



# Towards Ontology-Enhanced Representation Learning for Large Language Models

Francesco Ronzano and Jay Nanavati  
IQVIA Advanced NLP Team

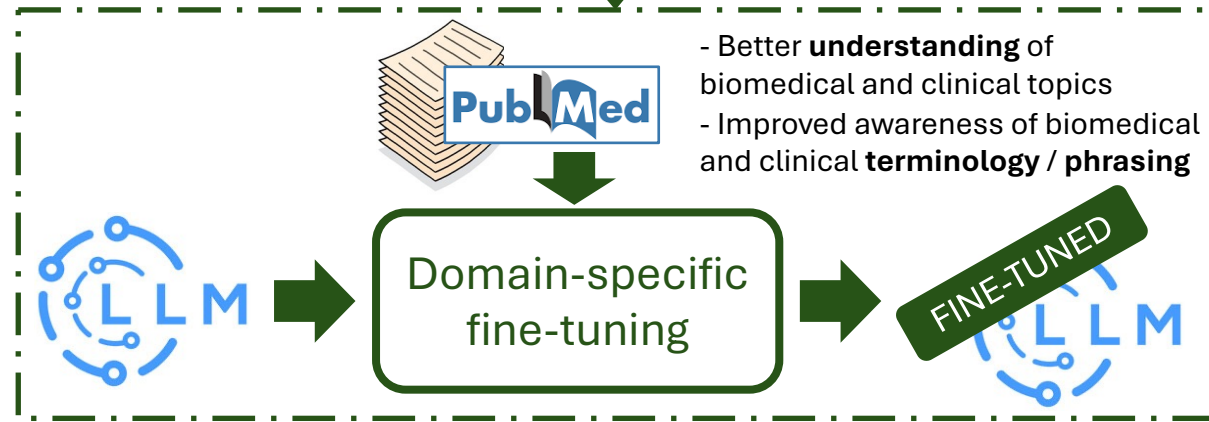
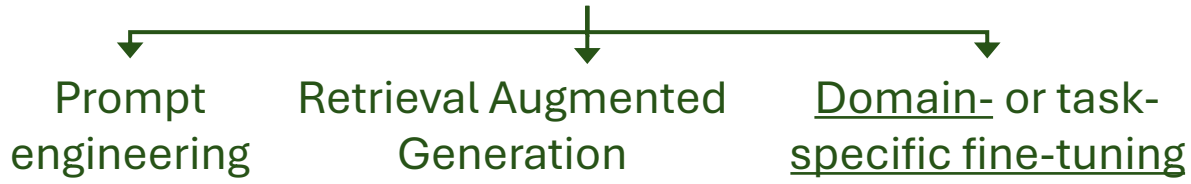


# Agenda

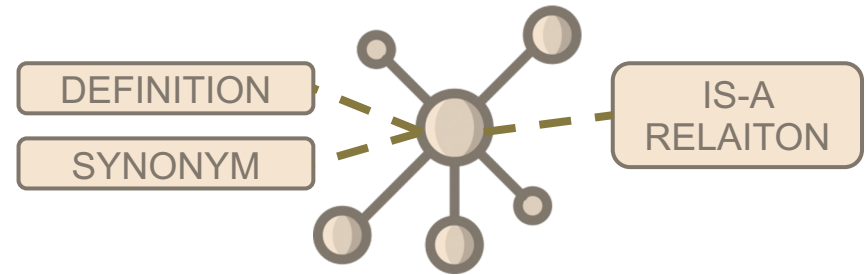
- + Exploiting ontologies to teach Large Language Models
- + The ontological knowledge infusion approach
  - Choosing the *source ontology* and the *target embedding-Large Language Model*
  - Generating *synthetic definitions of concepts*
  - Selecting *positive and negative pairs of concept definitions*, driven by the ontology
  - Fine-tuning the embedding-Large Language Model by *contrastive representation learning*
- + Evaluation: do embedding-Large Language Models better understand diseases, after infusing the MONDO disease ontology?
- + Key learnings
- + Next steps

# Exploiting ontologies to teach Large Language Models

## DOMAIN SPECIALIZATION OF LARGE LANGUAGE MODELS



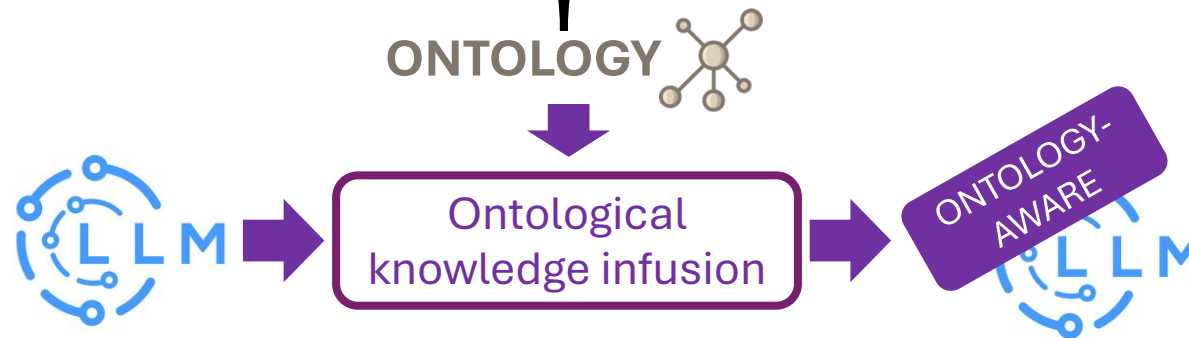
## ONTOLOGY



BioPortal

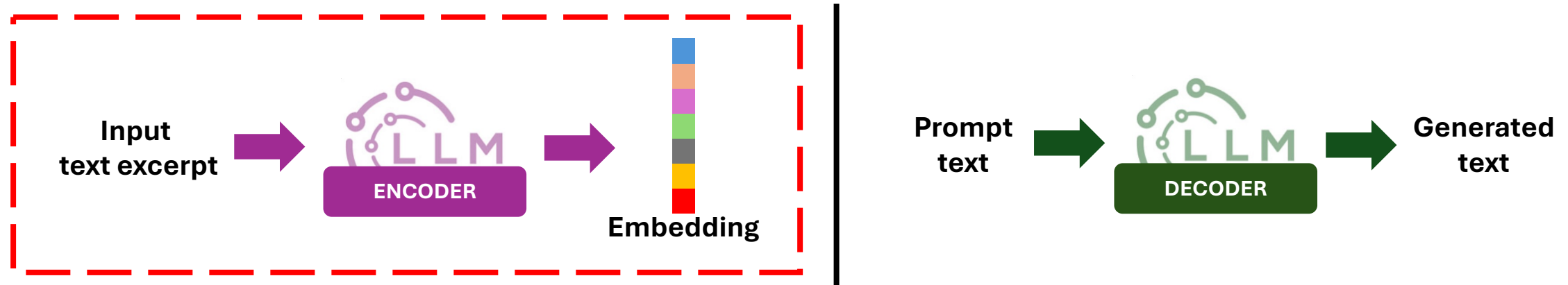
the world's most comprehensive repository of biomedical ontologies

- 1,033 biomedical ontologies
- 15,5 million concepts
- 36,200 properties
- 100 million mappings between ontologies



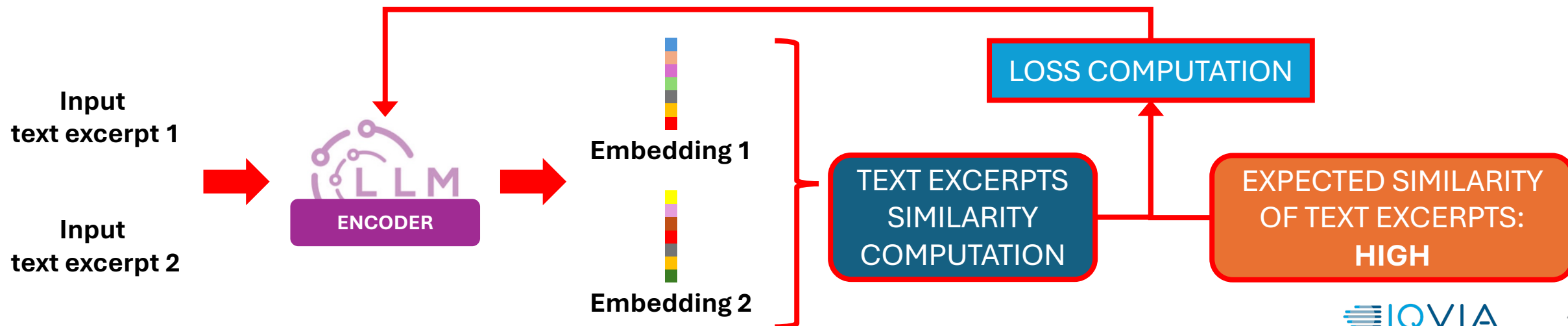
# The ontological knowledge infusion approach

Which Large Language Model architecture do we target to infuse ontological knowledge?



Which fine-tuning framework is exploited to support ontological knowledge infusion?

Contrastive learning framework

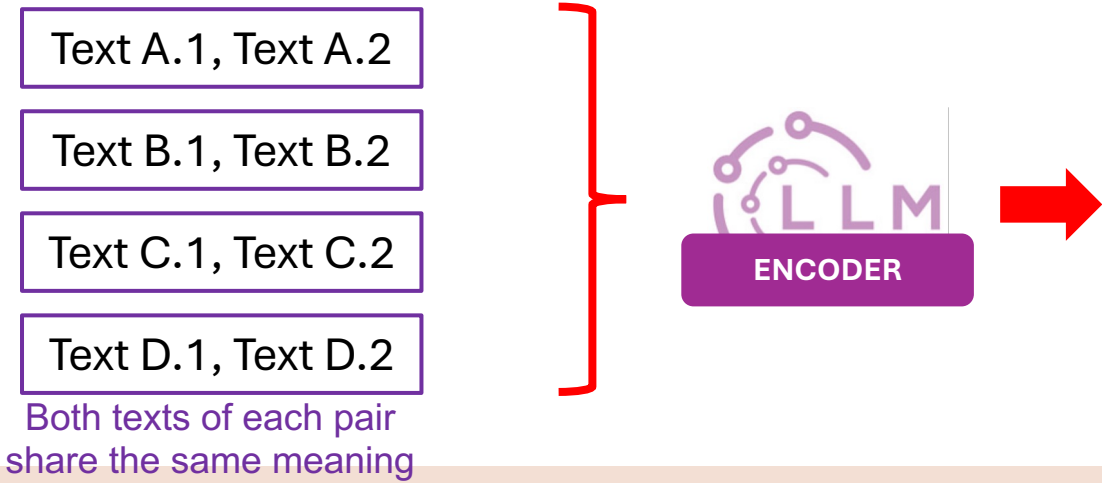


# The ontological knowledge infusion approach

Text A.1: A car is a four-wheeled road vehicle.  
Text A.2: A car is a mean of transport moving on wheels.

## Which fine-tuning framework is exploited to support ontological knowledge infusion? (cont.)

Contrastive training objective: **InfoNCE**



Text similarity matrix

	A.2	B.2	C.2	D.2
A.1	sim(A.1, A.2)	sim(A.1, B.2)	sim(A.1, C.2)	sim(A.1, D.2)
B.1	sim(B.1, A.2)	sim(B.1, B.2)	sim(B.1, C.2)	sim(B.1, D.2)
C.1	sim(C.1, A.2)	sim(C.1, B.2)	sim(C.1, C.2)	sim(C.1, D.2)
D.1	sim(D.1, A.2)	sim(D.1, B.1)	sim(D.1, C.2)	sim(D.1, D.2)

sim: similarity function among text embeddings

### COMPUTATION OF CATEGORICAL CROSS-ENTROPY LOSS

*Predicted probability distribution*

softmax(	sim(A.1, A.2)	sim(A.1, B.2)	sim(A.1, C.2)	sim(A.1, D.2)	)
softmax(	sim(B.1, A.2)	sim(B.1, B.2)	sim(B.1, C.2)	sim(B.1, D.2)	)
softmax(	sim(C.1, A.2)	sim(C.1, B.2)	sim(C.1, C.2)	sim(C.1, D.2)	)
softmax(	sim(D.1, A.2)	sim(D.1, B.2)	sim(D.1, C.2)	sim(D.1, D.2)	)

*Expected probability distribution*

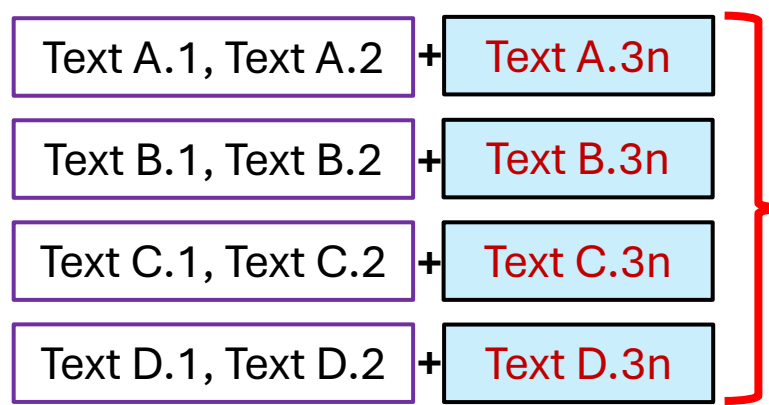
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

# The ontological knowledge infusion approach

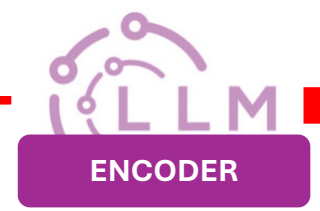
Text A.1: A car is a four-wheeled road vehicle.  
Text A.2: A car is a mean of transport moving on wheels.  
Text A.3n: Two over six wheels of that bus were damaged because of the accident with a car.

## Which fine-tuning framework is exploited to support ontological knowledge infusion? (cont.)

Contrastive training objective: **InfoNCE** with (hard-)negative texts



Both texts of each pair share the same meaning



Text similarity matrix

	A.2	B.2	C.2	D.2	A.3n	B.3n	C.3n	D.3n
A.1	sim(A.1, A.2)	sim(A.1, B.2)	sim(A.1, C.2)	sim(A.1, D.2)	sim(A.1, A.3n)	sim(A.1, B.3n)	sim(A.1, C.3n)	sim(A.1, D.3n)
B.1	sim(B.1, A.2)	sim(B.1, B.2)	sim(B.1, C.2)	sim(B.1, D.2)	sim(B.1, A.3n)	sim(B.1, B.3n)	sim(B.1, C.3n)	sim(B.1, D.3n)
C.1	sim(C.1, A.2)	sim(C.1, B.2)	sim(C.1, C.2)	sim(C.1, D.2)	sim(C.1, A.3n)	sim(C.1, B.3n)	sim(C.1, C.3n)	sim(C.1, D.3n)
D.1	sim(D.1, A.2)	sim(D.1, B.1)	sim(D.1, C.2)	sim(D.1, D.2)	sim(D.1, A.3n)	sim(D.1, B.3n)	sim(D.1, C.3n)	sim(D.1, D.3n)

sim: similarity function among text embeddings

### COMPUTATION OF CATEGORICAL CROSS-ENTROPY LOSS

Predicted probability distribution

softmax(

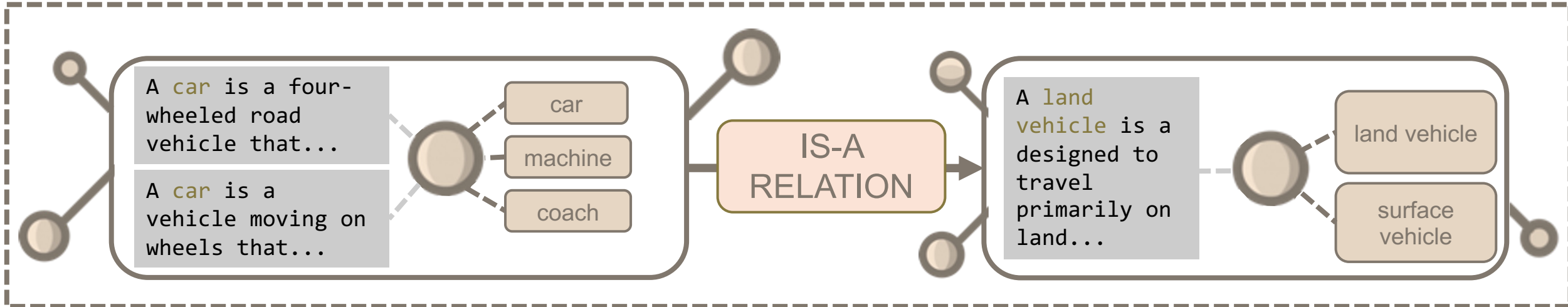
sim(A.1, A.2)	sim(A.1, B.2)	sim(A.1, C.2)	sim(A.1, D.2)	sim(A.1, A.3n)	sim(A.1, B.3n)	sim(A.1, C.3n)	sim(A.1, D.3n)
sim(B.1, A.2)	sim(B.1, B.2)	sim(B.1, C.2)	sim(B.1, D.2)	sim(B.1, A.3n)	sim(B.1, B.3n)	sim(B.1, C.3n)	sim(B.1, D.3n)
sim(C.1, A.2)	sim(C.1, B.2)	sim(C.1, C.2)	sim(C.1, D.2)	sim(C.1, A.3n)	sim(C.1, B.3n)	sim(C.1, C.3n)	sim(C.1, D.3n)
sim(D.1, A.2)	sim(D.1, B.1)	sim(D.1, C.2)	sim(D.1, D.2)	sim(D.1, A.3n)	sim(D.1, B.3n)	sim(D.1, C.3n)	sim(D.1, D.3n)

Expected probability distribution

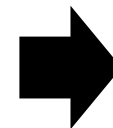
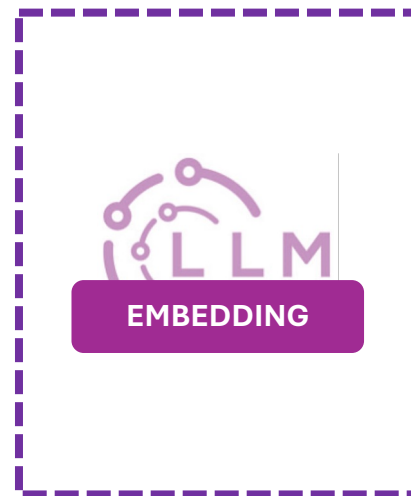
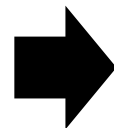
1	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0
0	0	1	0	0	0	0	0
0	0	0	1	0	0	0	0

# The ontological knowledge infusion approach

## STEP 1: Choosing the source ontology and the target embedding-Large Language Model



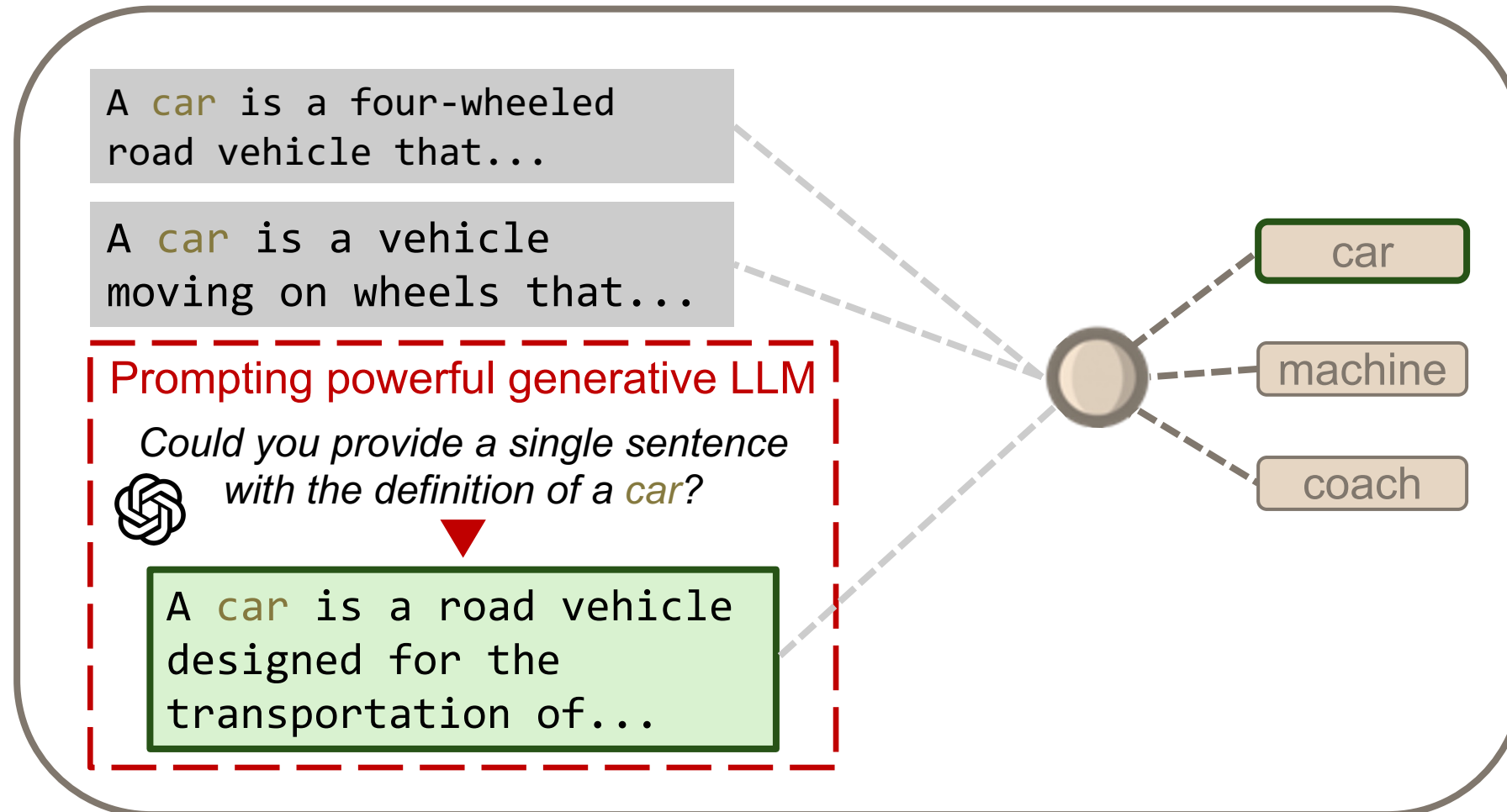
Input text excerpt



Text embedding

# The ontological knowledge infusion approach

## STEP 2: Generating **synthetic definitions of concepts**

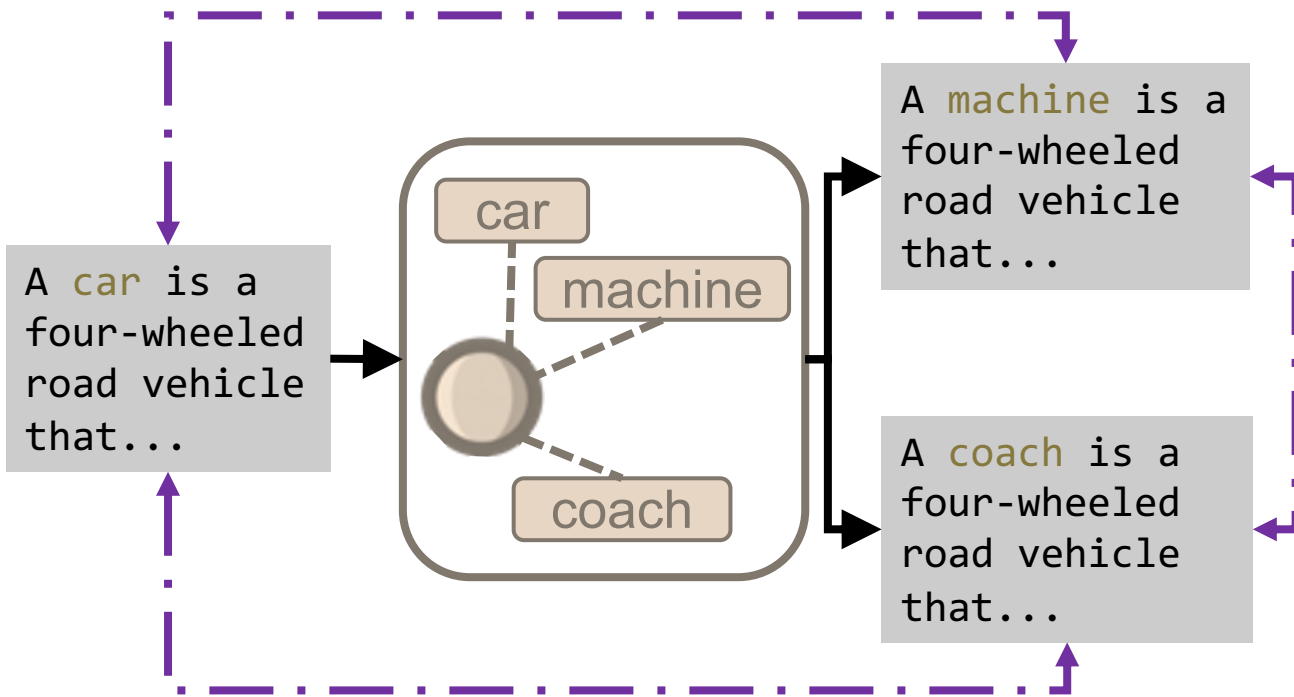




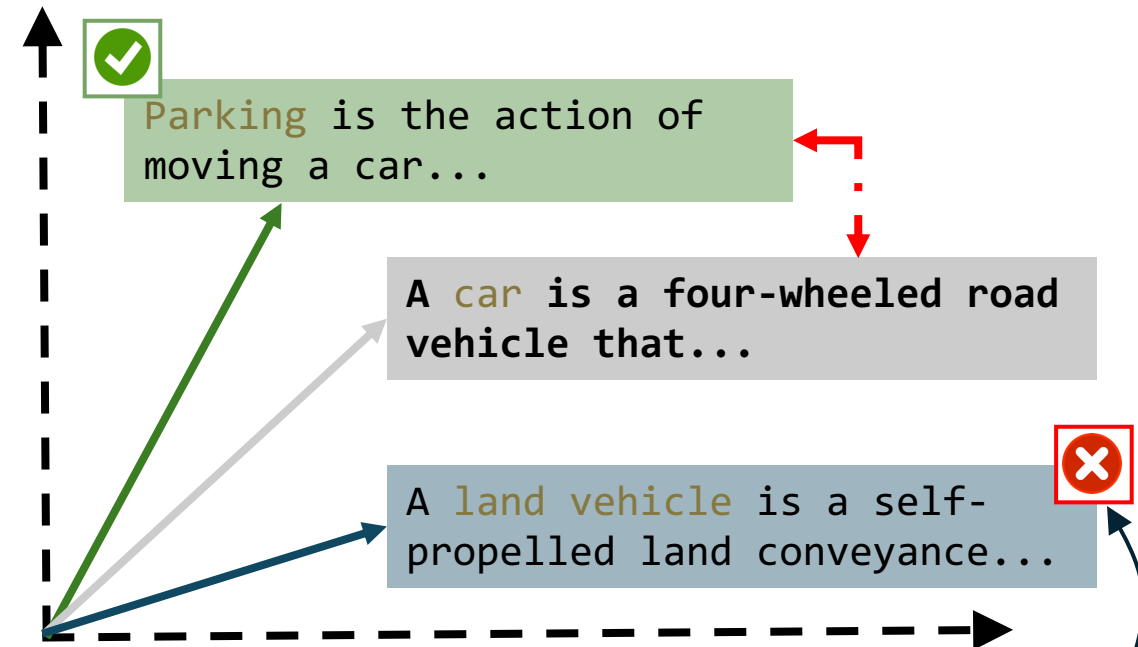
# The ontological knowledge infusion approach

**STEP 3:** Selecting **positive** and **negative** pairs of concept definitions, driven by the ontology

**POSITIVE PAIR** of semantically related texts



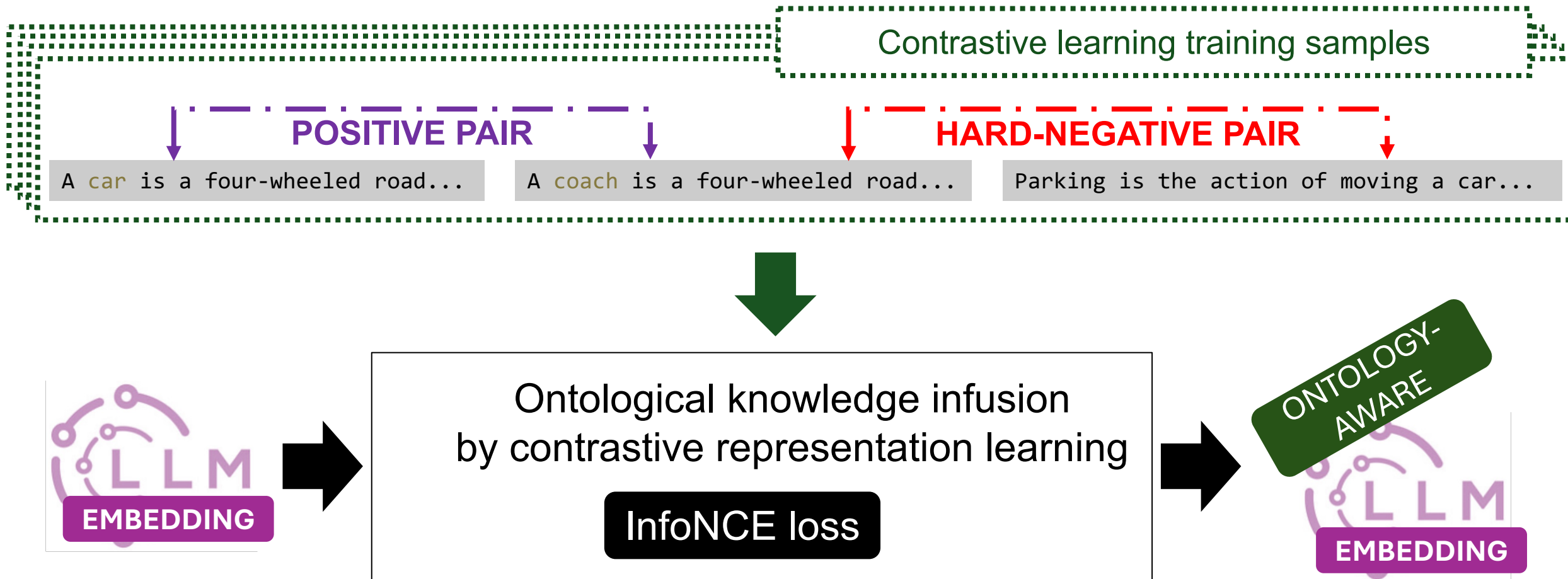
**HARD-NEGATIVE PAIR** of texts



Do not pair definitions of ancestor or descendant concepts

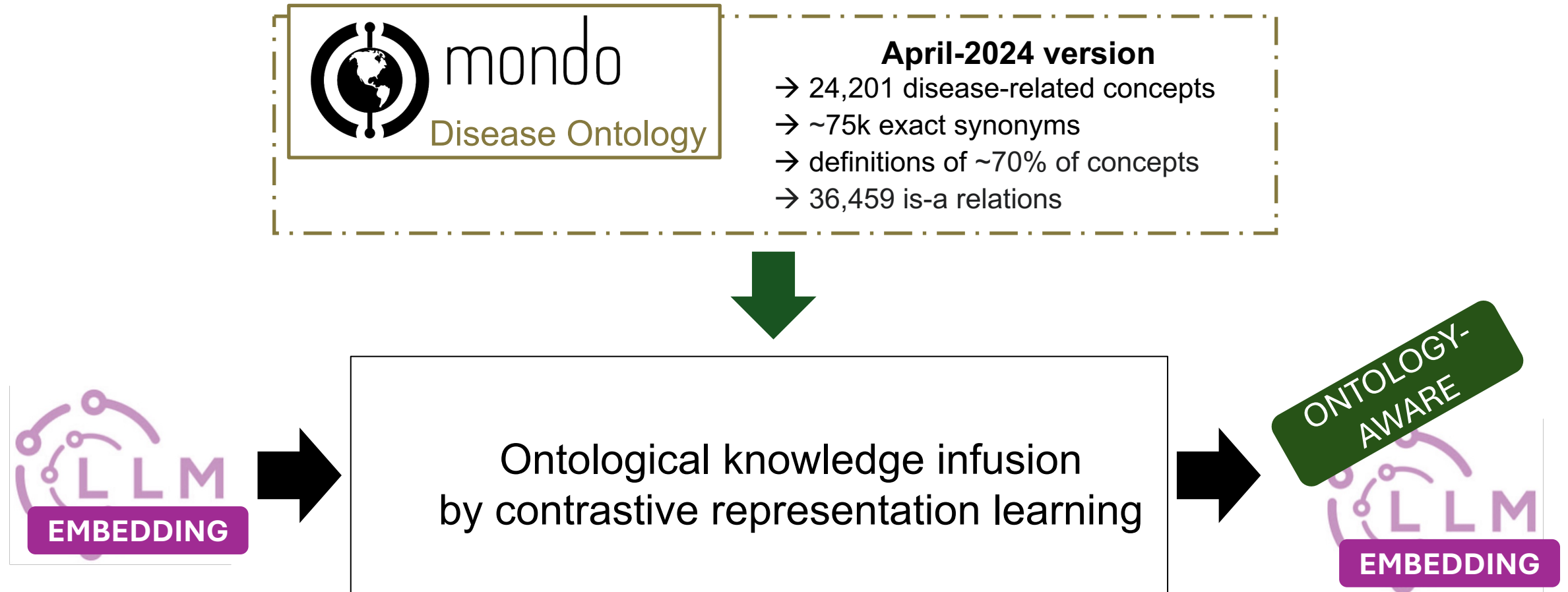
# The ontological knowledge infusion approach

**STEP 4:** Fine-tuning the embedding-Large Language Model by **contrastive representation learning**



## EVALUATION

Do embedding-Large Language Models better understand diseases, after infusing the MONDO disease ontology?



# EVALUATION

## Do embedding-Large Language Models better understand diseases, after infusing the MONDO disease ontology?



publicly available on



**PubMedBERT:** 110M parameters, pre-trained from scratch using *abstracts from PubMed and full-text articles from PubMedCentral* with masked language modelling and next sentence prediction

**SapBERT:** 110M parameters, fine-tuned in a *contrastive learning framework* to increase the similarity of *pairs of synonyms of biomedical concepts*, from the UMLS meta-thesaurus

**GTEbase:** 110M parameters, fine-tuned by means of a *two-stages contrastive learning framework*: pre-training text pairs by weak supervision, a subsequent training on higher-quality annotated datasets

**GIST:** 100M parameters, *one of the best performing small embedding-LLMs in the MTEB leader-board*, fine-tuned by a *contrastive objective* relying on an dynamic selection strategy to identify in-batch negative samples

Specialized to the biomedical domain

Not specialized to a specific knowledge domain

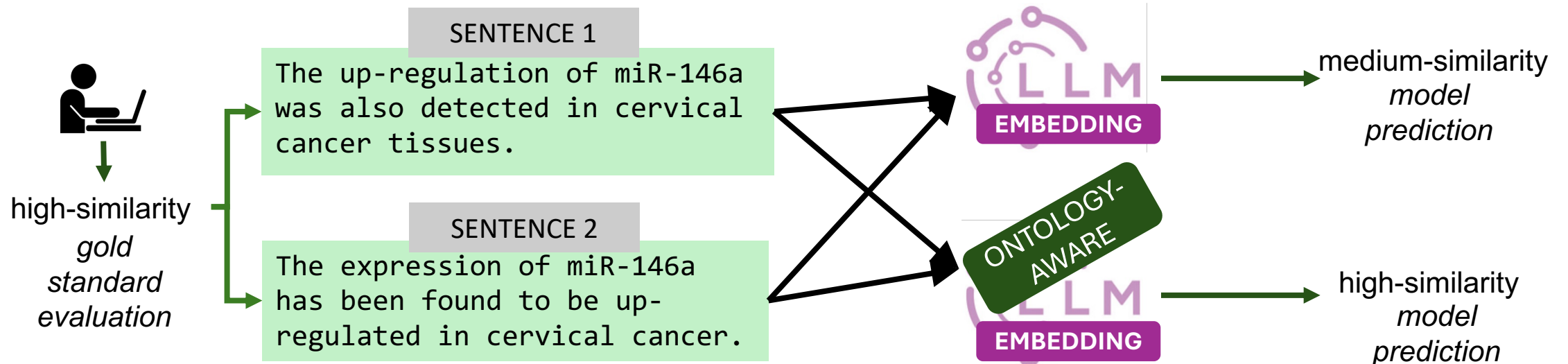
Pre-trained with standard objectives

Fine-tuned by contrastive objectives

## EVALUATION

Do embedding-Large Language Models better understand diseases, after infusing the MONDO disease ontology?

**TASK:** Sentence Similarity



**DATASETS:** **BIOSSES:** 100 sentence pairs from biomedical publications

**SEMEVAL:** about 12k pairs of sentences dealing with several distinct knowledge domains

## EVALUATION

Do embedding-Large Language Models better understand diseases, after infusing the MONDO disease ontology?

Embedding-Large Language Model		BIOSSES		STS 12		STS 13		STS 14		STS 15		STS 16	
		All	Dis	All	Dis	All	Dis	All	Dis	All	Dis	All	Dis
PubMedBERT	ORIGINAL	53.74	69.80	25.99	46.34	28.09	16.21	25.80	00.30	37.33	21.31	47.99	<b>80.33</b>
	ONTOLOGY-AWARE	<b>71.23</b>	<b>77.41</b>	<b>41.90</b>	<b>47.83</b>	<b>42.19</b>	<b>18.30</b>	<b>37.94</b>	<b>12.32</b>	<b>49.17</b>	<b>23.55</b>	<b>58.37</b>	72.78
SapBERT	ORIGINAL	81.86	83.21	70.89	68.84	79.23	35.73	70.37	47.64	77.85	56.99	76.71	89.73
	ONTOLOGY-AWARE	<b>85.45</b>	<b>84.79</b>	<b>72.31</b>	<b>79.99</b>	<b>80.66</b>	<b>46.04</b>	<b>72.44</b>	<b>52.07</b>	<b>79.79</b>	<b>64.05</b>	<b>77.58</b>	<b>92.86</b>
GTEbase	ORIGINAL	87.26	<b>90.30</b>	75.70	69.85	85.72	87.91	81.51	76.66	88.81	87.40	83.82	93.60
	ONTOLOGY-AWARE	<b>87.40</b>	89.62	<b>76.44</b>	<b>70.17</b>	<b>86.12</b>	<b>88.15</b>	<b>81.94</b>	<b>77.69</b>	<b>88.86</b>	<b>88.18</b>	<b>84.21</b>	<b>94.71</b>
GIST	ORIGINAL	87.96	89.66	76.15	63.88	87.85	88.64	83.39	74.52	89.43	<b>85.75</b>	85.35	<b>93.78</b>
	ONTOLOGY-AWARE	<b>88.86</b>	<b>92.05</b>	<b>76.69</b>	<b>65.94</b>	<b>87.99</b>	<b>89.26</b>	<b>83.64</b>	<b>75.45</b>	<b>89.56</b>	85.42	<b>85.69</b>	<b>93.78</b>

BIOMEDICAL  
SENTENCE SIMILARITY  
(in-domain)

NON-BIOMEDICAL  
SENTENCE SIMILARITY  
(out-of-domain)

**All:** all pairs of sentences considered

**Dis:** pairs of sentences mentioning diseases





arXiv > cs > arXiv:2405.20527

Computer Science > Computation and Language

[Submitted on 30 May 2024]

# Towards Ontology-Enhanced Representation Learning for Large Language Models

Taking advantage of the widespread use of ontologies to organise and harmonize knowledge across several distinct domains, this paper proposes a novel approach to improve an embedding-Large Language Model (embedding-LLM) of interest by infusing the knowledge formalized by a reference ontology: ontological knowledge infusion aims at boosting the ability of the considered LLM to effectively model the knowledge domain described by the infused ontology. The linguistic information (i.e. concept synonyms and descriptions) and structural information (i.e. is-a relations) formalized by the ontology are utilized to compile a comprehensive set of concept definitions, with the assistance of a powerful generative LLM (i.e. GPT-3.5-turbo). These concept definitions are then employed to fine-tune the target embedding-LLM using a contrastive learning framework. To demonstrate and evaluate the proposed approach, we utilize the biomedical disease ontology MONDO. The results show that embedding-LLMs enhanced by ontological disease knowledge exhibit an improved capability to effectively evaluate the similarity of in-domain sentences from biomedical documents mentioning diseases, without compromising their out-of-domain performance.

Available on ArXiv at: <https://arxiv.org/abs/2405.20527>

**IQVIA**

**Towards Ontology-Enhanced Representation Learning for Large Language Models**

*Francesco Ranzano and Jay Navavati – IQVIA Advanced NLP Team*

---

**OVERVIEW**

- **Ontologies** are extensively used to **organize and harmonize information** inside and across a wide range of domains and applications
- We propose a novel, automated approach to improve an embedding-Large Language Model of interest by infusing the linguistic (concept labels and definitions) and structural (taxonomic relations) knowledge formalized by a reference ontology, with the assistance of a powerful generative LLM
- Ontological knowledge infusion **boosts the ability of the considered embedding-Large Language Model to effectively deal with the knowledge domain described by the infused ontology, without compromising out-of-domain performance**

---

**ONTOLOGICAL KNOWLEDGE INFUSION APPROACH**

**STEP 1** Choice of an ontology of interest to be infused in a target embedding-Large Language Model

**STEP 2** Generation of **synthetic concept definitions** by prompting a powerful generative LLM

**STEP 3** Creation of pairs of **similar texts** by ontology-driven synonym substitution

**STEP 4** Ontology-driven selection of **hard-negative definitions** associated to each concept

**STEP 5** Fine-tuning of the embedding-Large Language Model by contrastive representation learning

---

**EXPERIMENTAL EVALUATION**

We rely on widespread **Sentence Similarity** benchmarks to evaluate the quality of text embeddings generated by four ~100M parameters embedding-Large Language Models before and after infusing the knowledge formalized by the MONDO Disease Ontology

Embedding-Large Language Model	BIOSSER		STS 12		STS 13		STS 14		STS 15		STS 16		
	All	Dis	All	Dis	All	Dis	All	Dis	All	Dis	All	Dis	
PubMedBERT	ORIGINAL	53.74	69.80	25.99	46.34	28.09	16.21	25.90	00.30	37.33	21.31	47.09	80.33
	ONTOLOGY-AWARE	71.23	77.41	41.90	47.83	42.59	18.30	37.94	12.32	49.17	23.55	58.37	72.78
SapBERT	ORIGINAL	81.86	83.21	70.89	68.84	79.23	35.73	70.37	47.64	77.85	56.90	76.71	89.73
	ONTOLOGY-AWARE	85.45	84.79	72.31	79.59	80.66	46.04	72.44	52.67	79.79	64.05	77.58	92.84
GTEbase	ORIGINAL	67.26	90.30	75.70	69.85	85.72	87.91	81.51	76.66	88.81	87.40	83.82	93.60
	ONTOLOGY-AWARE	87.40	89.62	76.44	79.17	86.12	88.15	81.94	77.89	88.86	88.18	84.21	94.71
GIST	ORIGINAL	87.96	89.66	76.15	63.88	87.85	88.64	83.30	74.52	89.43	85.78	85.35	93.78
	ONTOLOGY-AWARE	88.88	92.05	76.69	65.94	87.59	89.26	83.64	75.41	89.86	85.42	85.69	93.78

All pairs of sentences considered  
Dis: pairs of sentences mentioning diseases

BIOMEDICAL SENTENCE SIMILARITY (in-domain)
NON-BIOMEDICAL SENTENCE SIMILARITY (out-of-domain)

Data-centric Machine Learning Research Workshop @ International Conference on Machine Learning 2024  
© 2024. All rights reserved. IQVIA is a registered trademark of IQVIA Inc. in the United States, the European Union, and various other countries.

# Key learnings

- **Ontologies** are extensively used to **organize and harmonize information** across distinct domains and applications
- Knowledge resources like **ontologies** can be effectively exploited to both **create** and **effectively exploit high-quality textual data** useful to train Large Language Models, in data-hungry scenarios
- **Ontologies** can be effectively used to **prompt (powerful) generative Large Language Models** to drive the focused creation of **additional textual data** to support Large Language Models training
- **Contrastive learning** constitutes an effective technique to enrich the **latent knowledge** embedded inside a *Large Language Model* by relying on the **explicitly knowledge** formalized by an *ontology*

## Tools exploited



### Sentence Transformers

- Python library useful to access, use, and train text and image embedding models
- Provide customizable classis useful to train (batching, loss function, etc.) and evaluate sentence embeddings



### Hugging Face

- Used as a repository of embedding Large Language Models



### OpenAI GPT-3.5-Turbo

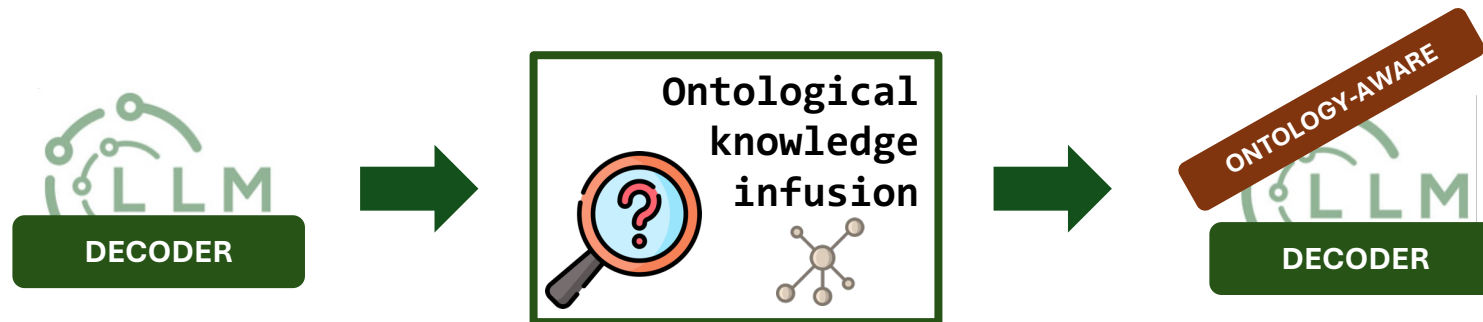
- Prompted to generate synthetic definitions of concepts from ontologies



# Next steps

- Evaluate **bigger embedding-Large Language Models**, eventually derived from generative models (e.g. by LLM2Vec)
- Consider **distinct / multiple ontologies** to quantify the effectiveness of ontological knowledge infusion under distinct scenarios
- Explore **alternative strategies for ontology-driven training data generation**
- **Extend evaluation to additional tasks**, besides sentence similarity

...and **extend the proposed ontological knowledge infusion approach to generative Large Language Models**





**Thanks for your attention**  
**Any questions?**

**Towards Ontology-Enhanced  
Representation Learning  
for Large Language Models**

Francesco Ronzano and Jay Nanavati  
IQVIA Advanced NLP Team