

# Embedding Complex Knowledge: From Geometric to Language Models

Jiaoyan Chen

Department of Computer Science

The University of Manchester

3<sup>rd</sup> March, 2026

KMI Seminar, The Open University

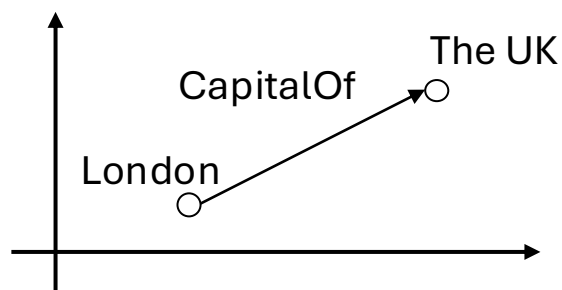


# Embedding Symbolic Knowledge

- **Vector** or **parameter**-based representation of symbolic knowledge
- Why?
  - Knowledge inference with **uncertainty** (e.g., incompleteness, approximation & prediction, induction of schema & rule)
  - Similarity-based **matching** across modalities (e.g., retrieval, alignment and resolution)
  - **Inject knowledge** into parameter-based models (e.g., tuning LLM)
  - Kinds of **downstream applications** with machine learning and statistical models

# Knowledge Graph Embedding

- Originate from **word embeddings**
  - Represent word cooccurrence and correlation by a neural network; e.g., similarity: (cat, kitten) > (cat, dog)
- Mostly aim at facts of RDF triples e.g., <London, CapitalOf, The UK>
  - **Geometric modeling**: TransE, TransR, TransH, ..



TransE: modeling relation as a translation mapping

- **Sequence learning**: RDF2Vec, ...
- **Graph propagation**: R-GCN, ...

# Ontology (Description Logic) Embedding

- How to represent more complex ontologies of Description Logic (DL) in Euclidean space?

$\mathcal{T} = \{\text{Father} \sqsubseteq \text{Parent} \sqcap \text{Male}, \text{Mother} \sqsubseteq \text{Parent} \sqcap \text{Female},$   
 $\text{Child} \sqsubseteq \exists \text{hasParent.Father}, \text{Child} \sqsubseteq \exists \text{hasParent.Mother},$   
 $\text{hasParent} \sqsubseteq \text{relatedTo}\}$   
 $\mathcal{A} = \{\text{Father}(\text{Alex}), \text{Child}(\text{Bob}), \text{hasParent}(\text{Bob}, \text{Alex})\}$

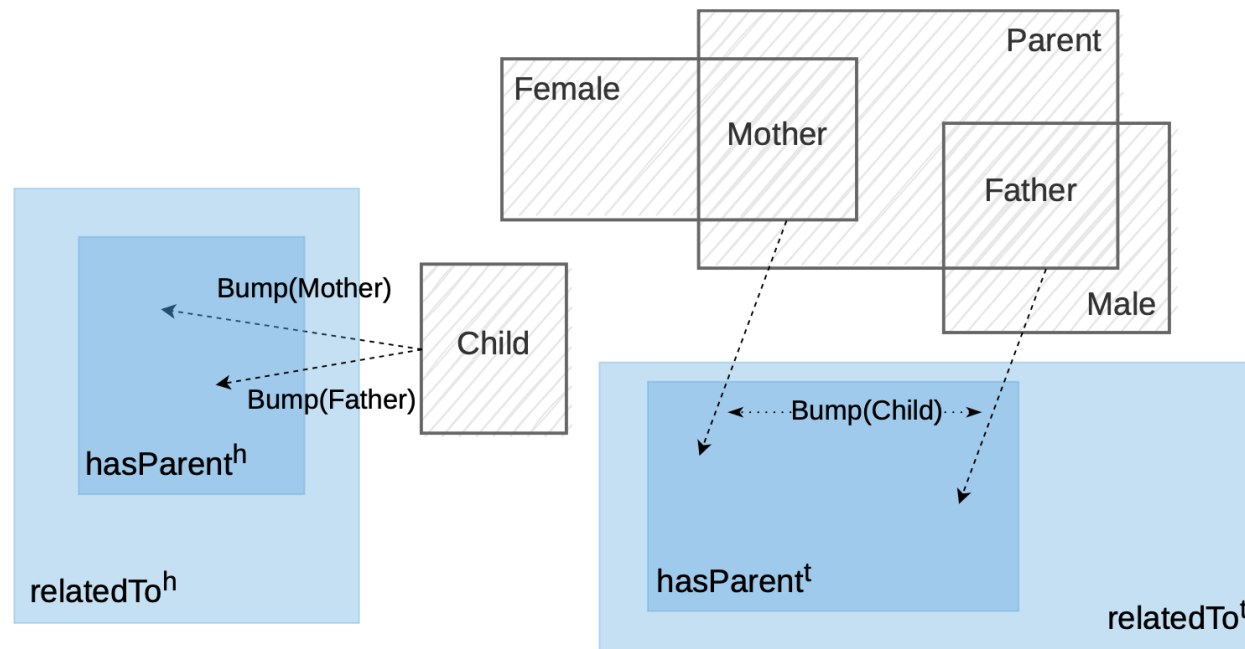
A toy famil ontology in DL  $\mathcal{EL}^{++}$  which allows complex concept construction:

$$\perp \mid \top \mid A \mid C \sqcap D \mid \exists r. C \mid \{a\}$$

- Embedding: Region-based
  - Individual – Point
  - Concept – Ball, Box, ...

# Ontology (Description Logic) Embedding

- Example: Box<sup>2</sup>EL
  - Individual: one n-point
  - Concept: one n-box
    - Conjunction, subsumption, membership
  - Relation: two n-boxes (head & tail)
    - Composition, subsumption
  - Concept interaction: bumping vector
    - Existential quantification  
 $Child \sqsubseteq \exists hasParent. Father$



Representation of the family ontology in Box<sup>2</sup>EL

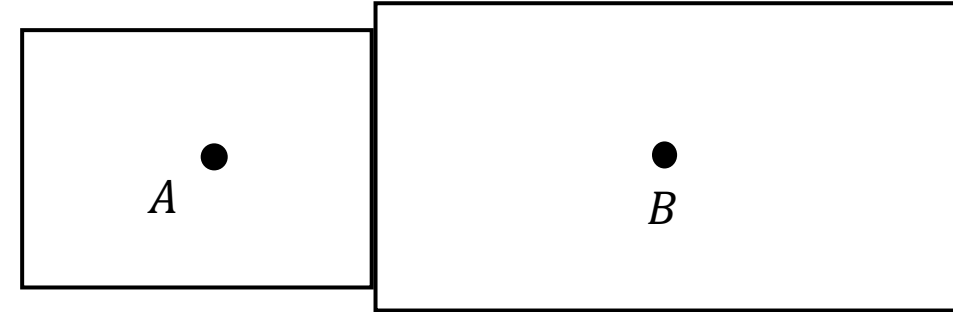
Jackermeier, Mathias, Jiaoyan Chen, and Ian Horrocks. "Dual box embeddings for the description logic EL++." *Proceedings of the ACM Web Conference 2024*. 2024.

# Ontology (Description Logic) Embedding

- Box<sup>2</sup>EL training:

- Element-wise distance of two boxes:

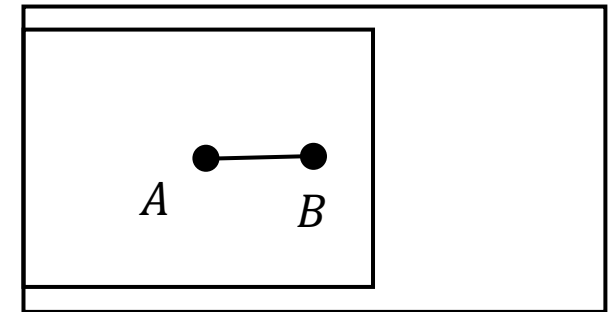
$$d(A, B) = |c(A) - c(B)| - o(A) - o(B)$$



The horizontal dimension distance is 0

- Score/loss of concept subsumption (inclusion of two boxes):

$$\mathcal{L}_{\subseteq}(A, B) = \begin{cases} \|\max\{0, d(A, B) + 2o(A) - \gamma\}\| & \text{if } B \neq \emptyset \\ \max\{0, o(A)_1 + 1\} & \text{otherwise,} \end{cases}$$



The horizontal dimension loss is 0

Box A moves **left**, the horizontal dimension loss > 0, penalize the embeddings by training;

Box A moves **right**, the horizontal dimension loss is still 0

# Ontology (Description Logic) Embedding

- Box<sup>2</sup>EL training: loess/scores for axioms of each normal form (NF)
  - NF1:  $C \sqsubseteq D$        $\mathcal{L}_1(C, D) = \mathcal{L}_{\sqsubseteq}(\text{Box}(C), \text{Box}(D))$ .
  - NF2:  $C \sqcap D \sqsubseteq E$        $\mathcal{L}_2(C, D, E) = \mathcal{L}_{\sqsubseteq}(\text{Box}(C) \cap \text{Box}(D), \text{Box}(E))$
  - NF3:  $C \sqsubseteq \exists r. D$        $\mathcal{L}_3(C, r, D) = \frac{1}{2} \left( \mathcal{L}_{\sqsubseteq}(\text{Box}(C) + \text{Bump}(D), \text{Head}(r)) \right. \\ \left. + \mathcal{L}_{\sqsubseteq}(\text{Box}(D) + \text{Bump}(C), \text{Tail}(r)) \right)$ .
  - NF4:  $\exists r. C \sqsubseteq D$        $\mathcal{L}_4(r, C, D) = \mathcal{L}_{\sqsubseteq}(\text{Head}(r) - \text{Bump}(C), \text{Box}(D))$
  - NF5:  $C \sqcap D \sqsubseteq \perp$        $\mathcal{L}_5(C, D) = \|\max\{\mathbf{0}, -(\mathbf{d}(\text{Box}(C), \text{Box}(D)) + \gamma)\}\|$

# Ontology (Description Logic) Embedding

- Box<sup>2</sup>EL training: Loess for axioms of each normal form (NF)

- NF6:  $r \sqsubseteq s$        $\mathcal{L}_6(r, s) = \frac{1}{2} \left( \mathcal{L}_{\sqsubseteq}(\text{Head}(r), \text{Head}(s)) + \mathcal{L}_{\sqsubseteq}(\text{Tail}(r), \text{Tail}(s)) \right)$

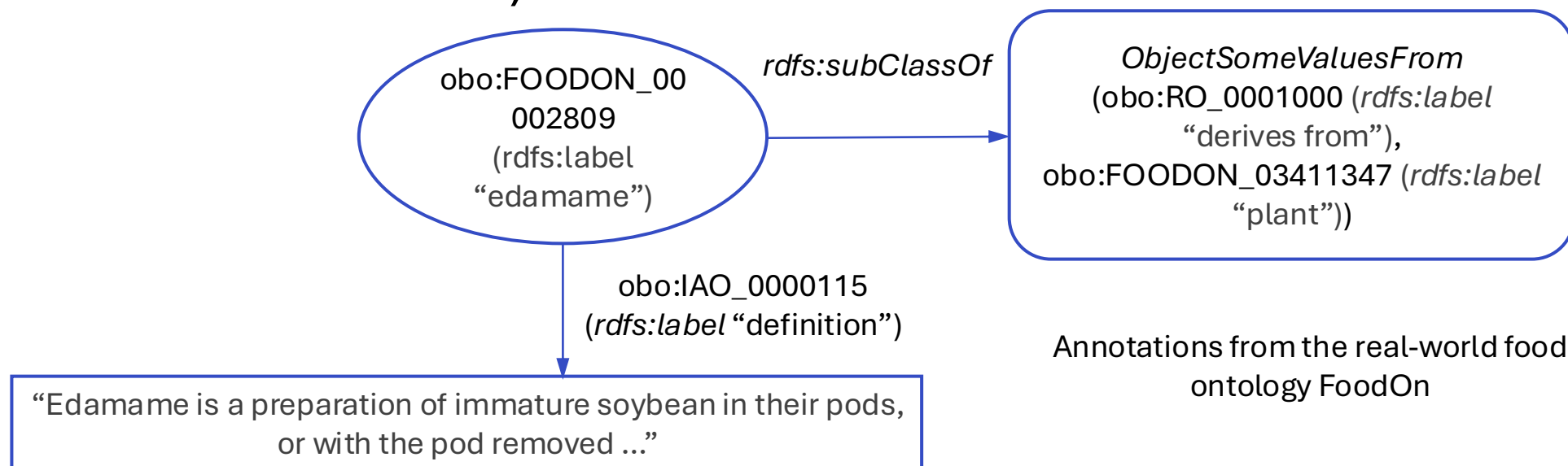
- NF7:  $r_1 \circ r_2 \sqsubseteq s$        $\mathcal{L}_7(r_1, r_2, s) = \frac{1}{2} \left( \mathcal{L}_{\sqsubseteq}(\text{Head}(r_1), \text{Head}(s)) + \mathcal{L}_{\sqsubseteq}(\text{Tail}(r_2), \text{Tail}(s)) \right)$

ABox can be transformed into TBox:

$$\begin{aligned} C(a) &\rightsquigarrow \{a\} \sqsubseteq C \\ r(a, b) &\rightsquigarrow \{a\} \sqsubseteq \exists r. \{b\} \end{aligned}$$

# Text-aware Ontology Embedding

- OWL ontology includes more than formal semantics (e.g., labels, textual definitions)



- How to jointly embed the informal textual knowledge and the formally defined knowledge?

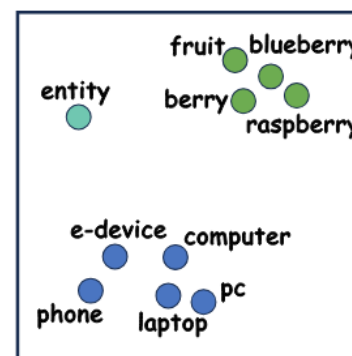
# Text-aware Ontology Embedding

- Non-contextual: word2vec
  - **OWL2Vec\***: Train a Word2Vec model from an OWL ontology, following RDF2Vec
  - Corpus (sentences) extraction via seriation (Manchester OWL Syntax), walking on the graph, OWL to RDF projection, etc.
- Contextual: Transformer-based language models
  - Common task-specific paradigm: pre-train then fine-tune
  - E.g., **BERTMap** which fine-tunes a BERT alike language model for ontology alignment

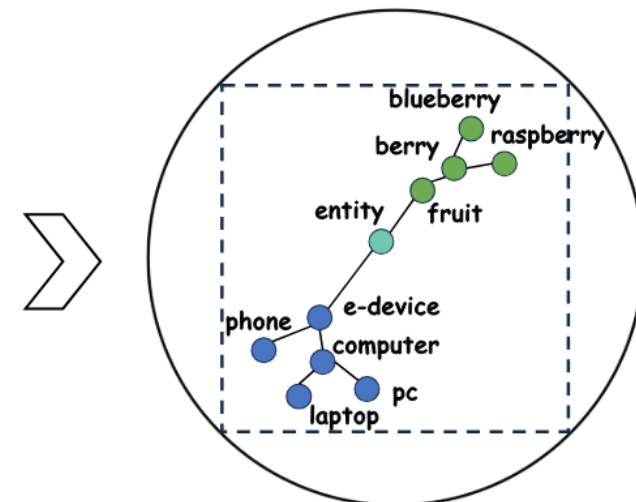
They both lose formally defined semantics!

# Text-aware Ontology Embedding

- LM as hierarchy encoder (**HiT**)
  - Re-train a BERT alike LM by an ontology
  - Force the LM's concept encodings to a hierarchy in a hyperbolic space (Poincare ball)
    - Motivated by its efficiency for representing hierarchies



Concept's text embedding in Euclidean Space by the last layer (tanh activation) of an LM, which is in a  $d$ -dimensional hyper-cube



Concept's text embedding by an ontology retrained LM, which is in a Poincare Ball of radius  $\sqrt{d}$  that circumscribe the hyper cube

# Text-aware Ontology Embedding

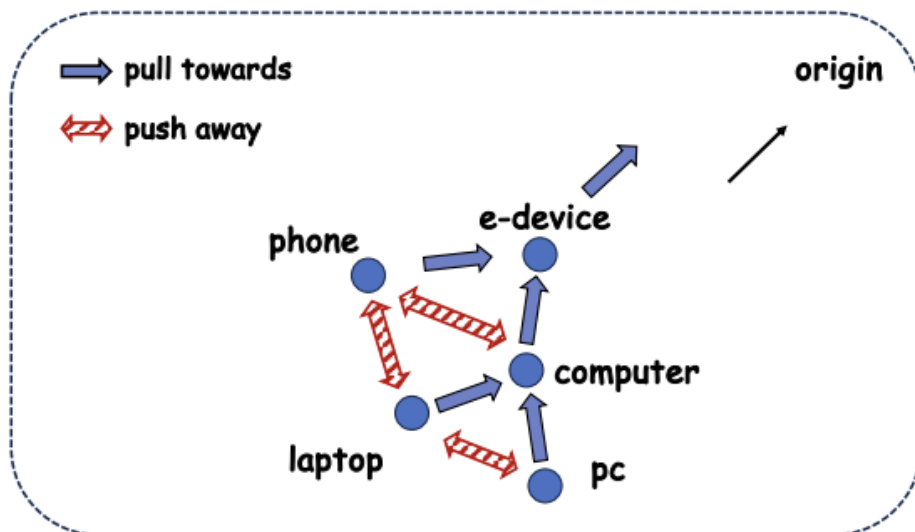
- Training of **HiT**

- Hyperbolic clustering loss: clustering related concepts and distancing unrelated concepts (a concept is close to its positive parent and distanced from its negative parent)

$$\mathcal{L}_{cluster} = \sum_{(e, e^+, e^-) \in \mathcal{D}} \max(d_c(e, e^+) - d_c(e, e^-) + \alpha, 0)$$

- Hyperbolic centripetal loss: make parent closer to origin

$$\mathcal{L}_{centri} = \sum_{(e, e^+, e^-) \in \mathcal{D}} \max(\|e^+\|_c - \|e\|_c + \beta, 0)$$



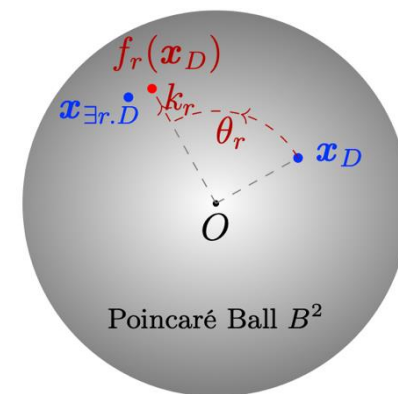
- Subsumption inference with HiT embeddings

- Consider both contrastive and centripetal losses

$$s(e_1 \sqsubseteq e_2) = -(d_c(e_1, e_2) + \lambda(\|e_2\|_c - \|e_1\|_c))$$

# Text-aware Ontology Embedding

- LM as  $\mathcal{EL}^{++}$  ontology encoder (**OnT**)
  - Extend HiT to complex concepts E.g.,  $\exists isParentOf.Person$
- Solution #1: verbalization
  - $\exists isParentOf.Person \rightarrow$  “something that is parent of some person”
  - Embedding notion:  $\exists r.D \rightarrow \mathbf{x}_{\exists r.D}$
  - Does not represent roles’ logics e.g., if  $A \sqsubseteq B$  then  $\exists r.A \sqsubseteq \exists r.A$
- Solution #2: Relation by rotation
  - A rotation function to represent a relation:  $\exists r.D \rightarrow f_r(\mathbf{x}_D)$
  - Learning:  $\mathbf{x}_{\exists r.D} < f_r(\mathbf{x}_D), f_r(\mathbf{x}_D) < \mathbf{x}_{\exists r.D}$



# Summary

- Geometric models
  - **Box<sup>2</sup>EL** (concept as a box and a bump vector, relation as two boxes)
- Language models
  - Non-contextual (OWL2Vec\*)
  - Transformer-based
    - **BERTMap**: fine-tune a pre-trained encoder model with an additional classifier for concept equivalence matching
    - **HiT**: Continuously train a pre-trained encoder model for representing hierarchical concepts in a Poincare ball
    - **OnT**: Extend HiT for conserving relationship between concepts

# Discussion on Application

- Link prediction in knowledge intensive domains
  - E.g., **protein-protein interaction prediction** with the Gene Ontology as evaluated in the Box<sup>2</sup>EL paper
- Ontology construction and curation
  - **Subsumption completion**, including named concepts and complex concepts in either formal representation or natural language (evaluated in Box<sup>2</sup>EL, HiT and OnT papers)
  - **Subsumption matching and new concept placement**; can be combined with LLM and agentic AI, where such embeddings can efficiently get candidates and reduce space space via indexing
- Hierarchical retrieval
  - E.g., most queries to SNOMED CT are **out-of-vocabulary**, which require to return most close concepts with a subsumption relationship instead of equivalence

- Acknowledgement
  - Key contributors: Yuan He (Oxford, now in Amazon), Hui Yang (Manchester), Mathias Jackermeier (Oxford), Ian Horrocks (Oxford)
  - Main funders: EPSRC (OntoEm & ConCur), Samsung Research UK
- Feel free to contact me
  - [jiaoyan.chen@manchester.ac.uk](mailto:jiaoyan.chen@manchester.ac.uk)
- Q&A