# Towards Ontology-Enhanced Representation Learning for Large Language Models

**Francesco Ronzano**          FRANCESCO.RONZANO@IQVIA.COM
*IQVIA*

**Jay Nanavati**          JAY.NANAVATI@IQVIA.COM
*IQVIA*

## Abstract

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.

**Keywords:** knowledge infusion, biomedical ontology, large language model, contrastive learning, representation learning, sentence similarity

## 1 Introduction

The availability of high-quality textual data is essential to boost the ability of LLMs to effectively understand, model and reason about the semantics of a text. With this respect, recently, *Synthetic Textual Data Augmentation* methods have been gaining an increasing relevance: very large, pre-trained LLMs have been often exploited to generate, restructure or annotate textual data that in turns is exploited to enhance smaller, domain specific LLMs (Ding et al., 2024; Tan et al., 2024). Besides synthetic data, the structured information included in a diverse set of knowledge resources has also been used to boost LLMs giving rise to several proposals of *Knowledge-resource Driven LLM-enhancement* techniques (Hu et al., 2023): for example, knowledge graphs represent the most relevant type of knowledge resources exploited to this purpose (Yang et al., 2024; Yasunaga et al., 2022).

Knowledge resources like ontologies are extensively used to organise and harmonize information inside and across a wide range of distinct domains and applications (Patel and Debnath, 2024). Various methods have been proposed to improve machine learning model performance by relying on ontologies and vice-versa (Kulmanov et al., 2020): recent examples include ontology-driven interaction with and fine-tuning of generative LLMs (Palagin et al., 2023; Baldazzi et al., 2023), ontology-based refinement of knowledge graph queries (Allemang and Sequeda, 2024) and exploitation of LLMs to create or enrich ontologies (Ciatto et al., 2024; Mateiu and Groza, 2023).

In this stream of works, the main contribution proposed and evaluated by this paper is **a novel, automated approach to infuse external knowledge, formalized by an ontology of interest, into an embedding-LLM (i.e. text encoder)**[1]. This is achieved by leveraging on both: (i) the linguistic and structural information formalized by an ontology (*Knowledge-resource Driven LLM-enhancement*) and (ii) a powerful generative LLM (i.e. GPT-3.5-turbo) to perform *Synthetic Textual Data Augmentation*. By using the generative LLM, a rich set of real and/or synthetic definitions are gathered for all the concepts specified by the considered ontology. These definitions are then exploited to create training samples useful to fine-tune a the target embedding-LLM by a contrastive learning framework: training samples (i.e. pairs of similar and dissimilar definitions) are generated by following a principled approach, aimed at maximizing their effectiveness for fine-tuning. Once fine-tuning finalizes, the vectorial representations of texts generated by the embedding-LLM will incorporate the knowledge formalized by the considered ontology.

## 2 Related work

Recently, several contrastive representation learning approaches have been proposed to improve the quality of text embeddings by exploiting collections of pairs of related or similar texts (e.g. query-answer, texts conveying the same meaning, etc.) to fine tune embedding-LLMs: the quality of embeddings is improved by increasing the similarity of the LLM-generated vectorial representations of semantically close texts (Hadsell et al., 2006). Examples of contrastive learning frameworks include Sentence-BERT (Reimers and Gurevych, 2019) where a dual-encoder network architecture, coupled with multiple loss functions is used to fine-tune text embeddings in a supervised way. SimCSE (Gao et al., 2021) proposes to use distinct LLM dropout masks as data augmentation strategies to generate pairs of similar emeddings for unsupervised fine-tuning of embedding-LLMs in a contrastive objective. Schick and Schütze (2021) use generative LLMs to create labelled text pairs useful for unsupervised fine-tuning of embedding-LLMs. Also Wang et al. (2023) relies on generative LLMs (i.e. GPT-3.5-turbo and GPT-4) to generate synthetic training data spanning over multiple tasks and languages: this data is then exploited to fine-tune Mistral-7b, a decoder-only LLM to generate better emeddings. Su et al. (2022) extend the text excerpt to be embedded with free-text instructions describing the task the embedding will be used for. Overall, a common paradigm exploited to fine-tune embedding-LLMs by contrastive learning relies on a contrastive pre-training phase that exploits collections of text pairs generated semi-automatically by weak supervision, followed by a fine-tuning phase where LLMs are improved by relying on higher-quality annotated datasets (Li et al., 2023; Wang

---

1. Implementation available on GitHub repository: `https://github.com/iqvianlp/llm-onto-infuse/`.

et al., 2022). Besides data augmentation approaches, two additional key ingredients useful to boost the effectiveness of contrastive learning frameworks exploited to generate better text embeddings are (i) the strategy to select of text pairs that will constitute training samples and (ii) the choice of the training objective (i.e. loss function) (Wang and Dou, 2023). In Liu et al. (2020), synonyms of biomedical concepts retrieved from the UMLS metathesaurus (Bodenreider, 2004) are exploited to fine-tune embedding-LLMs by contrastive learning. In comparison, our ontological knowledge infusion approach aims at enhancing embedding-LLMs by: (i) taking advantage of a richer set of linguistic and structural features of ontologies (beyond synonymy); (ii) exploiting the text of whole sentences (instead of noun phrases) to fine-tune embedding-LLMs; (iii) proposing a novel, automated approach to create training text pairs.

## 3 Workflow to infuse ontological knowledge in embedding-LLMs

The proposed approach to infuse ontological knowledge in embedding-LLMs relies on linguistic features - *synonym terms* and *definitions* - and structural features - *taxonomic relations* - that characterize the set of concepts defined by an ontology. As better detailed in the next Sections, these features, shared by most ontologies, are exploited to support and drive: (i) the *generation of text excerpts describing the concepts of the ontology* and (ii) the *effective aggregation of these text excerpts into pairs of similar or dissimilar ones*, to be exploited to fine-tune embedding-LLMs.

### 3.1 Fine-tuning embedding-LLMs by contrastive learning

The the textual information (i.e. definitions of ontological concepts) gathered by relying on a reference ontology is infused in an embedding-LLMs of choice by fine-tuning such LLM in a contrsative objective.

#### 3.1.1 Contrastive learning architecture and training objective

In our experiments we rely on the contrastive learning framework described by Chen et al. (2020) to perform ontological knowledge infusion. In principle, we can infuse ontological knowledge in any embedding-LLM ($EMB$) capable of generating, given a text $t$, the corresponding dense vectorial representation $h_t = EMB(t) \in R^n$, where $n$ represents the dimension of $h_t$, the text embedding[2]. Let suppose to have at our disposal a collection of $I$ pairs of semantically related texts $(t_i, t_i^+)$, where $0 < i <= I$. Considering batches of $N$ pairs of semantically related texts, the corresponding pairs of embeddings $(h_i, h_i^+)$ can be computed by relying on the LLM ($EMB$). The categorical cross-entropy loss is exploited to favour, for each embedding $h_i$, the identification of (i.e. prediction of the class associated to) the associated positive embedding $h_i^+$: samples from other embedding pairs in the same batch are considered as noise. This training objective, referred to as **InfoNCE loss** (Oord

---

2. Depending on the specific scenario and embedding-LLM considered, the embedding $h_t$ can be generated by distinct strategies, including pooling of single-token embeddings or by considering special-purpose tokens of the embedding-LLM (e.g. the CLS token).

et al., 2018), is described by the following formula:

$$loss_i = -log \frac{e^{sim(h_i,h_i^+)/\tau}}{\sum_{i=1}^{N} e^{sim(h_i,h_j^+)/\tau}} \tag{1}$$

where $N$ is the batch size and $\tau$ is the temperature. $sim(p,q)$ is the similarity function between the embeddings $p$ and $q$: it is common practice to use cosine similarity.

It has been shown that when InfoNCE loss is exploited, the quality of learned embeddings improves sensibly if in each batch, for each pair of semantically related texts $(t_i, t_i^+)$, one (or more) **hard negative texts** are included (Chen et al., 2017; Gao et al., 2021). Given a pair of positive texts, an associated hard negative sample $w$ is a text that is semantically distinct from the texts $t_i$ and $t_i^+$, even if its embedding $h_w = EMB(w)$ is characterized by a high semantic similarity with the embeddings of any positive text (i.e. $h_i$ and $h_i^+$). Therefore, if for each positive text pair one or more hard negative texts are selected, each training sample exploited by the considered contrastive learning framework and training objective would be represented by the tuple of texts $(t_i, t_i^+, w_i^{HN_1}, ..., w_i^{HN_K},)$ where $t_i$ and $t_i^+$ represent the pair of positive texts, while $w_i^{HN_k}$, with $0 < i <= K$, are the associated $K$ hard negative samples[3].

### 3.1.2 ONTOLOGY-DRIVEN CREATION OF TRAINING SAMPLES

This Section describes the novel procedure we devise to create training samples useful to infuse ontological knowledge into an embedding-LLM of our choice, in a contrastive objective. After generating synthetic definitions of the concepts included in the considered ontology by prompting a generative LLM, we exploit these definitions in order to create training samples, thus selecting positive text pairs as well as associated hard negative texts.

**Prompting LLMs to generate synthetic concept definitions**: our ontological knowledge infusion approach is based on the availability of textual contents describing the concepts formalized by the ontology of choice: in the current setting, we focus on concept definitions[4]. Ontologies could include definitions of (part of) their concepts. To guarantee the availability of at least one definition associated to each concept, generative LLMs are prompted to create synthetic definitions: for each synonym of an ontology concept, a one-sentence synthetic definition of that concept is collected (see part (a) of Figure 1). The structure of the definition-generation prompt is highly dependent on both the features and domain of the considered ontology and the generative LLM of choice.

**Creation of semantically related pairs of texts by synonym substitution**: as illustrated in Section 3.1.1, the contrastive learning framework exploited to infuse ontological knowledge in embedding-LLMs relies on training samples constituted by pairs of semantically related texts. To create such text pairs, we rely on both the definitions (real and synthetic ones) and the synonyms of the concepts of the ontology (see part (b) of Figure 1).

For each concept, we select all its real or synthetic definitions that mention one and only one of its synonyms. Then for each selected definition, we generate similar definitions by replacing the mentioned synonym with a distinct synonym of the same concept. Therefore,

---

3. If $K$ is zero, hard negative sampling is not exploited.
4. As future work, we would like to evaluate the knowledge infusion effectiveness of other types of textual information associated to concepts, distinct from definitions.
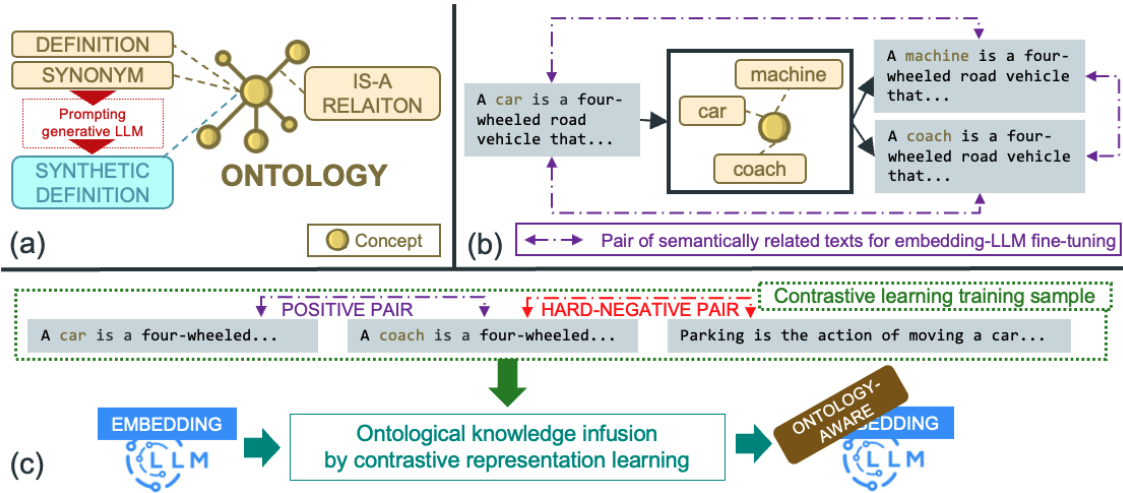
Figure 1: (a) structure of ontologies, generation of synthetic concept definitions; (b) creation of pairs os semantically related sentences by synonym substitution; (c) overview of the ontological knowledge infusion approach.

in our contrastive learning settings, we consider as pairs of semantically related texts (i.e. training samples) all possible pairs of definitions of the same concept, one obtained by performing synonym substitution over the text of the paired definition[5]. This approach would not allow the generation of any positive text pair for concepts with a single synonym: to remedy this, we automatically generate an additional, novel synthetic synonym for such concepts, by appending one of the synonyms of its parent concept to its actual synonym.

**Ontology-driven selection of hard negative samples**: we rely on both embedding similarity and the taxonomic concept-to-concept relations specified by the ontology to select one hard negative text associated to each pair of semantically related texts (i.e. each pair of definitions of the same concept; see part (c) of Figure 1). Given the positive text pair $(t_i, t_i^+)$ including two definitions of the concept $C_t$, the associated hard negative text $w$ should verify the following conditions:

- considering the taxonomic relations defined by the ontology, the concept $C_w$ should not be an ancestor or a descendant of the concept $C_t$;

- the embedding of the definition $w$ (i.e. $h_w$) should represent, among all the embeddings of definitions of concepts distinct from $C_t$, the one with the highest semantic similarity with the embeddings of the texts $t_i$ and $t_i^+$ (i.e. $h_{t_i}$ and $h_{t_i^+}$)[6].

---

5. In our experiments, we observed that the exploitation of training samples that include pairs of definitions not obtained by synonym substitution - e.g. (real definition, synthetic definition) - hinders the performance of our ontological knowledge infusion approach, probably since both definitions would highlight distinct traits of the same concept, characterising the same concept through distinct perspectives.

6. Similarity scores are averaged across $t_i$ and $t_i^+$; we use cosine similarity in our experiments.

## 4 Infusing disease knowledge relying on MONDO ontology

We showcase and evaluate our ontological knowledge infusion approach by infusing the disease knowledge formalized by a widespread and rich biomedical ontology into four distinct flavours of embedding-LLMs. More specifically, we consider MONDO[7] (Vasilevsky et al., 2022), an ontology that aims at globally harmonizing the characterization of diseases by unifying the information contained in multiple knowledge resources and data models. We used the April-2024 version of MONDO that defines 24,201 disease-related concepts characterized by almost 75 thousands synonyms[8]. MONDO includes a definition of almost 70% of its concepts. Concepts are hierarchically organized through 36,459 is-a relations.

We prompt GPT-3.5-turbo to generate a single-sentence definition from each synonym of concepts from the MONDO ontology (see Appendix A): 57,692 synthetic concept definitions are thus generated[9], with at least one definition of each concept.

By relying on the synonym substitution procedure described in Section 3.1.2, considering both ontology-provided and synthetic concept definitions, about 400 thousands pairs of semantically related definitions are generated. Each one of these pairs of definitions is extended with an hard negative definition selected by following the procedure explained in the latest part of Section 3.1.2: as a result, the ontology-driven creation of 400 thousand training samples ready to infuse the disease knowledge formalized by MONDO ontology in any embedding-LLM of choice is finalized.

We exploit the contrastive learning framework described in Section 3.1.2 to infuse ontological knowledge in the following four embedding-LLMs:

- **PubMedBERT** (Gu et al., 2021): a BERT-like LLM (110M parameters), pre-trained from scratch using abstracts from PubMed and full-text articles from PubMedCentral through masked token prediction and next sentence prediction;

- **SapBERT** (Liu et al., 2020): relies on PubMedBERT as base model, is fine-tuned in a contrastive learning framework to increase the similarity of pairs of synonyms of biomedical concepts, retrieved from the UMLS meta-thesaurus;

- **GTEbase** Li et al. (2023): an encoder LLM (110M parameters) created fine-tuning BERT-base-uncased by means of a two-stages contrastive learning framework: a pre-training phase relies on collections of text pairs generated semi-automatically by weak supervision. A subsequent training phase improves the LLM by higher-quality annotated datasets;

- **GIST** (Solatorio, 2024): at the time of writing, represents one of the best performing small embedding-LLMs (about 100M) in the MTEB leader-board[10], fine-tuned by a contrastive objective relying on an dynamic selection strategy to identify in-batch negative samples.

Each one of these embedding-LLMs, all having a comparable number of parameters, is characterized by specific peculiarities that make it interesting with respect to our evaluations: **PubMedBERT** is not fine-tuned exploiting a contrastive objective and is pre-trained on texts from the same domain of the MONDO ontology; **SapBERT** is fine-tuned

---

7. https://mondo.monarchinitiative.org/
8. The main name of each concept together with its EXACT synonyms are considered; obsolete concepts have been ignored.
9. This number is lower than the actual number of MONDO concept synonyms (about 75 thousands) since we applied specific synonym filtering rules resumed in Appendix D.
10. https://huggingface.co/spaces/mteb/leaderboard

on biomedical-synonym knowledge encoded in the UMLS meta-thesaurus; **GTEbase** is fine-tuned on a huge set of annotated datasets, including biomedical ones; **GIST** is a robust LLM, highly-scoring in widespread embedding benchmarks. Appendix B provides a detailed description of the LLM fine-tuning settings.

## 4.1 Evaluation: datasets and results

We evaluate the quality of the embeddings generated by the four LLMs previously introduced, before and after infusing knowledge from the MONDO ontology. We choose sentence similarity (STS) to perform our evaluation since this task is commonly extensively used to measure the quality of text embeddings. We considered the following STS annotated datasets: **BIOSSES** (Soğancıoğlu et al., 2017) including 100 sentence pairs from biomedical publications and **the five test sets of SemEval Sentence Similarity challenges** released yearly from 2012 to 2016 (Agirre et al., 2012, 2013, 2014, 2015, 2016). Spearman's correlation is exploited to evaluate the performance of the four embedding-LLMs. Table 1 and Table 2 show BIOSSES and SemEval results, respectively. Appendix C provides detailed information on both datasets.

| Embedding LLM | BIOSSES | |
|---|---|---|
| | All | Dis |
| $PubMedBERT_{orig}$ | 53.74 | 69.80 |
| $PubMedBERT_{kinf}$ | **71.23** | **77.41** |
| $SapBERT_{orig}$ | 81.86 | 83.21 |
| $SapBERT_{kinf}$ | **85.45** | **84.79** |
| $GTEbase_{orig}$ | 87.26 | **90.30** |
| $GTEbase_{kinf}$ | **87.40** | 89.62 |
| $GIST_{orig}$ | 87.96 | 89.66 |
| $GIST_{kinf}$ | **88.86** | **92.05** |

Table 1: BIOSSES STS Spearman correlation scores, before ($orig$) and after ($kinf$) infusing disease-related ontological knowledge in embedding-LLMs. Evaluation scores computed considering: whole dataset (All), sentences with disease mentions (Dis).

## 4.2 Discussion

Our ontological knowledge infusion approach consistently improves the sentence similarity performance of embedding-LLMs across a wide range of evaluation datasets, as demonstrated in Table 1 and Table 2. By infusing disease knowledge, this novel method enhances both domain-specific (i.e. customized to effectively model biomedical texts) and general-purpose embedding-LLMs, with a more pronounced impact on domain-specific models like SapBERT. The infusion of domain knowledge improves *in-domain* sentence similarity performance of embedding-LLMs, evaluated against biomedical sentence pairs (Table 1 and 'Dis' columns of Table 2): at the same time, ontological knowledge infusion does not deteriorate the capability of embedding-LLMs to effectively evaluate the similarity of *out-of-domain* sentences, from domains distinct from biomedicine ('All' columns of Table 2).

| Embedding | STS12 | | STS13 | | STS14 | | STS15 | | STS16 | |
|---|---|---|---|---|---|---|---|---|---|---|
| LLM | All | Dis | All | Dis | All | Dis | All | Dis | All | Dis |
| $PMBERT_{orig}$ | 25.99 | 46.34 | 28.09 | 16.21 | 25.80 | 00.30 | 37.33 | 21.31 | 47.99 | **80.33** |
| $PMBERT_{kinf}$ | **41.90** | **47.83** | **42.19** | **18.30** | **37.94** | **12.32** | **49.17** | **23.55** | **58.37** | 72.78 |
| $SapBERT_{orig}$ | 70.89 | 68.84 | 79.23 | 35.73 | 70.37 | 47.64 | 77.85 | 56.99 | 76.71 | 89.73 |
| $SapBERT_{kinf}$ | **72.31** | **79.99** | **80.66** | **46.04** | **72.44** | **52.07** | **79.79** | **64.05** | **77.58** | **92.86** |
| $GTEbase_{orig}$ | 75.70 | 69.85 | 85.72 | 87.91 | 81.51 | 76.66 | 88.81 | 87.40 | 83.82 | 93.60 |
| $GTEbase_{kinf}$ | **76.44** | **70.17** | **86.12** | **88.15** | **81.94** | **77.69** | **88.86** | **88.18** | **84.21** | **94.71** |
| $GIST_{orig}$ | 76.15 | 63.88 | 87.85 | 88.64 | 83.39 | 74.52 | 89.43 | **85.75** | 85.35 | **93.78** |
| $GIST_{kinf}$ | **76.69** | **65.94** | **87.99** | **89.26** | **83.64** | **75.45** | **89.56** | 85.42 | **85.69** | **93.78** |

Table 2: SemEval STS Spearman correlation scores, before ($orig$) and after ($kinf$) infusing disease-related ontological knowledge in embedding-LLMs. Evaluation scores computed considering: whole dataset (All), sentences with disease mentions (Dis). PMBERT is the abbreviated name for PubMedBERT.

The four embedding-LLMs studied have comparable parameter counts, but the extent of improvement varies based on the pre-training and fine-tuning strategies they underwent before ontological knowledge infusion. LLMs built using more basic approaches, such as PubMedBERT, exhibit greater performance gains after ontological knowledge infusion. In contrast, more advanced LLMs like GTEbase and GIST show smaller but consistent improvements across most evaluation scenarios. SapBERT, an embedding-LLM pre-trained on biomedical data using a novel, synonymy-based pre-training approach, demonstrates strong baseline performance. However, our ontological knowledge infusion method further enhances its sentence similarity capabilities, highlighting the effectiveness of our approach in improving even state-of-the-art domain-specific models.

## 5 Conclusions and future work

In this paper we presented and evaluated a novel approach to infuse the knowledge formalized by ontologies in embedding-LLMs with the aim of improving the ability of such embedding-LLMs to effectively model the knowledge domain described by that ontology. We showcased the effectiveness of our approach by infusing the knowledge formalized by the disease ontology MONDO into four representative flavours of embedding-LLMs.

As future directions of research we would like to explore a wider range of scenarios and approaches useful to perform ontology-driven knowledge infusion in embedding-LLMs by: (i) comparing a wider range of LLM flavours, possibly with bigger sizes and distinct architectures; (ii) evaluating the effectiveness of ontological knowledge infusion when we consider ontologies describing distinct domains with different granularities; (iii) exploring alternative LLM-prompting strategies to generate textual data by relying on ontological knowledge; (iv) considering additional evaluation tasks, besides sentence similarity, thus making evaluations more comprehensive, spanning a wider set of usage scenarios of embedding-LLMs.

## Reproducibility statement

The paper provides the information needed to reproduce the proposed ontological knowledge infusion approach and the related evaluations, in particular:

- **Ontology-driven training sample creation procedure**: Section 3.1.2 provides the procedural description of the distinct steps that contribute to the creation of training data to support ontological knowledge infusion. Related to this, Appendix A specifies the prompt used to generate synthetic definitions of ontological concepts by relying on GPT-3.5-turbo;

- **LLM fine-tuning information**: Section 3.1.1 explains the fine-tuning approach we adopted, relying on a contrastive objective. Details concerning fine-tuning process and values of hyper-parameters are specified in Appendix B;

- **Evaluation datasets, metrics and procedure**: as described in Section 4.1, widespread public sentence similarity datasets (together with related standard evaluation metrics) are exploited to evaluate the effectiveness of ontology-driven knowledge infusion. More information on the contents of each evaluation dataset is provided in Appendix C.

To foster reproducibility of the results of this paper, the implementation of the proposed ontological knowledge infusion framework, together with the exploited evaluation procedures and references to evaluation datasets are available at `https://github.com/iqvianlp/llm-onto-infuse/`.

## Broader Impact Statement

This work presents a novel approach to infuse knowledge from ontologies into embedding-based Large Language Models (LLMs), improving their ability to model and reason about the domain described by the ontology. By leveraging linguistic and structural information in ontologies, the proposed method generates concept definitions used to fine-tune LLMs through contrastive learning. Potential positive impacts include enhancing LLMs' factual knowledge and reasoning capabilities in domains like biomedicine, law, and finance, supporting the development of reliable domain-specific applications. However, negative impacts may arise from amplifying biases and information unbalance, potentially in present in the ontologies. Care must be taken to use authoritative, regularly updated ontologies and to maintain human oversight, as LLMs can still make errors. This promising approach could positively impact knowledge work across specialized domains, but responsible deployment practices are crucial. Further research can explore knowledge infusion from diverse ontologies into various LLM architectures to fully realize the potential benefits.

## Acknowledgments and Disclosure of Funding

# References

Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. Semeval-2012 task 6: A pilot on semantic textual similarity.* sem 2012: The first joint conference on lexical and computational semantics—. In *Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012), Montréal, QC, Canada*, pages 7–8, 2012.

Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. * sem 2013 shared task: Semantic textual similarity. In *Second joint conference on lexical and computational semantics (* SEM), volume 1: proceedings of the Main conference and the shared task: semantic textual similarity*, pages 32–43, 2013.

Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. Semeval-2014 task 10: Multilingual semantic textual similarity. In *Proceedings of the 8th international workshop on semantic evaluation (SemEval 2014)*, pages 81–91, 2014.

Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Inigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, et al. Semeval-2015 task 2: Semantic textual similarity, english, spanish and pilot on interpretability. In *Proceedings of the 9th international workshop on semantic evaluation (SemEval 2015)*, pages 252–263, 2015.

Eneko Agirre, Carmen Banea, Daniel Cer, Mona Diab, Aitor Gonzalez Agirre, Rada Mihalcea, German Rigau Claramunt, and Janyce Wiebe. Semeval-2016 task 1: Semantic textual similarity, monolingual and cross-lingual evaluation. In *SemEval-2016. 10th International Workshop on Semantic Evaluation; 2016 Jun 16-17; San Diego, CA. Stroudsburg (PA): ACL; 2016. p. 497-511.* ACL (Association for Computational Linguistics), 2016.

Dean Allemang and Juan Sequeda. Increasing the llm accuracy for question answering: Ontologies to the rescue! *arXiv preprint arXiv:2405.11706*, 2024.

Teodoro Baldazzi, Luigi Bellomarini, Stefano Ceri, Andrea Colombo, Andrea Gentili, and Emanuel Sallinger. Fine-tuning large enterprise language models via ontological reasoning. In *International Joint Conference on Rules and Reasoning*, pages 86–94. Springer, 2023.

Olivier Bodenreider. The unified medical language system (umls): integrating biomedical terminology. *Nucleic acids research*, 32(suppl_1):D267–D270, 2004.

Ting Chen, Yizhou Sun, Yue Shi, and Liangjie Hong. On sampling strategies for neural network-based collaborative filtering. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 767–776, 2017.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.

Giovanni Ciatto, Andrea Agiollo, Matteo Magnini, and Andrea Omicini. Large language models as oracles for instantiating ontologies with domain-specific knowledge. *arXiv preprint arXiv:2404.04108*, 2024.

Bosheng Ding, Chengwei Qin, Ruochen Zhao, Tianze Luo, Xinze Li, Guizhen Chen, Wenhan Xia, Junjie Hu, Anh Tuan Luu, and Shafiq Joty. Data augmentation using llms: Data perspectives, learning paradigms and challenges. *arXiv preprint arXiv:2403.02990*, 2024.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021.

Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021.

Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006.

Linmei Hu, Zeyi Liu, Ziwang Zhao, Lei Hou, Liqiang Nie, and Juanzi Li. A survey of knowledge enhanced pre-trained language models. *IEEE Transactions on Knowledge and Data Engineering*, 2023.

Maxat Kulmanov, Fatima Zohra Smaili, Xin Gao, and Robert Hoehndorf. Machine learning with biomedical ontologies. *biorxiv*, pages 2020–05, 2020.

Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*, 2023.

Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. Self-alignment pretraining for biomedical entity representations. *arXiv preprint arXiv:2010.11784*, 2020.

Patricia Mateiu and Adrian Groza. Ontology engineering with large language models. *arXiv preprint arXiv:2307.16699*, 2023.

Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

Oleksandr Palagin, Vladislav Kaverinskiy, Anna Litvin, and Kyrylo Malakhov. Ontochatgpt information system: Ontology-driven structured prompts for chatgpt meta-learning. *arXiv preprint arXiv:2307.05082*, 2023.

Archana Patel and Narayan C Debnath. A comprehensive overview of ontology: Fundamental and research directions. *Current Materials Science: Formerly: Recent Patents on Materials Science*, 17(1):2–20, 2024.

Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.

Pedro Ruas and Francisco M Couto. Nilinker: attention-based approach to nil entity linking. *Journal of Biomedical Informatics*, 132:104137, 2022.

Timo Schick and Hinrich Schütze. Generating datasets with pretrained language models. *arXiv preprint arXiv:2104.07540*, 2021.

Gizem Soğancıoğlu, Hakime Öztürk, and Arzucan Özgür. Biosses: a semantic sentence similarity estimation system for the biomedical domain. *Bioinformatics*, 33(14):i49–i58, 2017.

Aivin V Solatorio. Gistembed: Guided in-sample selection of training negatives for text embedding fine-tuning. *arXiv preprint arXiv:2402.16829*, 2024.

Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen-tau Yih, Noah A Smith, Luke Zettlemoyer, and Tao Yu. One embedder, any task: Instruction-finetuned text embeddings. *arXiv preprint arXiv:2212.09741*, 2022.

Zhen Tan, Alimohammad Beigi, Song Wang, Ruocheng Guo, Amrita Bhattacharjee, Bohan Jiang, Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. Large language models for data annotation: A survey. *arXiv preprint arXiv:2402.13446*, 2024.

Nicole A Vasilevsky, Nicolas A Matentzoglu, Sabrina Toro, Joseph E Flack IV, Harshad Hegde, Deepak R Unni, Gioconda F Alyea, Joanna S Amberger, Larry Babb, James P Balhoff, et al. Mondo: Unifying diseases for the world, by the world. *medRxiv*, pages 2022–04, 2022.

Hao Wang and Yong Dou. Sncse: Contrastive learning for unsupervised sentence embedding with soft negative samples. In *International Conference on Intelligent Computing*, pages 419–431. Springer, 2023.

Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*, 2022.

Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368*, 2023.

Linyao Yang, Hongyang Chen, Zhao Li, Xiao Ding, and Xindong Wu. Give us the facts: Enhancing large language models with knowledge graphs for fact-aware language modeling. *IEEE Transactions on Knowledge and Data Engineering*, 2024.

Michihiro Yasunaga, Jure Leskovec, and Percy Liang. Linkbert: Pretraining language models with document links. *arXiv preprint arXiv:2203.15827*, 2022.

## Appendix A. GPT-3.5-turbo prompts to generate synthetic definitions

In order to generate synthetic definitions of concepts from the MONDO ontology, for each synonym of each concept (i.e. for each *MONDO_CONCEPT_SYNONYM*), GPT-3.5-turbo has been prompted by means of the following dialogue-prompt:

- `SYSTEM: You are an expert in clinical and biomedical sciences.`

- `USER: Could you provide a single sentence with the definition of` *MONDO_CONCEPT_SYNONYM*`?`

## Appendix B. Fine-tuning hyper-parameters

Ontological knowledge infusion has been performed by fine-tuning the considered embedding-LLMs by relying on the Sentence Transformers Python module available at `https://sbert.net/`. Eembedding-LLMs have been fine-tuned by relying on the InfoNCE loss (temperature $\tau$ equal to 0.05) with the following hyper-parameters:

- **batch size**: 24 training samples

- **learning rate**: 1e-8

- **learning rate scheduler**: constant learning rate after a warm-up period during which the learning rate increases linearly starting from 0 up to the set learning rate value (i.e. 1e-8)

- **weight decay**: 1e-4

Fine-tuning has been carried out for a maximum of two epochs: after the first epoch, the model with the best the Speareman correlation score computed against the whole BIOSSES dataset has been considered and fully evaluated. Cosine similarity has been consider in order to quantify the distance between pairs of text embeddings.

## Appendix C. Detailed information on evaluation datasets

Table 3 describes the size of the distinct Sentence Similarity evaluation datasets that have been exploited to evaluate the ontological knowledge infusion approach. For each dataset both the total number of sentence pairs and the number of sentence pairs where both sentences include one or more disease mentions is specified. Disease mentions have been spotted by relying on the disease Named Entity Recognition (NER) model introduced by Ruas and Couto (2022), available at `https://huggingface.co/pruas/BENT-PubMedBERT-NER-Disease`.

## Appendix D. Synonym filtering rules applied to MONDO ontology

We refine and filter the synonyms of concepts from MONDO ontology to be considered to generate synthetic concept definition by applying the following rules:

- if any, text in parenthesis is deleted from all the synonyms;

- duplicate synonyms are removed (case-insensitively);

| Dataset | Total number of sentence pairs (All) | Number of sentence pairs mentioning diseases (Dis) |
|---------|--------------------------------------|----------------------------------------------------|
| BIOSSES | 100 | 31 |
| STS12 | 3,108 | 34 |
| STS13 | 1,500 | 17 |
| STS14 | 3,750 | 59 |
| STS15 | 3,000 | 36 |
| STS16 | 1,186 | 15 |

Table 3: Number of sentence pairs included in the BIOSSES and SemEval Sentence Similarity datasets: total number of pairs (All) and number of pairs where both sentences include one or more disease mentions (Dis).

- for each concept, we consider just a single randomly-chosen synonym in case: (i) a pair of synonym is characterised by a string Levenshtein distance lower than 10 or (ii) a pair of synonyms differs just by word ordering (case- and punctuation-insensitively).