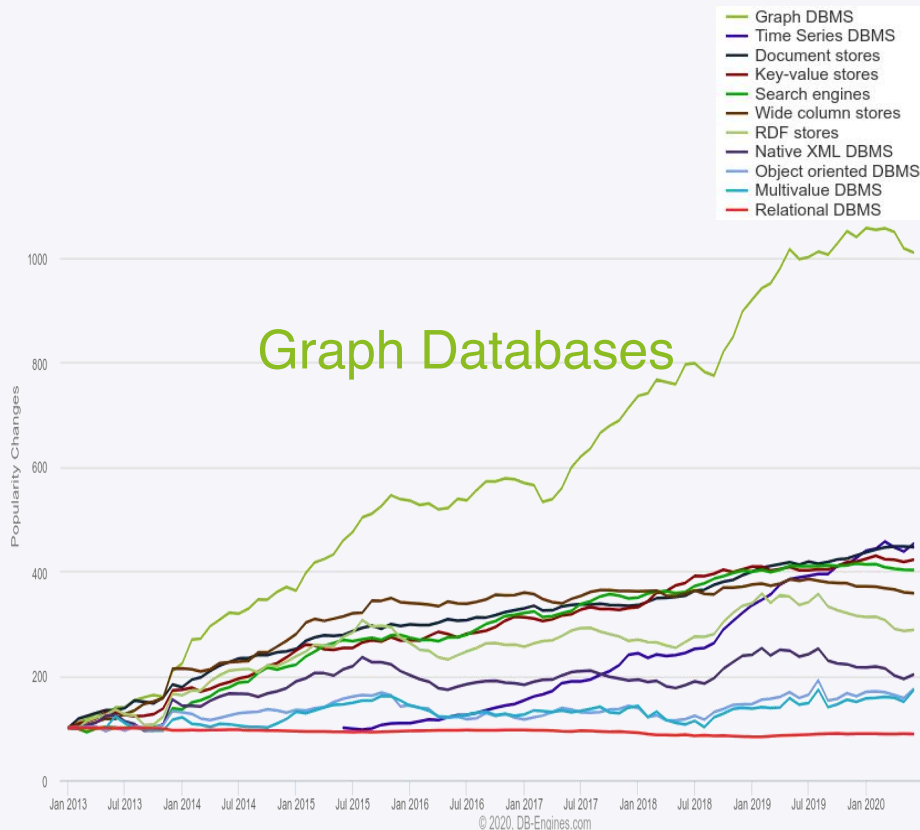
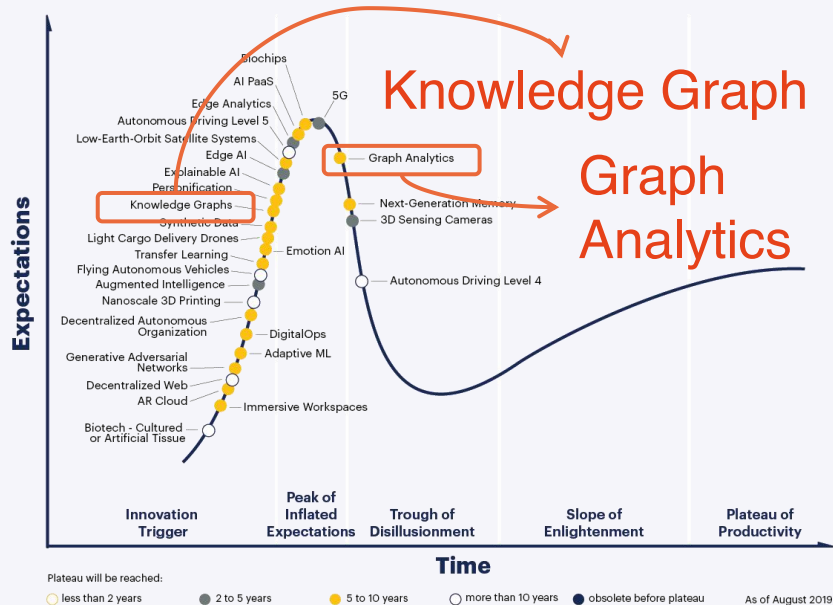
Self-Driving
Vehicles

Web Search

Knowledge Representation and Reasoning



Gartner Hype Cycle for Emerging Technologies, 2019



Knowledge Graph: What

No single formal definition, but it

Keeps real world entities.

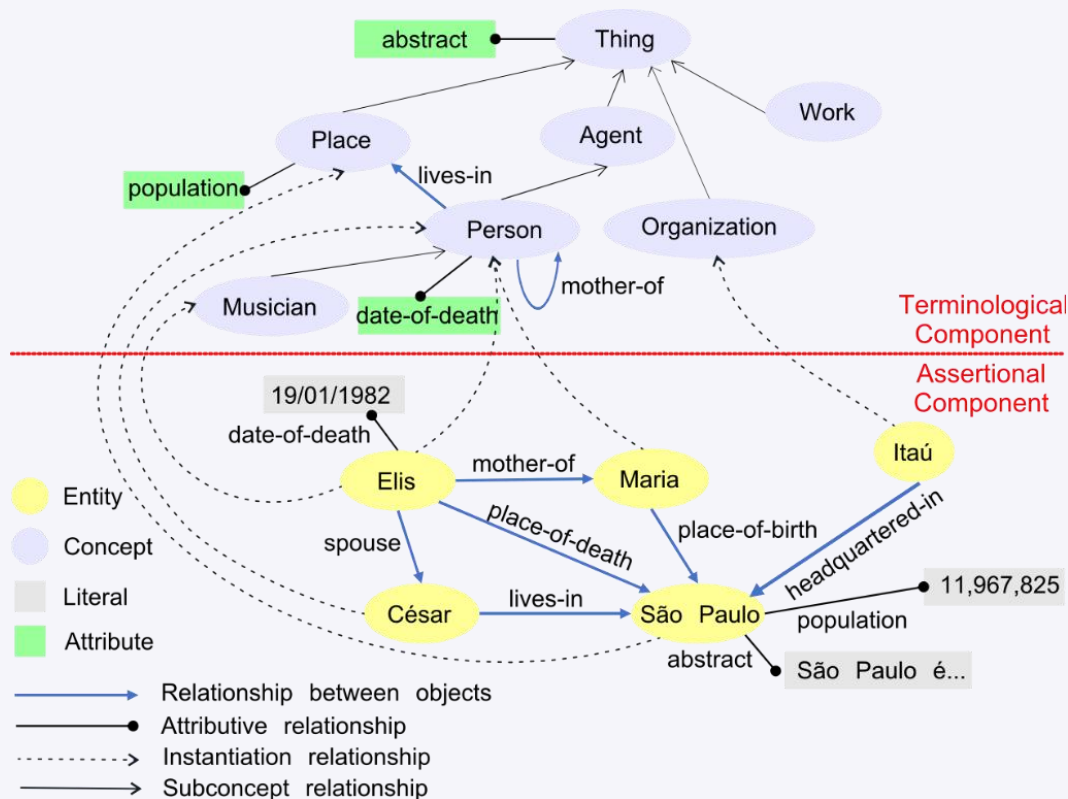
Provides relationships between them.

} Networked representation
for knowledge

May contain constraints and rules (ontology or schema).

Enables machine reasoning to infer unobserved knowledge.

Knowledge Graph: What



Knowledge Graph: Keywords and Areas

Keywords:

Graph database & triplestores, heterogeneous information networks, ontology, semantic networks, knowledge bases, knowledge based system.

Areas:

Artificial Intelligence: Knowledge Representation & Reasoning, Machine Learning, Natural Language Processing & Understanding.

Data Management: Data Integration, Data Modeling, and Information Retrieval.

Knowledge Graph: Examples



NELL



KBpedia



WIKIDATA



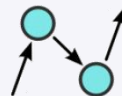
select: knowledge



GENE ONTOLOGY
Unifying Biology



BabelNet



SciGraph

Microsoft
Academic
Knowledge
Graph

General-purpose

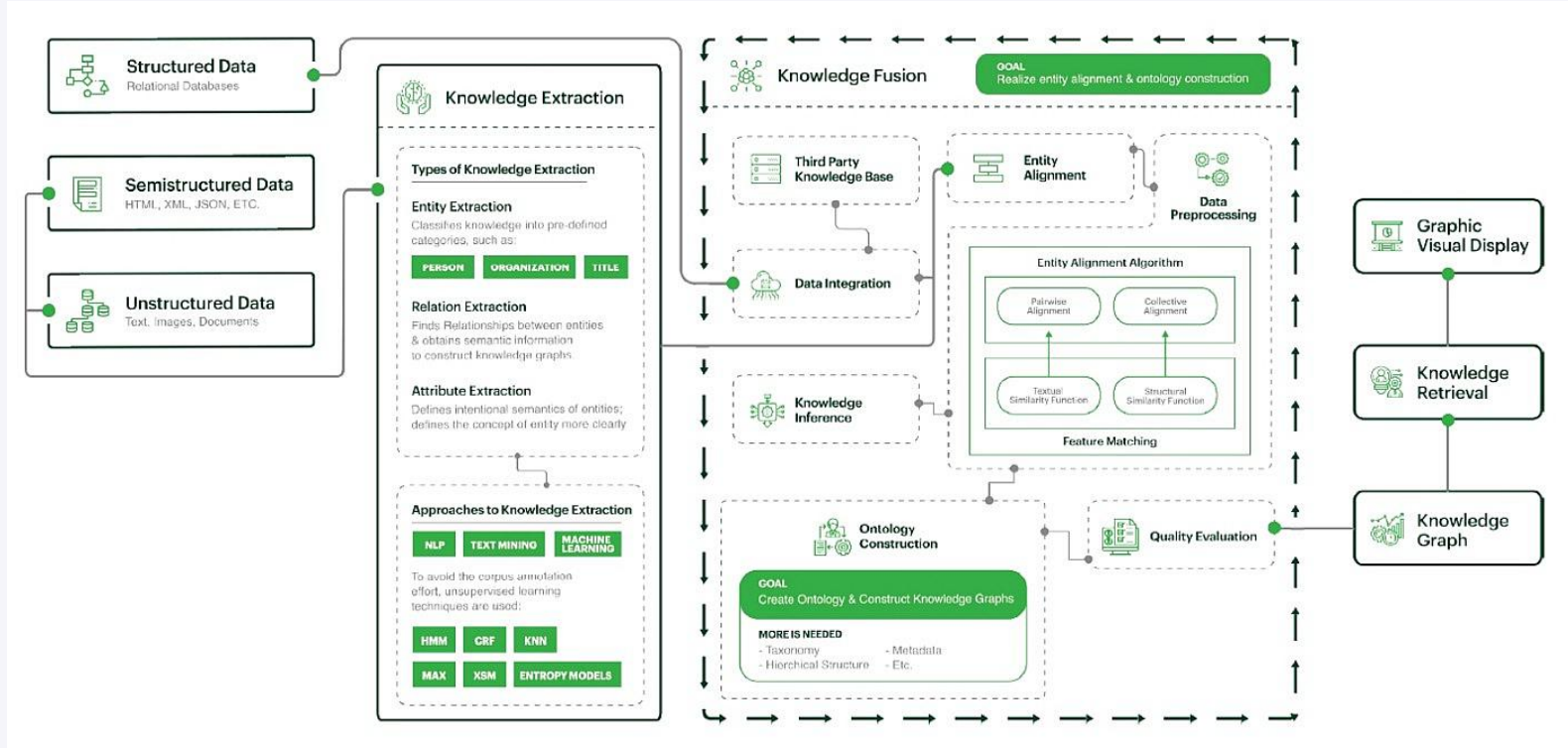
Product Graph
& E-commerce

Biomedical

Common-sense
& NLP

Academic

Knowledge Graph Cycle



Knowledge Graphs: Why

Conversational Agents

Data integration

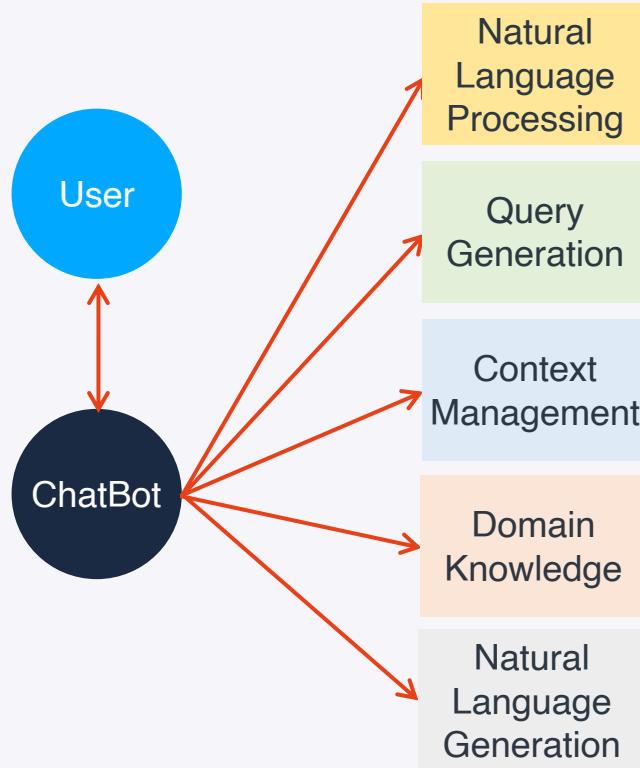
Fact checking and fake news detection

Question Answering

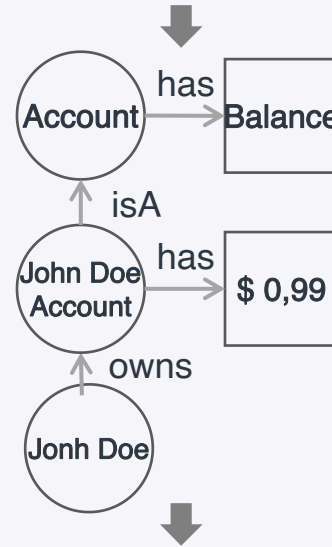
Recommender Systems

Search Engines

Conversational Agents



Jonh Doe: @bot, How much money do I have in my account?



Bot: @johndoe, You have \$0,99 in your bank account.

Data Integration: Covid19

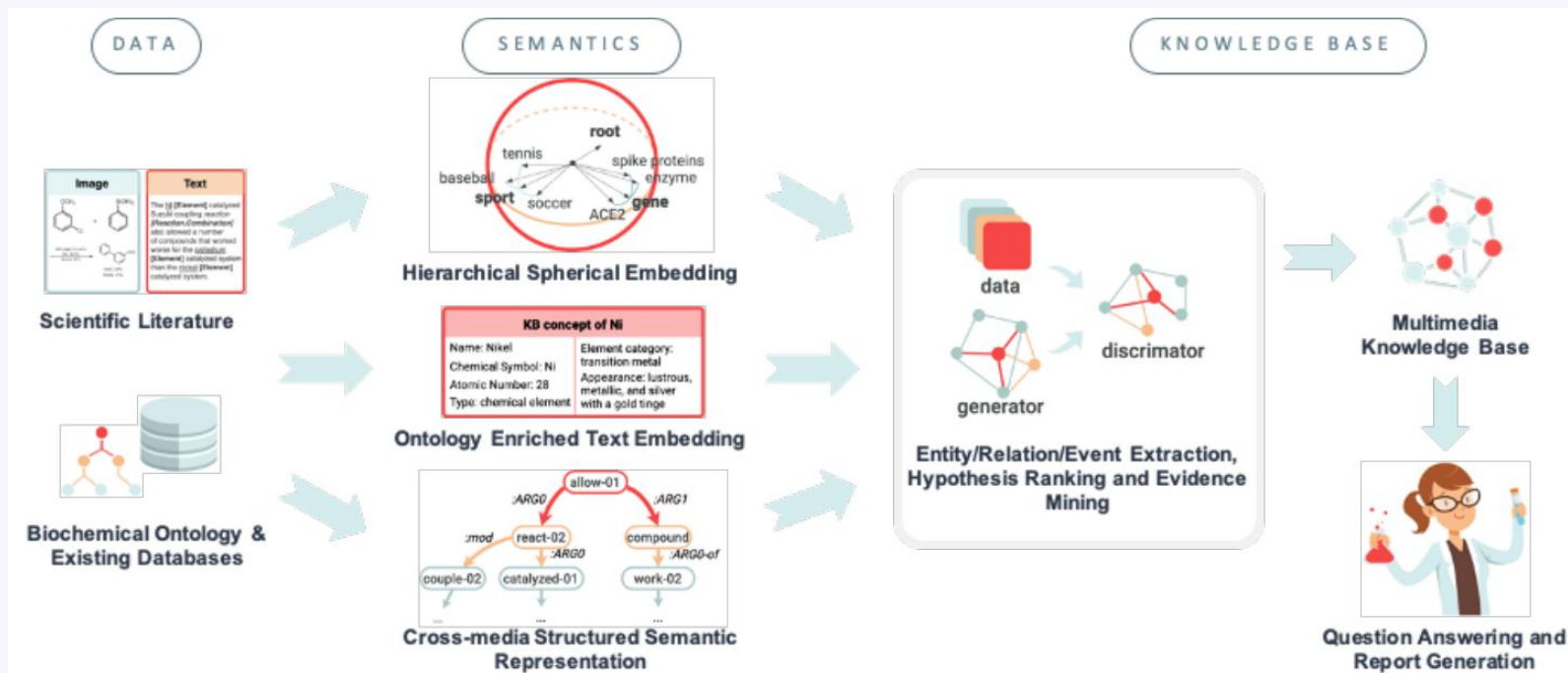
Data siloed: different databases, identifiers, formats, and licenses.

April 28th up to June 13th: 120K+ papers published on Pubmed related to Covid: **Quantity and Quality Challenges.**

Drug repurposing (a.k.a. repositioning, reprofiling, or re-tasking): investigating existing drugs for new therapeutic purposes:

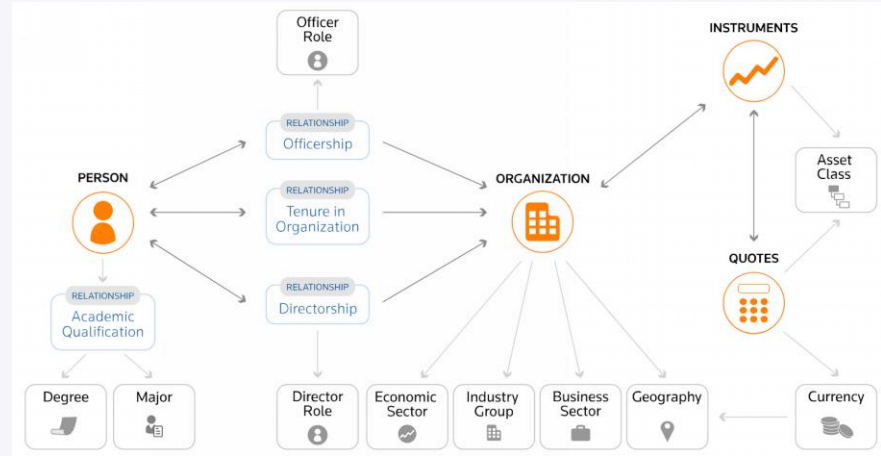
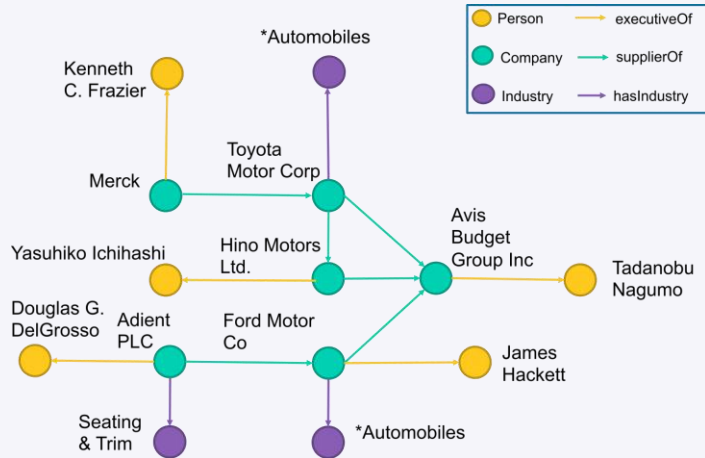
Minoxidil (Hypertension => Hair Loss), Aspirin (Analgesia => Colorectal Cancer).

Data Integration: Covid19

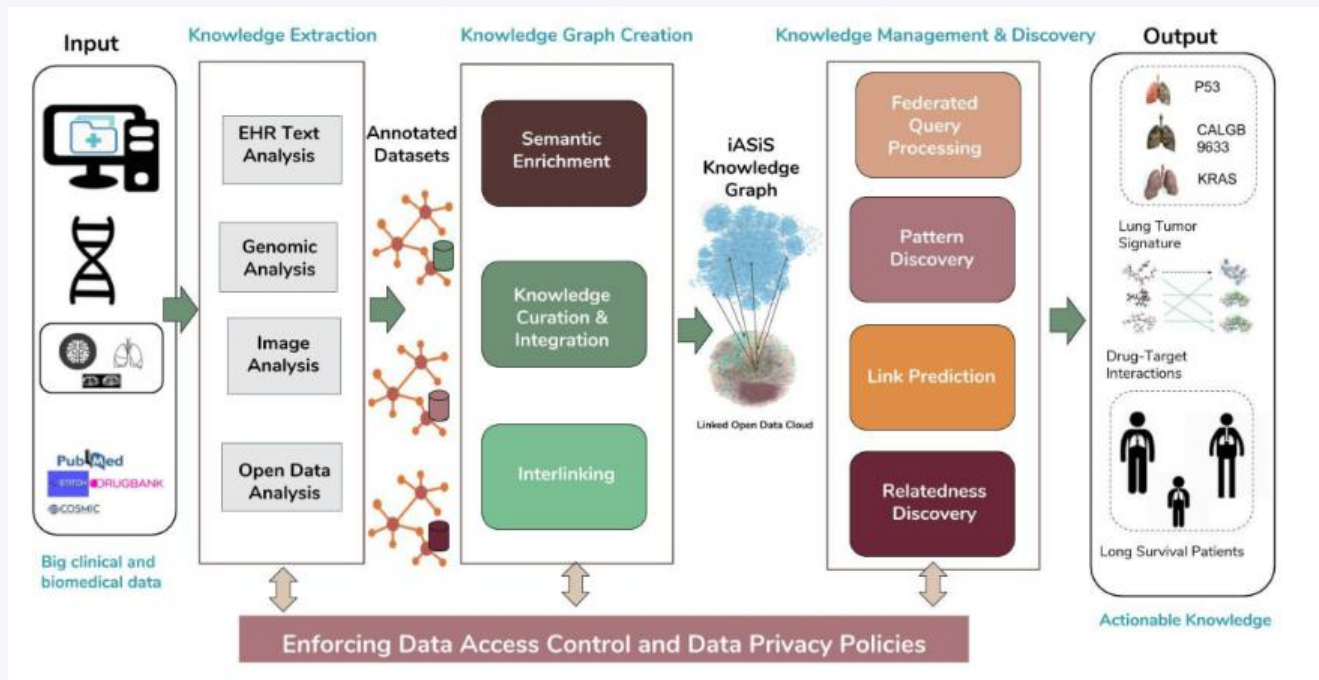


Data Integration: Financial Markets

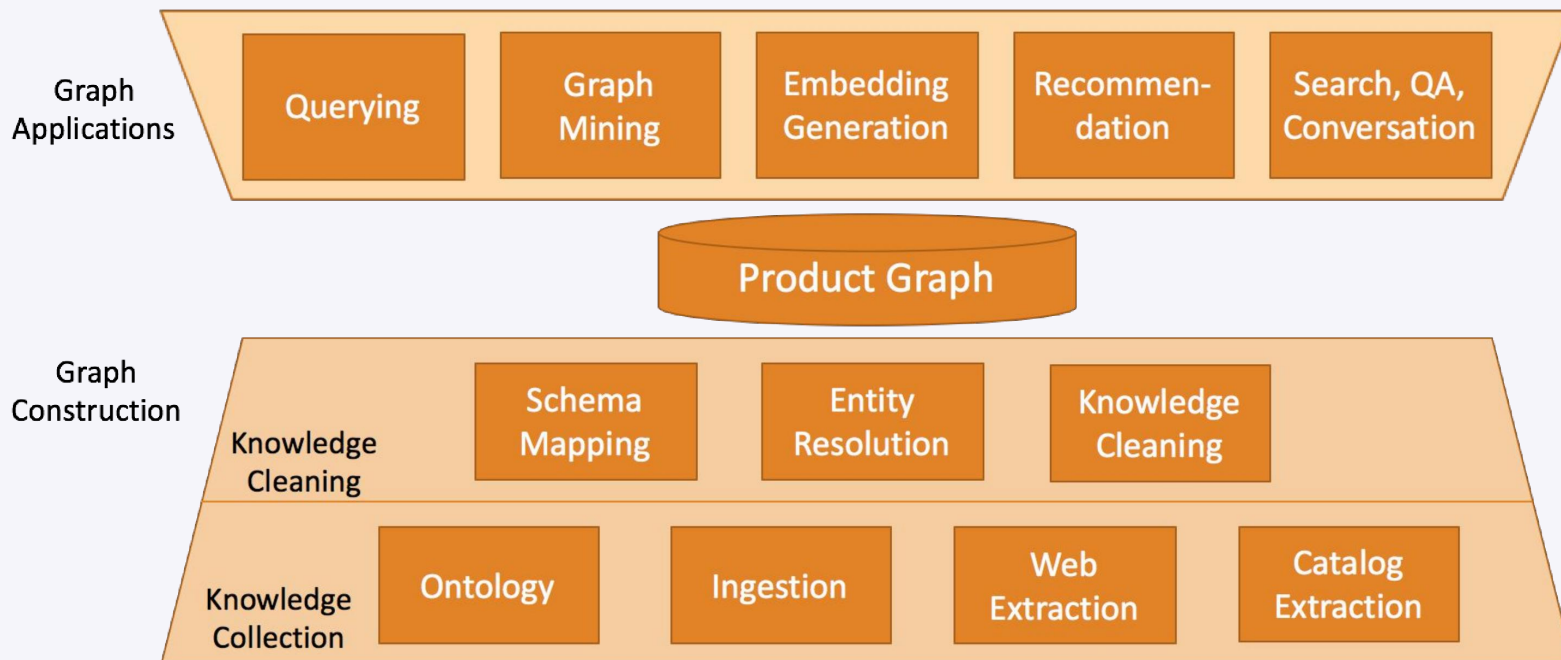
Bloomberg and Thomson Reuters (Refinitiv)



Data Integration: Personalized Medicine



Data Integration: Product Graph



Dong, L. **Challenges and Innovations in Building a Product Knowledge Graph**. KDD 2018.

Dong, L and Rekatsinas, T. **Data Integration and Machine Learning: A Natural Synergy**. Tutorial presented on SIGMOD 2018, VLDB 2018, and KDD 2019.

Fact Checking and Fake News Detection

Barack Obama secretly practices Islam?



Barack Obama



Columbia University



Association of American Universities



Canada



Stephen Harper



Calagary



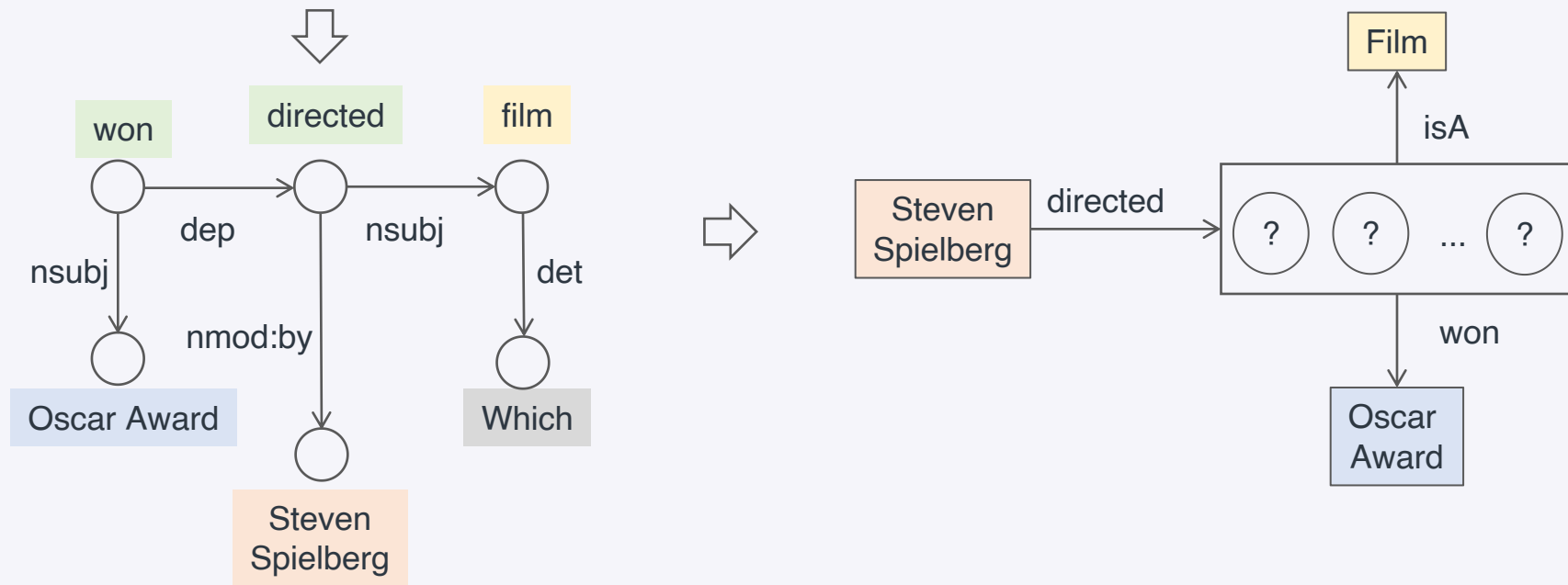
Naheed Neshi



Islam

Question Answering

Q: Which films directed by Steven Spielberg won the Oscar Award?

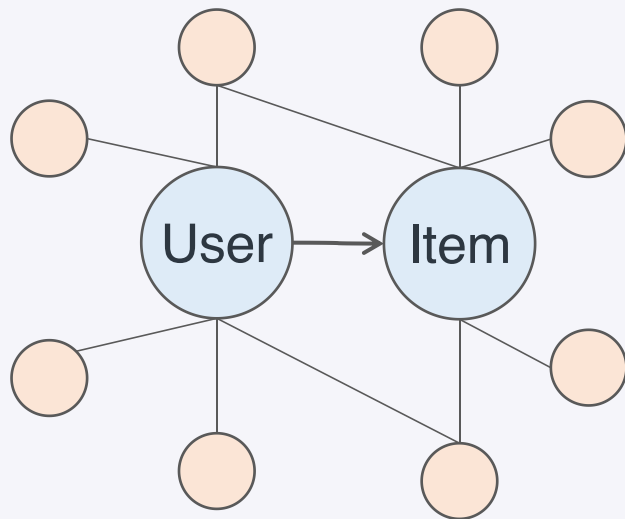


Recommender Systems



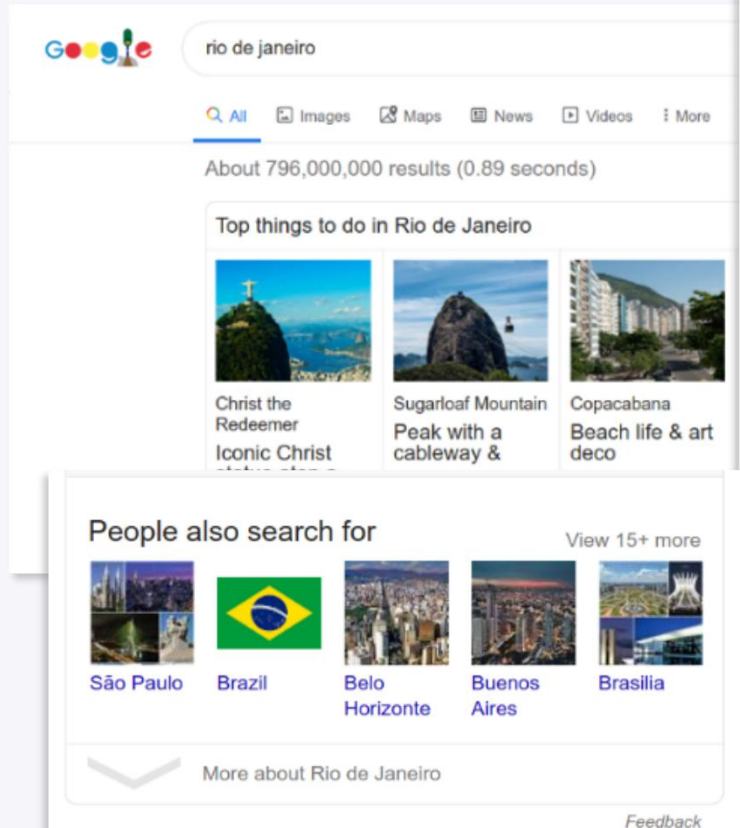
Bipartite graph: users and items (e.g., movies, music, shopping items).

Soften Cold Start Problem



Knowledge Graph

Search Engines



Google search results for "rio de janeiro". The search bar shows the query "rio de janeiro" and the Google logo. Below the search bar, there are tabs for "All", "Images", "Maps", "News", "Videos", and "More". The results show "About 796,000,000 results (0.89 seconds)". A section titled "Top things to do in Rio de Janeiro" features three images: Christ the Redeemer, Sugarloaf Mountain, and Copacabana. Below these images are captions: "Christ the Redeemer Iconic Christ", "Sugarloaf Mountain Peak with a cableway &", and "Copacabana Beach life & art deco". A section titled "People also search for" shows five images: São Paulo, Brazil, Belo Horizonte, Buenos Aires, and Brasília. Below these images are captions: "São Paulo", "Brazil", "Belo Horizonte", "Buenos Aires", and "Brasília". At the bottom, there is a "More about Rio de Janeiro" link and a "Feedback" link.

Google

rio de janeiro

Q All Images Maps News Videos More

About 796,000,000 results (0.89 seconds)

Top things to do in Rio de Janeiro

Christ the Redeemer
Iconic Christ

Sugarloaf Mountain
Peak with a cableway &

Copacabana
Beach life & art
deco

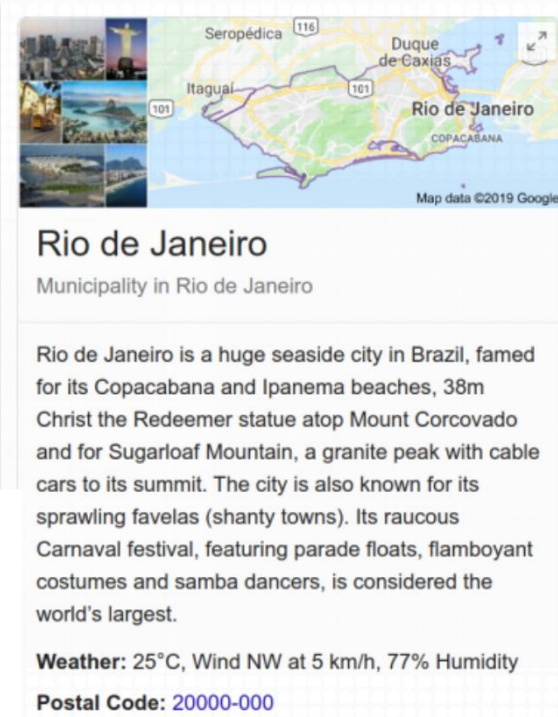
People also search for

São Paulo Brazil Belo Horizonte Buenos Aires Brasília

View 15+ more

More about Rio de Janeiro

Feedback



Information card for Rio de Janeiro. It features a map of the city and surrounding areas, including Seropédica, Duque de Caxias, Itaguaí, and Copacabana. The map is labeled "Rio de Janeiro" and "COPACABANA". Below the map, the title "Rio de Janeiro" is followed by the subtitle "Municipality in Rio de Janeiro". A paragraph describes the city: "Rio de Janeiro is a huge seaside city in Brazil, famed for its Copacabana and Ipanema beaches, 38m Christ the Redeemer statue atop Mount Corcovado and for Sugarloaf Mountain, a granite peak with cable cars to its summit. The city is also known for its sprawling favelas (shanty towns). Its raucous Carnival festival, featuring parade floats, flamboyant costumes and samba dancers, is considered the world's largest." Below the paragraph, the weather is listed as "Weather: 25°C, Wind NW at 5 km/h, 77% Humidity" and the postal code is "Postal Code: 20000-000".

Seropédica 116 Duque de Caxias 101 Itaguaí 101 Rio de Janeiro COPACABANA

Map data ©2019 Google

Rio de Janeiro

Municipality in Rio de Janeiro

Rio de Janeiro is a huge seaside city in Brazil, famed for its Copacabana and Ipanema beaches, 38m Christ the Redeemer statue atop Mount Corcovado and for Sugarloaf Mountain, a granite peak with cable cars to its summit. The city is also known for its sprawling favelas (shanty towns). Its raucous Carnival festival, featuring parade floats, flamboyant costumes and samba dancers, is considered the world's largest.

Weather: 25°C, Wind NW at 5 km/h, 77% Humidity

Postal Code: 20000-000

Labeled Property Graph

Graph Databases:

Ex.: Amazon Neptune, DGraph, Neo4J, JanusGraph, Memgraph, and TigerGraph.

Query languages:

Cypher, GCore, GS, Gremlin, and PGQL.

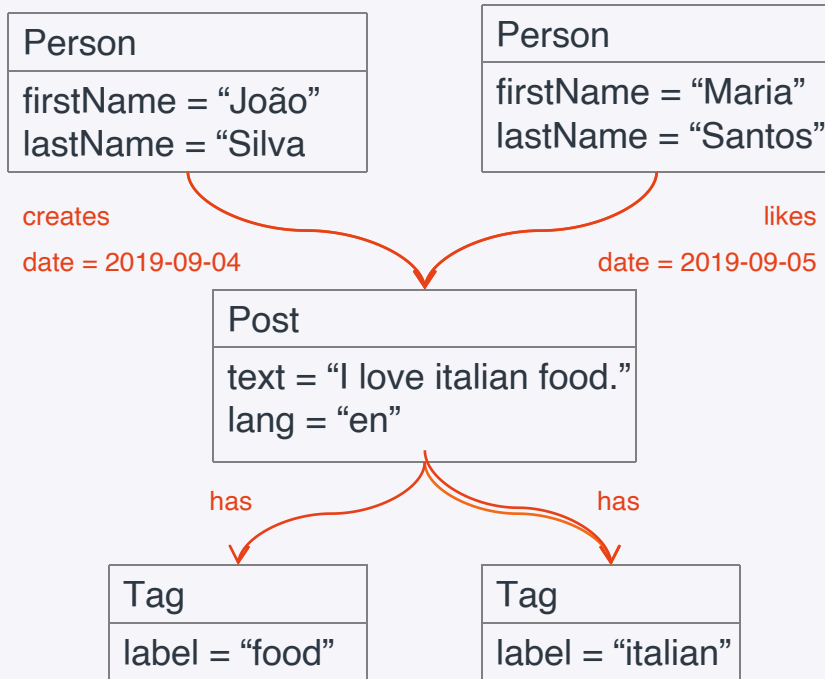
<https://aws.amazon.com/neptune/>

<https://neo4j.com/>

<https://janusgraph.org/>

<https://www.tigergraph.com/>

Labeled Property Graph



Get the tags associated with Maria's post preference.

```
match (:Person {firstName: "Maria"})  
->[:likes]->(:Post)->[:has]->(t: Tag)  
RETURN t.label as TagLabel
```


Resource Description Framework (RDF)

Nodes: Resources, literals or blank nodes

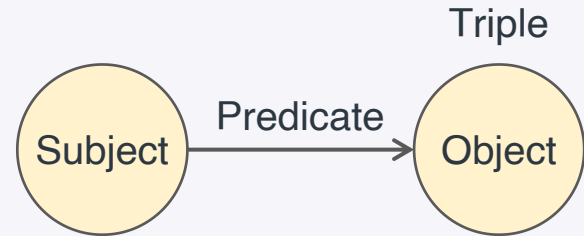
Edges: Predicates (resource)

Literal:

- Can be interpreted as datatypes

- Encoded as strings

- Represent data values



Resource Description Framework (RDF)

IRI/URI (Internationalized/Uniform Resource Identifier)

Serialization: JSON-LD, N-Triples, RDF/XML, and Turtle

Query Languages: SPARQL.

@prefix dbr: <http://dbpedia.org/resource>.

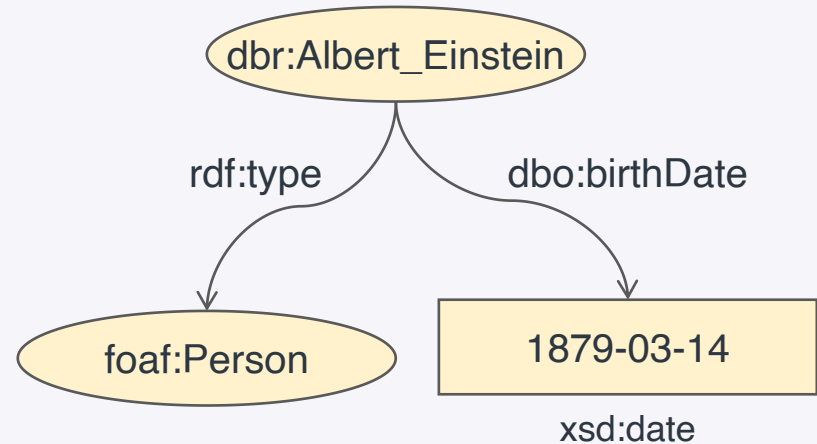
@prefix dbo: <http://dbpedia.org/ontology>.

@prefix foaf: <http://xmlns.com/foaf/0.1/> .

dbr:Albert_Einstein

rdf:type foaf:Person;

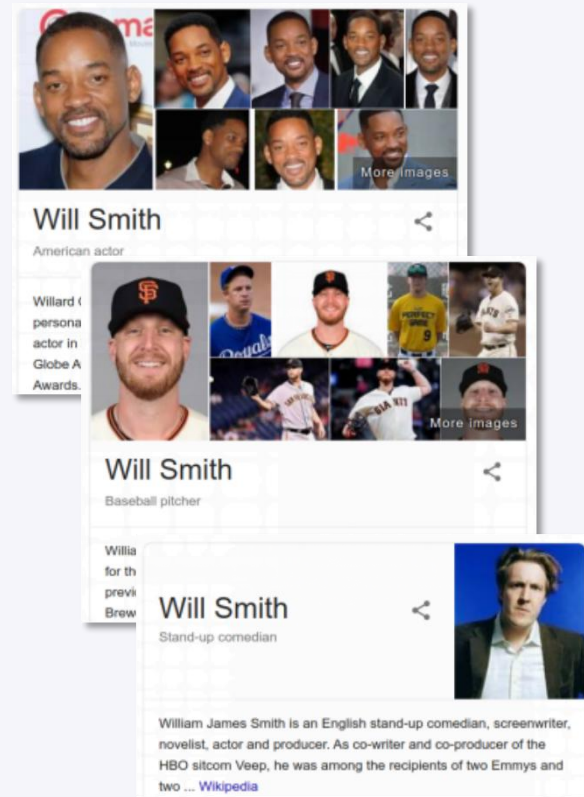
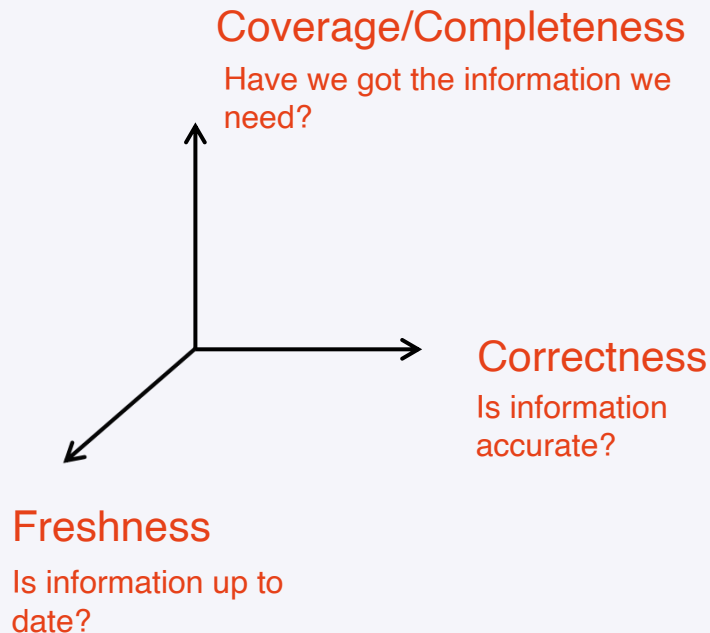
dbo:birthDate "1879-03-14"^^xsd:date.



Systems for Knowledge Bases/Graphs

AllegroGraph	https://allegrograph.com/
Atomgraph	https://atomgraph.com/
Amazon Neptune	https://aws.amazon.com/neptune/
Diffbot	https://www.diffbot.com/
Grakn	https://grakn.ai/
MarkLogic	https://www.marklogic.com/
Microsoft Cosmos	https://azure.microsoft.com/en-us/services/cosmos-db/
Ontotext GraphDB	https://www.ontotext.com/
Stardog	https://www.stardog.com/

Challenges



Tasks

Knowledge Base/Graph Construction: Extract and populate a KB with data extracted from a set of documents.

Knowledge Graph Completion (Reasoning): Infer (and discover) non-observed facts over relevant entities.

Will not be discussed:

- Data / Knowledge Fusion.

- Knowledge Graph Correction.

- KG Ontology alignment (matching)/ merging.

Named Entity Recognition (NER)

Pedro II of Brazil was the second and last monarch of the **Empire of Brazil**. He was born in **Rio de Janeiro** to **Emperor Pedro I of Brazil** and **Empress Maria Leopoldina**, who were married at that time.



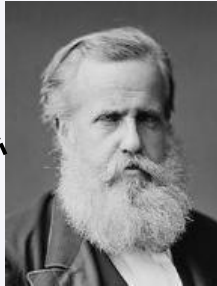
Person

State

City

Entity Linking

Pedro II of Brazil was the second and last monarch of the **Empire of Brazil**. He was born in **Rio de Janeiro** to **Emperor Pedro I of Brazil** and **Empress Maria Leopoldina**, who were married at that time.



Q156774



Q217230



Q939



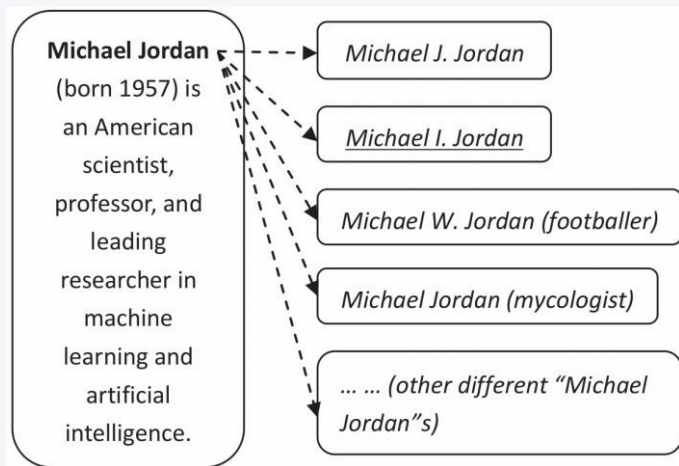
Q84239



Q8678

Entity Linking

- Disambiguation with(out) NER
- Candidate Entity Generation
- Candidate Entity Ranking
- Unlinkable Mention Prediction



The screenshot shows a Google search result for "Rachel Abrams", identified as an "American writer". A handwritten red arrow points from the word "Me" at the top to the name "Rachel Abrams". Another handwritten red arrow points from the text "could be me...?" to the same name. Below the name is a photo of Rachel Abrams. Further down, a list of biographical details is shown: "Born: January 2, 1951", "Died: June 7, 2013", "Spouse: Elliott Abrams (m. 1980–2013)", "Parents: Midge Decter", and "Children: Sarah Abrams, Jacob Abrams, Joseph Abrams". A handwritten red bracket on the left groups the "Born" and "Died" information, with the text "Not me" written next to it. A handwritten red arrow points from the text "Definitely not me" to the "Died: June 7, 2013" entry. At the bottom, a section titled "People also search for" displays five small portrait photos of other individuals.

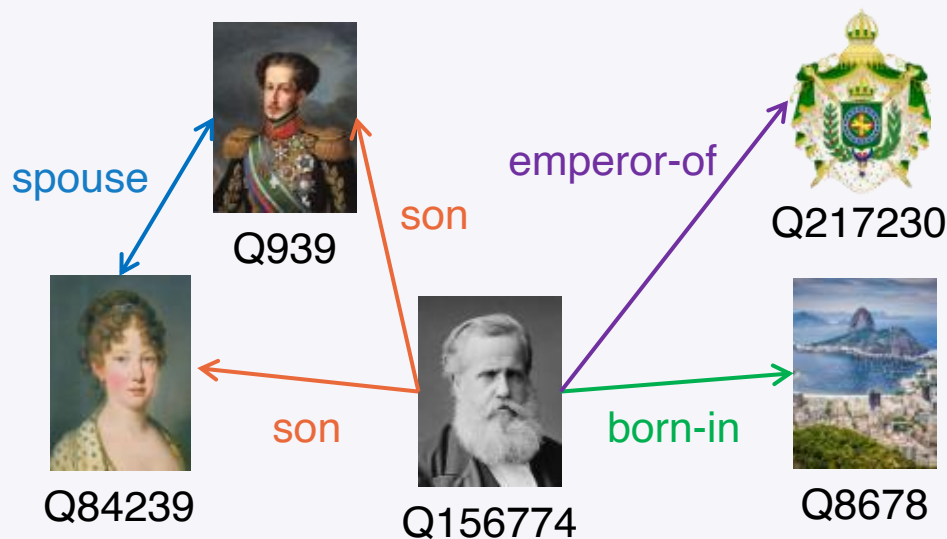
Relation Extraction

Entity Mention

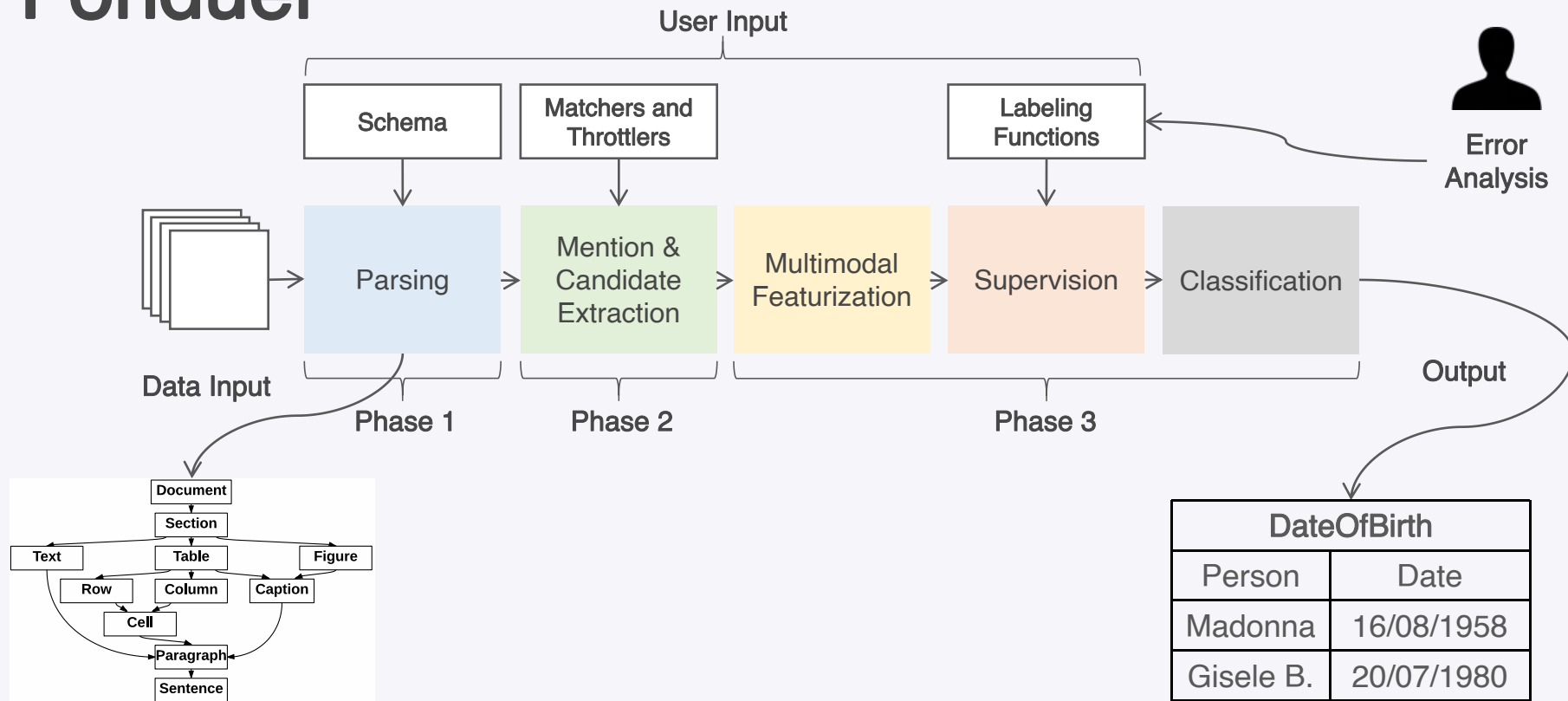
Pedro II of Brazil was the second and last monarch of the Empire of Brazil. He was born in Rio de Janeiro to Emperor Pedro I of Brazil and Empress Maria Leopoldina, who were married at that time.

Relation Candidate

(Pedro II of Brazil, born-in, Rio de Janeiro)



Fonduer



Fonduer

Mention and Candidate Extraction: Functions to return mentions to entities (matchers) and to decrease the number of relationship candidates (throttlers).

Multimodal Featurization: Associate textual, structural, visual and tabular features to relationship candidates.

Supervision and Classification

Labeling Function: Yields a label (-1, 0 or 1) for each candidate.

Data Programming: Estimate the true label for each candidate.

Multimodal BiLSTM: Estimate the true label for each candidate considering features.

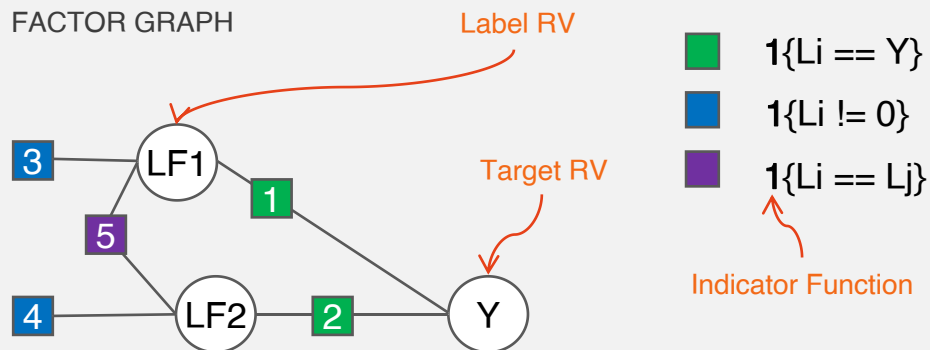
Data Programming

Estimate the ground truth based on labeling functions agreements and disagreements.



snorkel

FACTOR GRAPH



$$P_w(L, Y) = \frac{\exp(\sum_k w^T L)}{\int_w \exp(\sum_k w^T L)}$$

Ps.: k ranges on the candidate set.

<i>Spouse1</i>	<i>Spouse2</i>	<i>LF1</i>	<i>LF2</i>	<i>Y</i>
Tom Hanks	Rita Wilson	1	-1	?

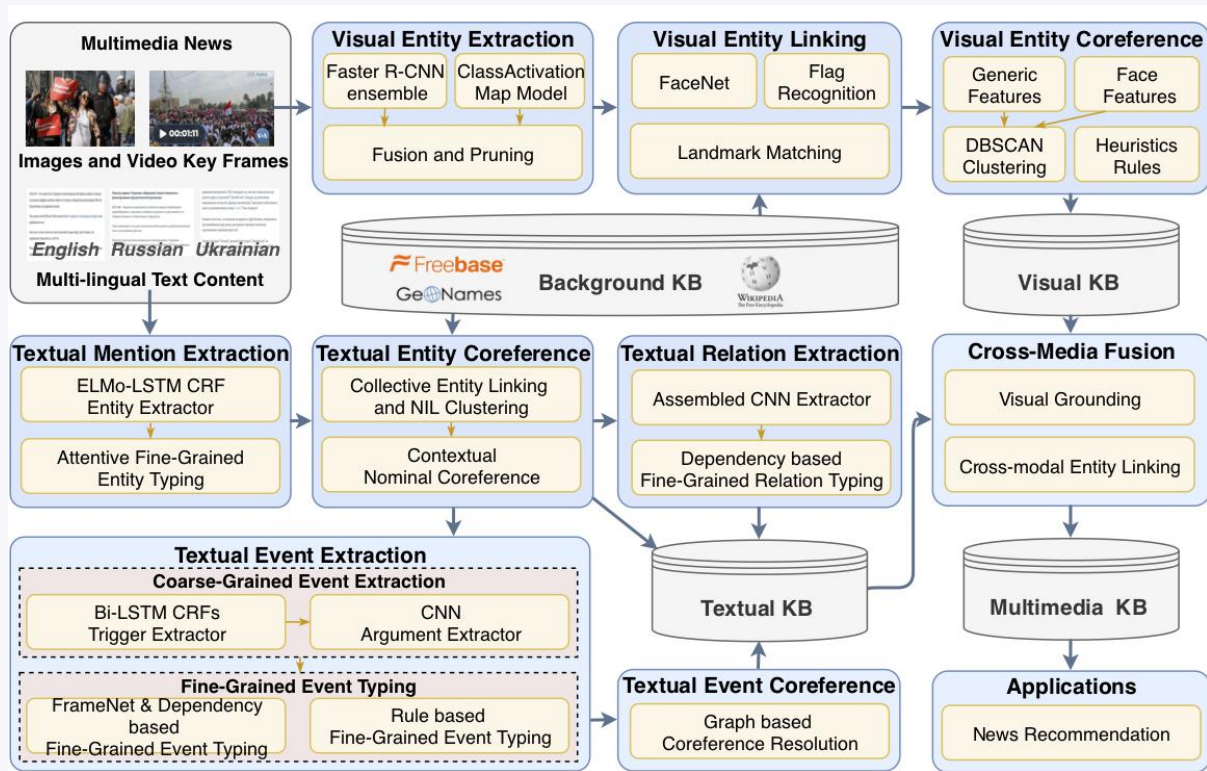
weight associated with the 1st factor

$$w = [w_1, w_2, w_3, w_4, w_5]$$

$$L = [?, ?, 1, 1, 0, 0]$$

label associated with the 2nd factor

GAIA



Home

Back

Query: Target: Януковича (Yanukovych)

Event Search

Number of Events: 2

Automated Summary:

Source Document

Translation from Ukrainian/Russian

Show Visual Knowledge Elements

Hide Visual Knowledge Elements

Столкновения 20 февраля стали одним из ключевых факторов, вынудивших Президента Украины Виктора Януковича пойти на подписание Соглашения об урегулировании политического кризиса на Украине, потеря доверия к самому Януковичу и к переформатированию парламентского большинства. 20 февраля постановление о запрете применения силы властью (Translation: The clashes on February 20 became one of the key factors forcing the President of Ukraine Viktor Yanukovich to sign the Agreement on the settlement of the political crisis in Ukraine, the loss of confidence in Yanukovich himself and the reformatting of the parliamentary majority, which issued a resolution on the evening of February 20 banning the use of force), что согласно всем имеющимся уликам те милиционеры и демонстранты, что стали жертвами снайперского огня, застрелены одними и теми же снайперами (Translation: that according to all available evidence, those policemen and demonstrators were shot by the same snipers)

Event Summary

Source Doc & Text Extraction Result

Visual Entity Linking

Visual Entity Extraction

Event Arguments

Knowledge Elements based Ranking Incorporating User Feedback

Similar Events

Dissimilar Events

Recommended Events

Date

Location

Attackers

Target

Instrument

Type of Attack

Source Doc Translation

HC000T6CP, 2011-01-19

Event

Time

Person

Organization

GeopoliticalEntity

Location

Facility

Vehicle

Weapon

Other

- The case of Kiev snipers
- Red Cross Volunteers of Ukraine provide first aid to a wounded man on Instutska the beginning of the eleventh hour on February 20, 2014
- Self-defenders carry out comrade on Instutskaya Street to the rear at the end of eleventh hour on February 20
- Mark 13 on a pierced bullet on Instutskaya Street, pasted by criminologists near the side opposite the Maidan
- The case of Kiev snipers question about the organizers and perpetrators of sniper Euromaidan participants and at the same time law enforcement officers in Kiev on 20, 2014, which killed 53 people (49 protesters and 4 law enforcement officers)

Instrument

Type of Attack

Date

Location

Attackers

Target

Instrument

Type of Attack

Event Type

Li, M. et al. GAIA: A Fine-grained Multimedia Knowledge Extraction System. ACL. 2020.
<http://blender.cs.illinois.edu/software/gaia-ie/>

38

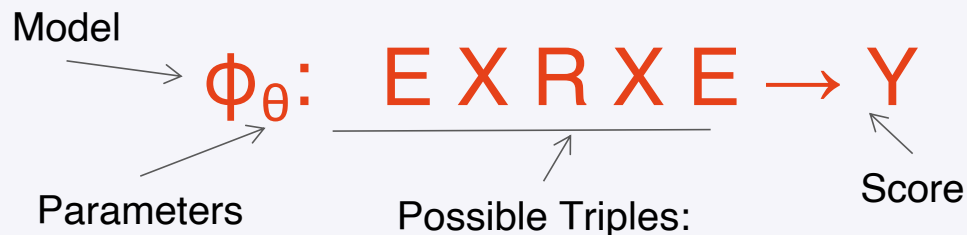
Knowledge Graph Completion

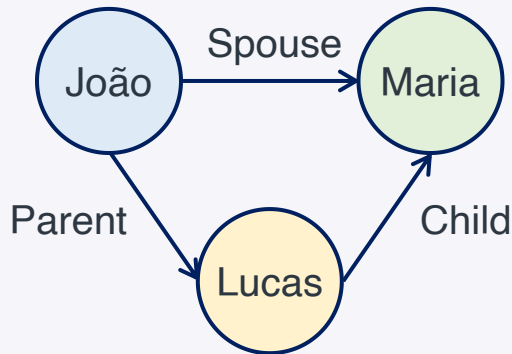
	Task	Query Example	Result Example
	Triple Classification	(Einstein, died-in, USA)?	(Yes, 90%)
Link Prediction	Tail Prediction	(Elvis Presley, starred-in, ?)	(1, Blue Hawaii, 3.23), (2, Change of Habit, 3.12), ...
	Head Prediction	(?, starred-in, Casablanca)	(1, Humphrey Bogart, 2.21), (2, Ingrid Bergman, 2.01), ...
	Relation Prediction	(Einstein, ?, Germany)	(1, born-in, 5.01), (2, died-in, 1.23),...
	Attribute Prediction	(Obama, nationality, ?)	(1, american, 2.21), (2, kenian, 1.02), ...
	Entity Classification	(Michael Jackson, isA, ?)	(1, singer, 6.20), (2, composer, 5.22),...

(ranking, answer, score)

Relational Machine Learning

Model Class	Triples are ...
Probabilistic Graphical Models	Interdependent.
Graphical Feature Model	Independent given observed features.
Latent Feature Models	Independent given latent features.





(Lucas, child, João)?

Possible Triples: **Not observed** + **Observed Triples**

(João, Spouse, João)
(João, Spouse, Maria)
(João, Spouse, Lucas)
(João, Parent, João)
(João, Parent, Maria)
(João, Parent, Lucas)
(João, Child, João)
(João, Child, Maria)
(João, Child, Lucas)

(Maria, Spouse, João)
(Maria, Spouse, Maria)
(Maria, Spouse, Lucas)
(Maria, Parent, João)
(Maria, Parent, Maria)
(Maria, Parent, Lucas)
(Maria, Child, João)
(Maria, Child, Maria)
(Maria, Child, Lucas)

(Lucas, Spouse, João)
(Lucas, Spouse, Maria)
(Lucas, Spouse, Lucas)
(Lucas, Parent, João)
(Lucas, Parent, Maria)
(Lucas, Parent, Lucas)
(Lucas, Child, João)
(Lucas, Child, Maria)
(Lucas, Child, Lucas)

Markov Logic Networks

Syntax: Weighted first-order formulas
Semantics: Templates for Markov Nets
Inference: Logical and Probabilistic
Learning: Statistical and Inductive
Logical Programming

Set of pairs (F_i, w_i) :

F_i : First-order logic formula;

w_i : A real number (weight).

n_i : Number of satisfied groundings of F_i in y (possible world).

Probabilistic Graphical Model:

One node (random variable) for each grounding atom.

Edges between nodes appearing at the same grounding formula.

Markov Logic Networks

Knowledge Base

Rules

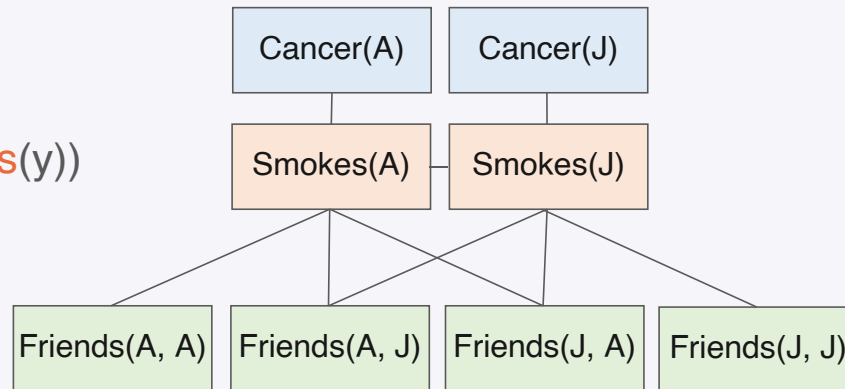
$\forall x \text{ Smokes}(x) \rightarrow \text{Cancer}(x)$

$\forall x \forall y \text{ Friends}(x,y) \rightarrow (\text{Smokes}(x) \leftrightarrow \text{Smokes}(y))$

Facts

Ana(A)

João(J)



$$P(\mathbf{Y} = \mathbf{y}) = \frac{1}{Z} \exp \left(\sum_i \underbrace{w_i}_{\text{Rule (feature) weights}} \underbrace{n_i(\mathbf{y})}_{\text{Number of times the rule is satisfied in the world } \mathbf{y}} \right) \quad Z = \sum_{\mathbf{y}} \exp \left(\sum_i w_i n_i(\mathbf{y}) \right)$$

Rule (feature) weights

Number of times the rule
is satisfied in the world \mathbf{y} .

Path Ranking Algorithm

Random walks of bounded length.

Feature Extraction

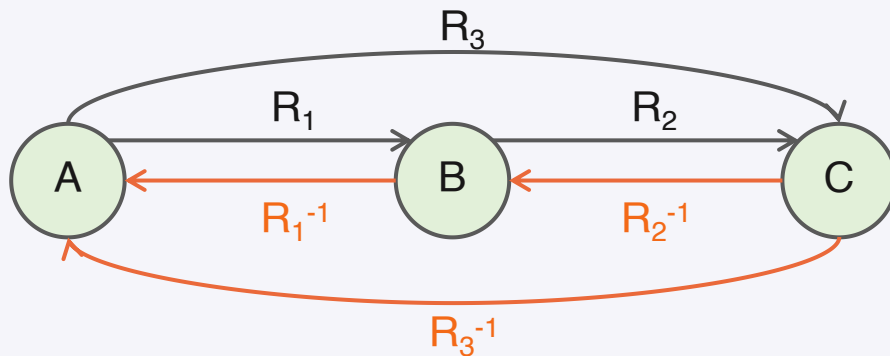
Build feature vectors for triples.

Vectors are based on relation paths $P_j(R_1, \dots, R_n)$.

Training:

Train an off-the-shelf machine learning model.

Relation Paths and Feature Extraction



$P_1(R_1, R_2)$
 $P_2(R_2^{-1}, R_1^{-1})$
 $P_3(R_3, R_2^{-1})$
 $P_4(R_3^{-1}, R_1)$
 $P_5(R_2, R_3^{-1})$
 $P_6(R_1^{-1}, R_3)$

Feature Vector of (A, R3, C):

Probabilities of reaching C from A by following given relation paths: [1/4, 0, 0, 0, 0, 0]

Other Approaches

Graphical Feature Models

Rule Mining (aka Association Rule Learning): RuleN, Rudik, AnyBURL, AMIE-3.

Probabilistic Graphical Models

Probabilistic Soft Logic (Hinge Markov Random Fields).

Meilicke, C. et al. **Fine-Grained Evaluation of Rule and Embedding-Based Systems for Knowledge Graph Completion**. ISWC 2018.

Ortona, S. et al. **RuDiK: Rule Discovery in Knowledge Bases**. PVLDB 2018.

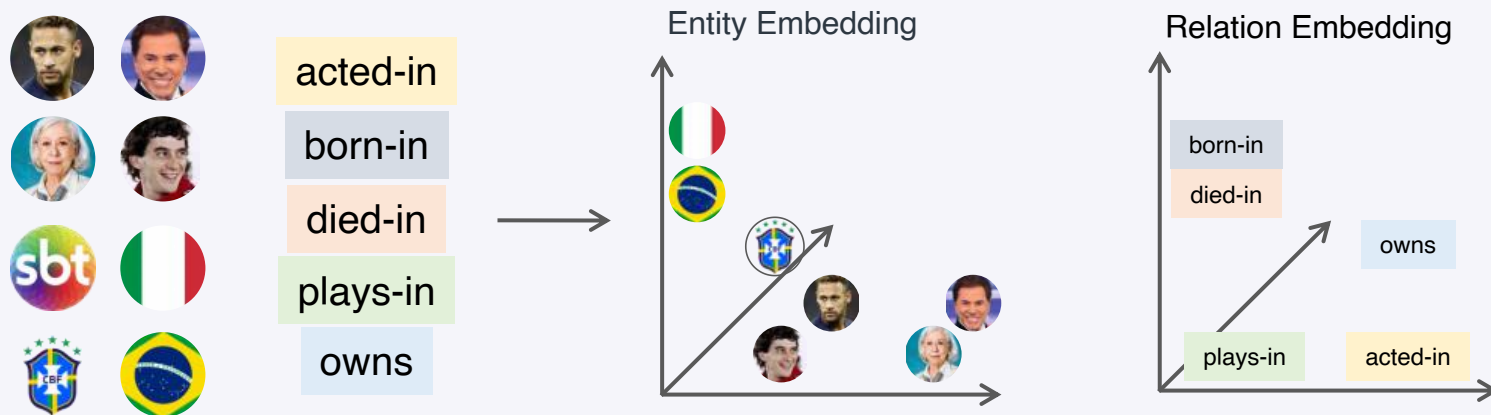
Meilicke, C. et al. **Anytime Bottom-Up Rule Learning for Knowledge Graph Completion**. IJCAI 2019

Lajus, J. et al. **Fast and Exact Rule Mining with AMIE 3**. ESWC 2020.

Bach, SH. et al. **Hinge-Loss Markov Random Fields and Probabilistic Soft Logic**. JMLR 2017.

Knowledge Graph Embedding (KGE)

Embed components of a Knowledge Graph including entities and relations into continuous vector spaces, so as to simplify the manipulation while preserving the inherent structure of the KG.



Anatomy of a KGE Model

Knowledge Graph (KG)

Triple corruption strategy
(e.g., negative sampling)

Cost (and loss) function

Scoring function for a triple

Optimization algorithm

Input: Observed triplets T , number of training epochs e , batch size b , number of corruptions c , model μ_{Θ} , and cost function \mathcal{J} .

INITIALIZE PARAMETERS Θ .

for $i = 1, \dots, e$ **do**

$T_i \leftarrow T$

while $|T_i| \neq 0$ **do**

$T^+ \leftarrow \text{SAMPLE}(T_i; b)$.

$B \leftarrow \bigcup_{t \in T^+} \langle t, \text{CORRUPT}(t; c) \rangle$.

UPDATE MODEL PARAMETERS ACCORDING TO $\nabla_{\Theta} \mathcal{J}(B)$.

$T_i \leftarrow T_i \setminus T^+$

end while

end for

Triplet Corruption

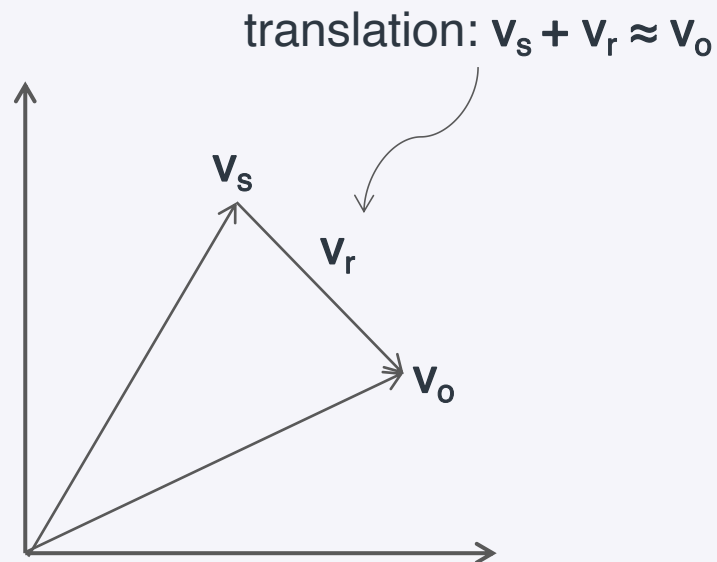
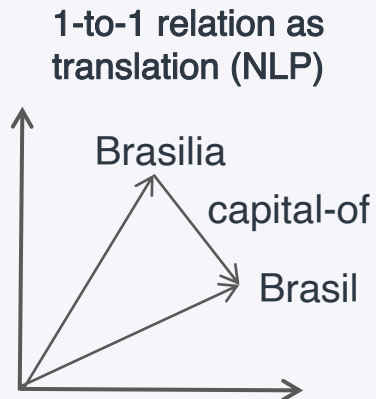
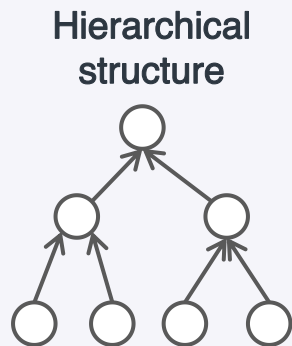
KGs usually don't include negative assertions.

- Manually label a set of negative examples.
- Use rules and constraints expressed in the KG.
- **Negative Sampling:** Sample triples from the unobserved set of possible triples.

$$(s, r, o) \Rightarrow (s, r, ?), (?, r, o)$$

Replace **?** with a random entity such that the resulting triple isn't in the KG.

TransE



TransE

Constraints on entity embedding: Prevent learning trivial representations.

Limitations on dealing with 1-N, N-1, N-M relations.

$$\phi_{(s,r,o)}^{\text{TransE}} := \|\mathbf{v}_s + \mathbf{v}_r - \mathbf{v}_o\|$$

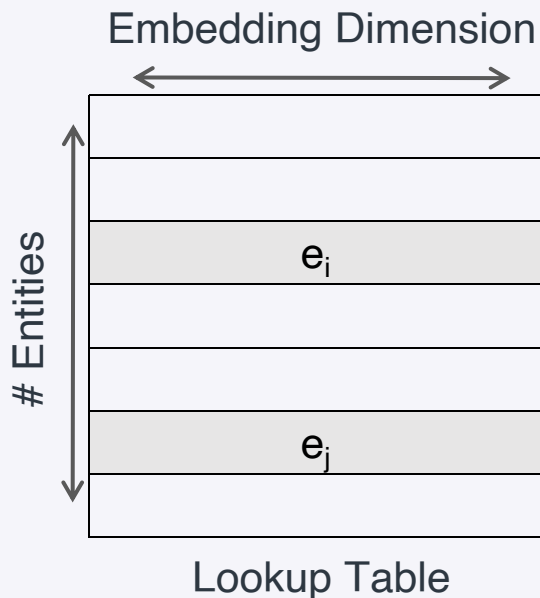
\mathbf{V} "Meryl Streep" + \mathbf{V} "starred-in" = \mathbf{V} "Death becomes her"

\mathbf{V} "Meryl Streep" + \mathbf{V} "starred-in" = \mathbf{V} "Doubt"

Movies are
very different:
cast, genre,
director, etc.



Shallow and Deep Models



Shallow Models

- Representation expressiveness depends on the embedding dimension.
- Transductiveness.



f

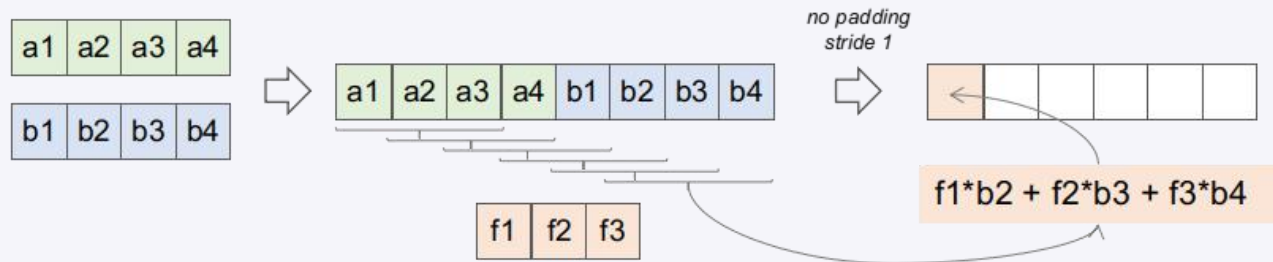


Deep models

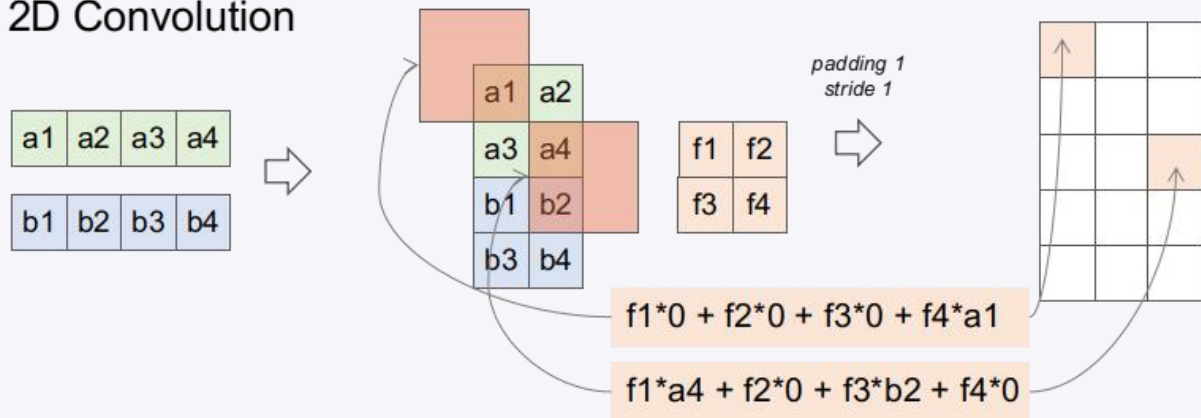
- Overfitting.
- Increase in time and space complexity.

ConvE

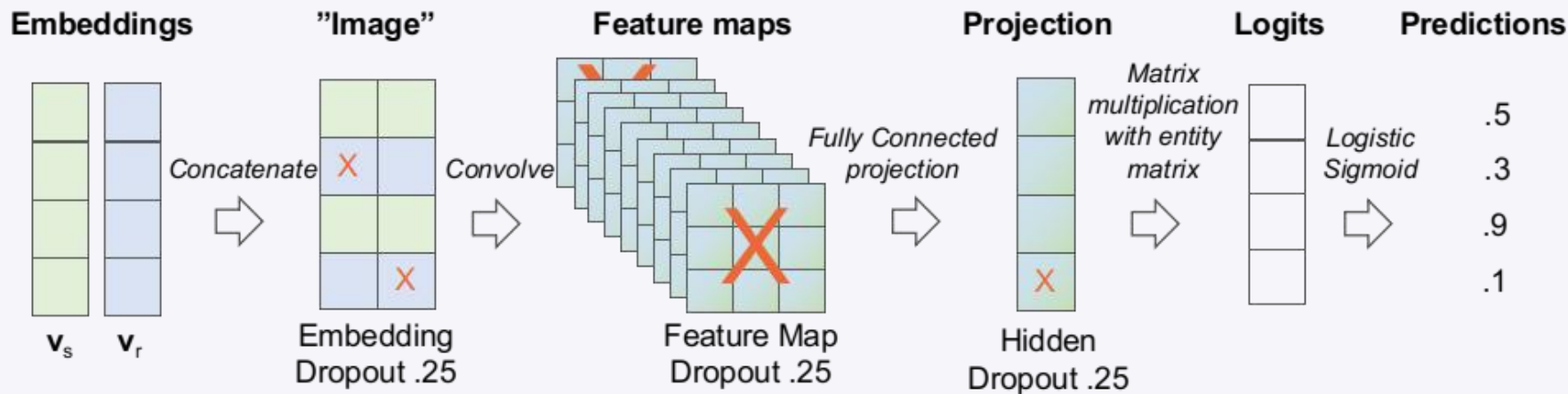
1D Convolution



2D Convolution

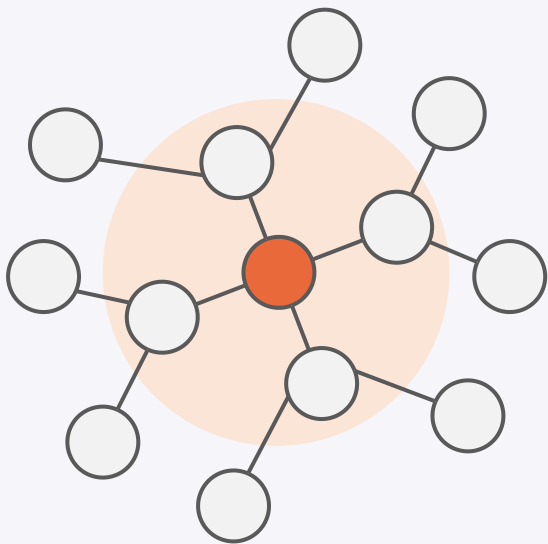


ConvE



$$\begin{aligned}
 \mathbf{M} &= \text{concat}(\bar{\mathbf{v}}_s, \bar{\mathbf{v}}_r) \\
 \mathbf{F} &= \text{conv}_\omega(\mathbf{M}) \\
 \mathbf{p} &= \text{vec}(\sigma_1(\mathbf{F}))\mathbf{W} \\
 \text{score} &= \langle \sigma_2(\mathbf{p}), \mathbf{v}_o \rangle
 \end{aligned}$$

Graph Convolution



Generalize the operation of *convolution* from grid data to graph data.

Main idea: Learn a representation for a node taking into account its neighbors representations.

Differentiable message-passing framework

$$\mathbf{h}_u^{(k)} = \gamma^{(k)} \left(\mathbf{h}_u^{(k-1)}, \bigoplus_{v \in \mathcal{N}(u)} \psi^{(k)} \left(\mathbf{h}_u^{(k-1)}, \mathbf{h}_v^{(k-1)}, \mathbf{e}_{v,u} \right) \right)$$

↑
u representation
at layer k

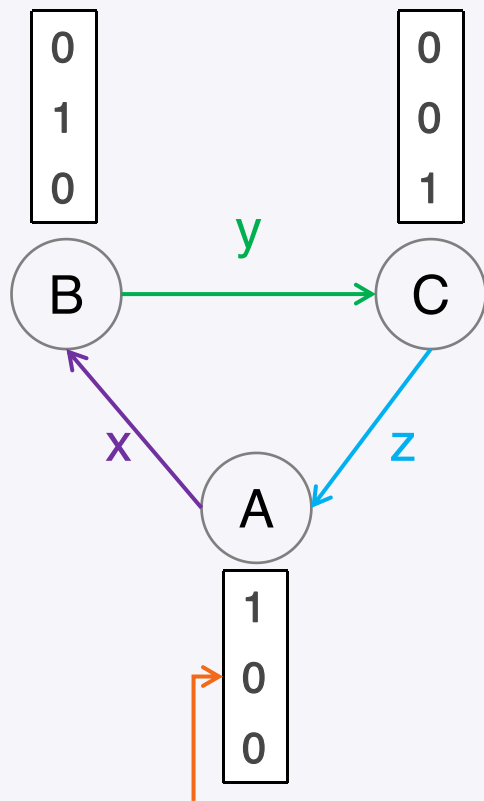
↑
(In)neighbors of u

↑
(v, u) edge
representation

Differentiable Functions: $\gamma^{(k)}$ \bigoplus $\psi^{(k)}$

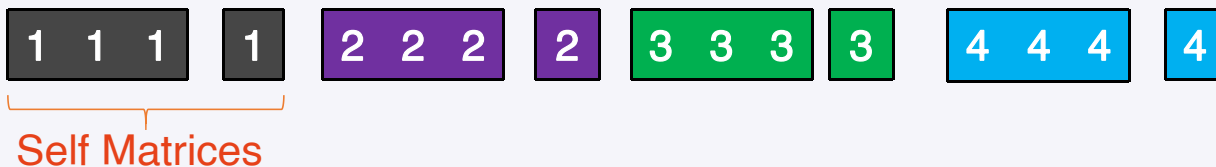
↑
Must be **Permutation Invariant**
(e.g., sum, max, deep sets NN)

**GNNs is a very hot
topic now!**



Initial Representation
One Hot Encode

Relation Matrices for 1 and 2 layer, resp.



$$h_A^{(1)} = \text{ReLu}(\begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 4 & 4 & 4 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}) = \begin{bmatrix} 5 \end{bmatrix}$$

$$h_A^{(2)} = \text{ReLu}(\begin{bmatrix} 1 & 5 \end{bmatrix} + \begin{bmatrix} 4 & 4 \end{bmatrix}) = \begin{bmatrix} 21 \end{bmatrix}$$

* $\text{ReLu}(x) = \max(x, 0)$

Attributive Relations and handling literals

KGs often include:

Numerical attributes: e.g., ages, dates, financial, and geoinformation.

Textual attributes: e.g., names, descriptions, and titles.

Images: e.g., profile photos, flags, and posters.

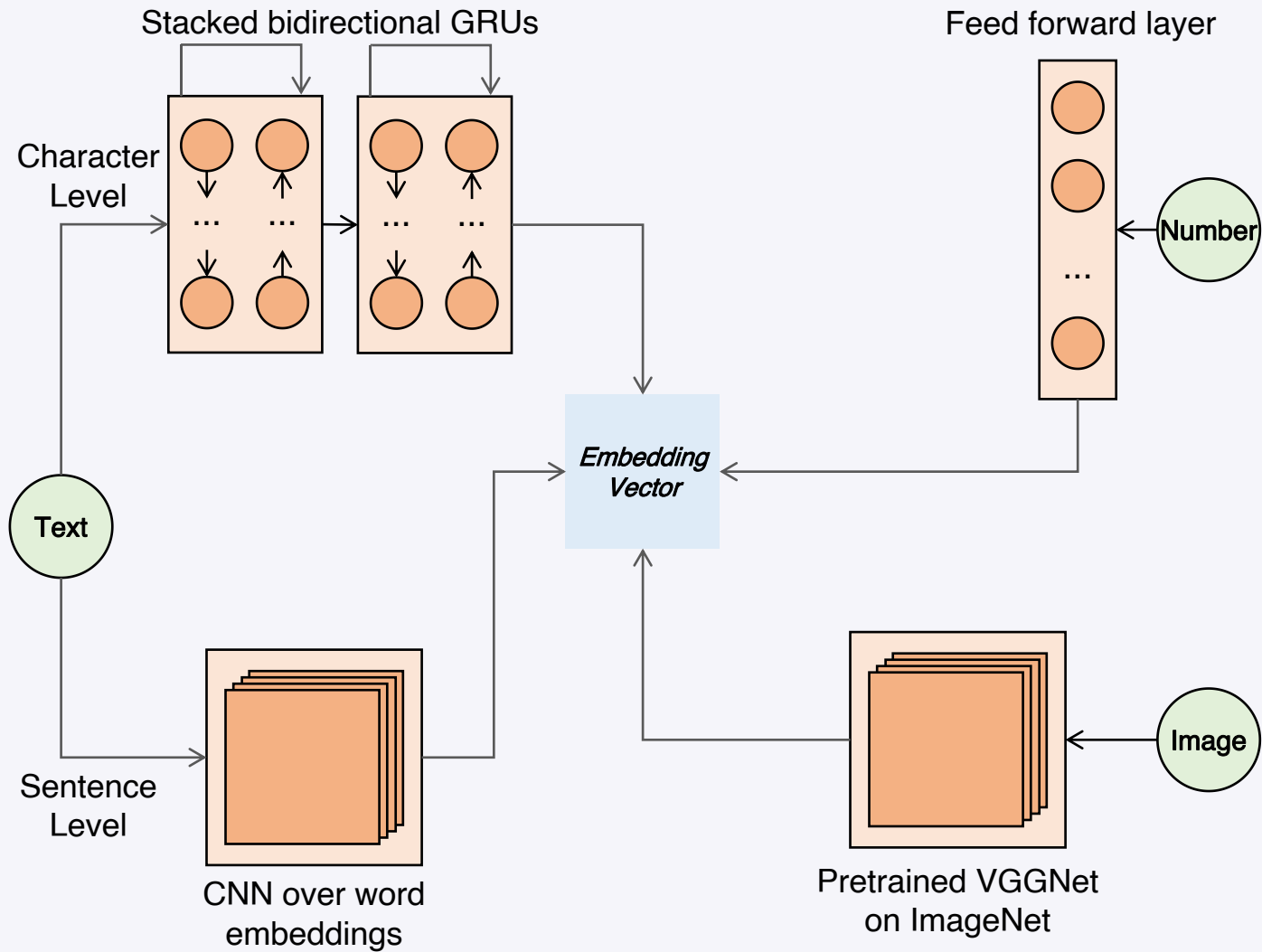
Useful for entities with few relationships

Naive Strategy: regard attributes as entities.

MKBE: Multimodal Knowledge Base Embeddings

Compositional encoding component: Different neural encoders for the variety of observed data.

Embedding Multimodal Data: Structured knowledge, numerical, text, and images.



Ontology Usage Example: JOIE

Instance-view, ontology-view, and cross view.

Cross-view association model:

- Cross-view grouping (CG).

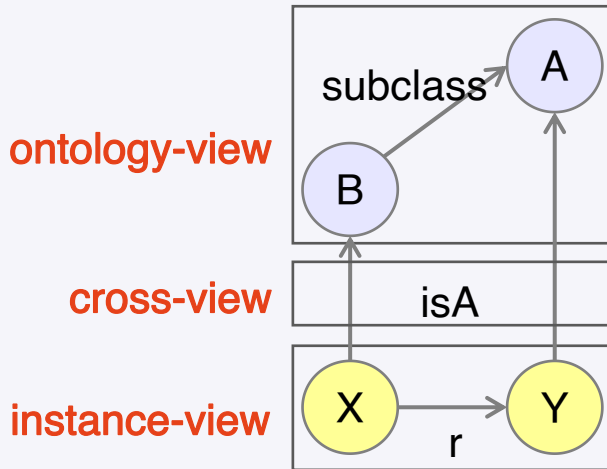
- Cross-view transformation (CT).

Intra-view model:

- Default.

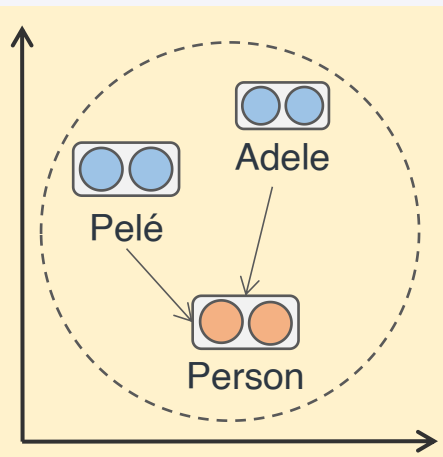
- Hierarchy-aware.

Specific Cost Functions.



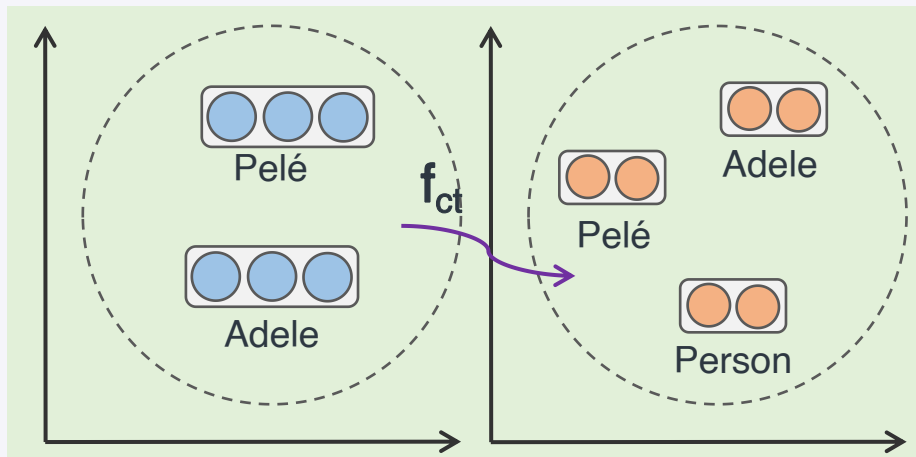
JOIE

entity space



1. An entity vector should be encompassed by its class vector radius.

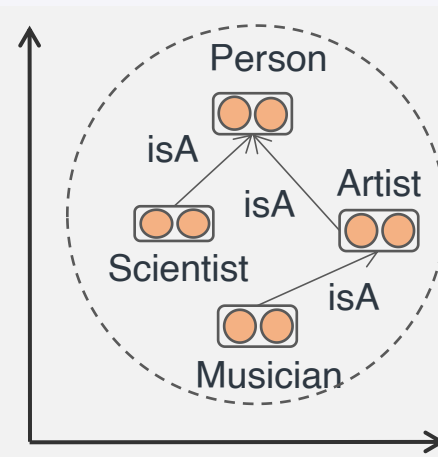
entity space



2. The model should be capable of mapping the entity space into the concept space.

concept space

concept space




3. The concept space should reflect the concept hierarchical structure.

Knowledge Hypergraphs

Triples often oversimplifies the complex nature of the data, in particular for hyper-relational data.

Retified Representation:

(Marie Curie, **educated-at**, University of Paris)



[academic-major: Physics]
[academic-degree: Master of Science]

Hyper-relational representation

education

subject: Marie Curie

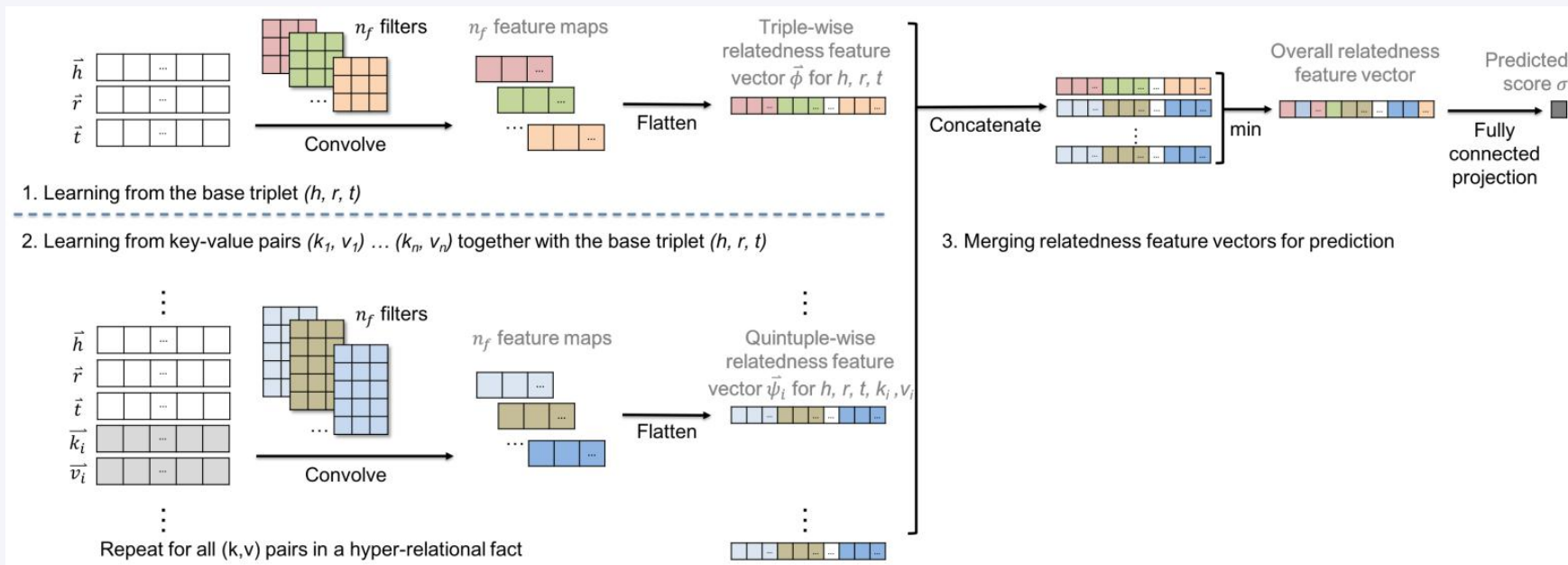
object: University of Paris

major: Physics

degree: Master of Science

HINGE

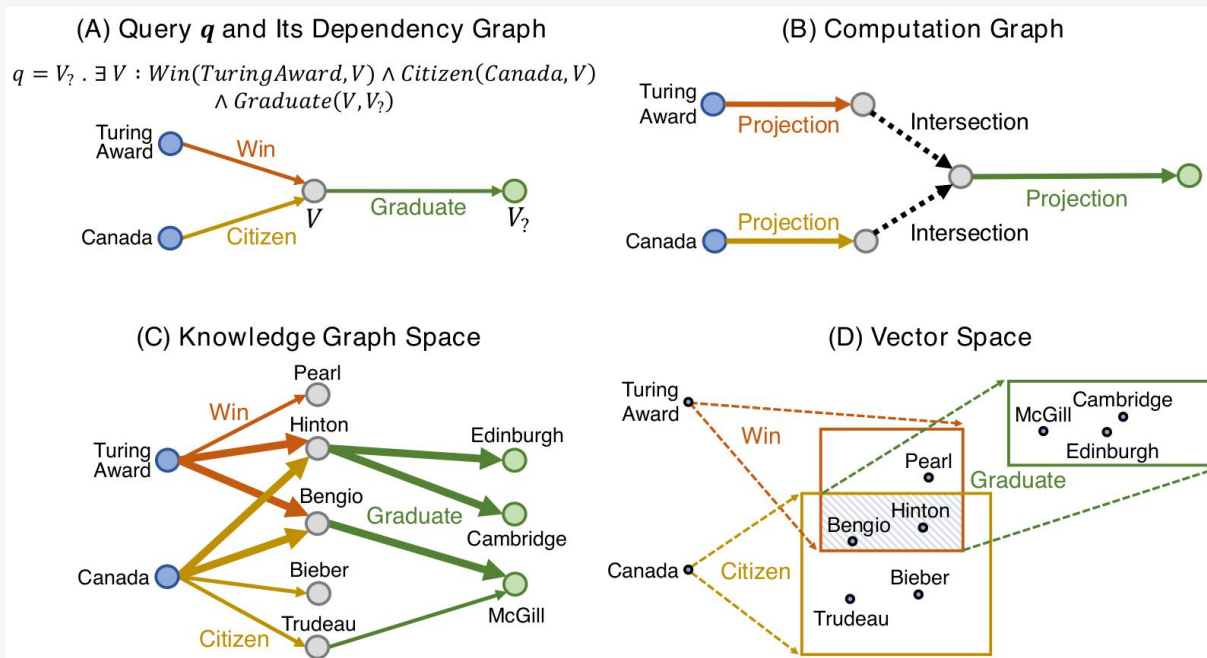
Triple (h, r, t) is associated with n key-value pairs (k_i, v_i)



Beyond Triplet Reasoning

What about answer complex KG queries on vector space by embedding logical operators?

Why: subgraph matching may be exponential and partially observed data.



Conferences

Artificial Intelligence:

AAAI	ECML PKDD	ICLR
ICML	IJCAI	NIPS

Data Management:

CIKM	ESWC	KDD	SIGMOD PODS
SIGIR	VLDB	WSDM	WWW

Natural Language Processing:

ACL	EMNLP
-----	-------

- Automated Knowledge Base Construction.
- Knowledge Graph Conference.
- International Workshop on Challenges and Experiences from Data Integration to Knowledge Graphs.
- Workshop on Knowledge Graph Technology and Applications.
- Workshop on Deep Learning for Knowledge Graphs.

KGE libraries and Systems

Ampligraph	Tensorflow	https://ampligraph.org/
DeepGraphLibrary	MXNet/Gluon, PyTorch, Tensorflow	https://www.dgl.ai/
Grakn KGLIB	Tensorflow	https://github.com/graknlabs/kglib
Graph Nets	Sonnet	https://github.com/deepmind/graph_nets
OpenKE	PyTorch	http://openke.thunlp.org
LibKGE	PyTorch	https://github.com/uma-pi1/kge
Pykg2vec	Tensorflow	https://github.com/Sujit-O/pykg2vec
PyKEEN	PyTorch	http://pykeen.readthedocs.io
PyTorch-BigGraph	PyTorch	https://torchbiggraph.readthedocs.io/en/latest
PyTorch Geometric	PyTorch	https://pytorch-geometric.readthedocs.io/en/latest/
StellarGraph	Tensorflow	https://www.stellargraph.io/

KB / KG construction

Research

- Continuously learning and self-correcting systems.
- Entity disambiguation and managing identity.
- Heterogeneous and multimodal information.
- KBC in specific domains.
- Knowledge graph alignment.
- Managing operations at scale.
- Multi-language knowledge bases.
- Virtual knowledge graphs.

KG Reasoning

Research

- Compare/Combine KGE to/with other approaches (e.g., PSL + Rule Mining).
- Dynamic Negative Sampling.
- Embed other structures: graphlets, paths, motifs, **queries**, etc.
- Few shot Learning.
- KGE and the lack of symbolic structures (e.g., rules and restrictions).
- KGE consistency (e.g., embed formal knowledge).
- KG sparsity and uncertainty.
- KG temporal and spatial dynamics.
- Inductive vs. transductive learning.
- Prediction calibration of KGE models.
- Representation: hypergraphs (n-ary relations), meta-properties, multidimension KGs.

Aprendizado de Máquina aplicado a Grafos de Conhecimento

Daniel N. R. da Silva, Artur Ziviani e Fabio Porto
dramos, ziviani, fporto@Incc.br

Machine Learning for Knowledge Graphs

Model Training and Evaluation

Triple Classification Protocol:

Test the model's ability to discriminate between true and false triples.

Triple (s,r,o) is classified as positive if its score exceeds a relation-specific decision threshold (learned on validation data).

Entity Ranking Protocol:

Assess model performance in terms of ranking answers to certain questions.

Evaluation Metrics

$$hits@k = \frac{1}{|Q|} \sum_{q \in Q} \mathbf{1}_{\leq k}(\text{rank}_q)$$

$$MRR = \frac{1}{|Q|} \sum_{q \in Q} \frac{1}{\text{rank}_q}$$

Regression, classification and **ranking metrics**: e.g., Hits@k, Mean Reciprocal Rank (MRR), and Precision Recall Curve.

Test triples:

(Neymar, born-in, Mogi)

Entities:

Kaká, Neymar, Mogi, and Gama

avg hits@1 = (1 + 0) / 2 = 0.5

avg mrr = (1 + 1/2) / 2 = 0.75

s	p	o	score	rank	
Neymar	born-in	Mogi	0.80	1	2
Neymar	born-in	Gama	0.70	2	
Neymar	born-in	Kaká	0.20	3	

P o i n t w i s e	<i>Square Error Loss</i>	$\frac{1}{2} \sum_{t \in T_{\text{train}}} (\phi_t - y_t)^2$	Closed World Assumption
	<i>Hinge Loss</i>	$\sum_{t \in T_{\text{train}}} [\lambda + (-1)^{y_t} \phi_t]_+$	
	<i>Logistic Loss</i>	$\sum_{t \in T_{\text{train}}} \log(1 + \exp((-1)^{y_t} \phi_t))$	
P a i r w i s e	<i>Hinge Loss</i>	$\sum_{t \in T_{\text{train}}^{(+)}} \sum_{t' \in T_{\text{train}}^{(-)}} [\lambda + \phi_{t'} - \phi_t]_+$	
	<i>Logistic Loss</i>	$\sum_{t \in T_{\text{train}}^{(+)}} \sum_{t' \in T_{\text{train}}^{(-)}} \log(1 + \exp(\phi_{t'} - \phi_t))$	

ϕ_t Score given for triple t

y_t Label (0 or 1) for triple t