

Contents lists available at IJIECM International Journal of Industrial Engineering and Construction Management



Journal Homepage: http://www.ijiecm.com/ Volume 1, No. 1, 2024

Self-Supervised Contrastive Embeddings for Semantic Enrichment: Contextual Alignment and Entity Linking

Simona Vasilica, Adela Bara

Department of Economic Informatics, Bucharest University of Economic Studies, Bucharest, Romania

ARTICLE INFO

Received: 2024/02/10 Revised: 2024/05/21 Accepted: 2024/06/15

Keywords:

Semantic enrichment; contrastive learning; entity linking; knowledge graphs; contextual alignment; retrieval; candidate generation

ABSTRACT

Semantic enrichment benefits from representations that respect both textual context and knowledge-graph structure. Building on the bibliometric baseline of the field [1], we propose a self-supervised contrastive framework that learns sentence- and mention-level embeddings aligned across three natural positive signals: (i) co-mention and coreference within documents, (ii) adjacency in a knowledge graph, and (iii) co-citation/co-reference at the article level. Without manual labels, the method supports two downstream tasks central to enrichment—contextual candidate generation and entity linking. On three domains (ontology/linked data, biomedical, social streams) our approach improves candidate-recall@50 by 6–12% and end-to-end linking F1 by 3–6% over strong neural baselines. Ablations isolate the contributions of graph-positive sampling and adaptive temperature. We release scripts to reproduce figures (loss curves, PR curves, embedding scatter) and tables (dataset summary, ablations), designed to compile with this template.

1. Introduction

Semantic enrichment attaches machine-interpretable structure to natural language content—entities, relations, types, and links into knowledge graphs (KGs)—so that downstream systems can retrieve, analyze, and reason more effectively. The research area spans linked data publication, ontology engineering, semantic annotation, entity linking, and more recent neural approaches to representation learning. A comprehensive bibliometric snapshot of this field identified the most active venues, nations, and author clusters, and surfaced recurring keyword overlays and topic co-occurrence patterns [1]. That descriptive baseline motivates work that does not only catalog the literature but also operationalizes cross-domain generalization in enrichment systems.

Motivation. In practical pipelines, two stages dominate performance: candidate generation (quickly retrieving plausible entities for a mention) and entity linking (selecting the correct entity given its context). These stages are often trained with supervised signals that are expensive to obtain and brittle under distribution

shift [23]. At the same time, the enrichment ecosystem contains abundant *implicit* supervision: co-mentions and coreference within articles, semantic proximity in KGs [6, 7], and article-level ties (co-citation, co-reference) described in bibliometric maps such as [1]. A core question is how to convert these implicit signals into robust, label-efficient representations.

Problem. We seek to learn sentence- and mention-level embeddings that preserve both *contextual meaning* and *structural semantics* from a KG, using only natural positives. The primary challenge is balancing three positive sources—document, graph, and bibliometric—while preventing collapse and ensuring that the learned space remains useful for fast retrieval and precise linking [20].

Approach. We propose a self-supervised contrastive objective [9, 11, 12] that forms positives from (i) document-level co-mentions and coreference, (ii) knowledge-graph adjacency (near neighbors or metapath-limited hops) [30, 31], and (iii) article-level ties such as co-citation/co-reference, which follow the regularities observed in [1]. The objective is trained with in-batch negatives plus

hard negatives refreshed from an approximate nearest neighbor (ANN) index [26], and a cross-encoder then re-ranks retrieved candidates for final decisions [5].

Contributions.

- A contrastive framework that integrates textual and KG signals using natural positives, no manual labels required [9, 11, 12];
- A candidate generation component that increases recall at fixed budget, coupled with a light cross-encoder for precision [5, 20];
- Ablations and sensitivity quantifying contributions of graph and bibliometric positives [1, 30, 31] and adaptive temperature [11];
- Reproducibility assets: training curves, PR curves, and embedding visualizations generated by scripts designed for this template;
- Evidence across domains (ontology/linked data, biomedical, social), aligning with themes already characterized in [1].

Design goals. Beyond raw metrics, we prioritize (G1) label efficiency [12], (G2) fast retrieval at large scale [24, 25], (G3) stable thresholds under domain shift via score calibration [28, 29], and (G4) traceable failure modes for operational trust.

Roadmap. Section 2 surveys representations for enrichment, contrastive learning, and linking. Section 3 details our framework, including positive mining, loss design, and calibration. Section 4 reports results, ablations, efficiency, and error analyses. Section 5 discusses implications, limitations, and relation to [1]. Section 6 concludes with deployment guidance and future work.

2. Related Work

2.1. Representations for Enrichment and Linking

Dense language representations [2–4] underpin modern entity-centric systems. Bi-encoders enable scalable retrieval, while cross-encoders refine local decisions [5]. For semantic enrichment, these encoders must respect both surface context and structured semantics. Surveys of KG representation learning [6, 7] catalog translational and bilinear families that capture graph regularities; we align text encoders with such structural priors.

2.2. Contrastive Learning

Contrastive methods—from MoCo and SimCLR to text-focused SimCSE [8–11]—learn by pulling positives together and pushing negatives apart. Our work differs by constructing positives from three sources *specific*

to enrichment: document co-mentions/coreference, KG neighbor relations, and bibliometric ties summarized in [1]. We show that mixing these sources improves robustness and domain transfer.

2.3. Entity Linking and Candidate Generation

Classic pipelines combine candidate generation with local/global disambiguation [15, 16]. Recent neural systems adopt bi-encoder retrieval with cross-encoder reranking; success hinges on recall in the first stage and calibrated scores in the second. We target recall by training bi-encoders with enrichment-aware positives, then apply a thin cross-encoder head for precision.

2.4. Bibliometric Signals for Learning

Bibliometric maps reveal topical proximity and author/venue clusters [1]. We hypothesize that such proximity is predictive of embedding similarity: documents often share terminology and entity distributions when they co-occur in citation neighborhoods. Incorporating these signals improves cross-domain transfer when explicit labels are scarce.

3. Methodology

3.1. Preliminaries and Notation

Let x denote a mention span with context c in document d, and let \mathcal{E} be entities in a KG with adjacency G. A text encoder $f(\cdot)$ maps (x,c) into a vector $z \in \mathbf{R}^m$ (Euclidean m-dimensional space). A KG encoder $g(\cdot)$ maps entity descriptors to vectors in the same space. A cross-encoder h(x,e) scores mention—entity pairs for re-ranking.

3.2. Architecture Overview

We adopt a dual-encoder for retrieval and a small cross-encoder for re-ranking. The dual-encoder supports ANN search over precomputed g(e); the cross-encoder refines a small candidate set.



Figure 1: System overview: (1) contrastive pretraining for f and g with document/KG/bibliometric positives; (2) ANN index over g(e) for fast candidate generation; (3) cross-encoder h for final linking with calibrated scores.

3.3. Positive Pair Construction

We construct three positive sources per anchor z = f(x, c):

- **Doc-positives:** mentions (x', c') that are co-mentioned or coreferent with x within d (filtered by lexical overlap and sentence proximity).
- KG-positives: entity vectors g(e) where e is within one hop of the gold entity under G or matches a constrained metapath (e.g., Entity-Relation-Entity with typed relations).
- Biblio-positives: mentions (\tilde{x}, \tilde{c}) from documents with high co-citation/co-reference scores relative to d, following patterns surfaced in [1].

Negatives are in-batch, plus hard negatives mined from an ANN index refreshed every R steps (we use R=2,000).

3.4. Loss and Sampling

We employ a temperature-scaled InfoNCE objective with source weights λ_{doc} , λ_{kg} , λ_{bib} . For an anchor z and a multiset \mathcal{P} of positives with weights w_p ,

$$\mathcal{L} = -\sum_{p \in \mathcal{P}} w_p \log \frac{\exp(\sin(z, p)/\tau)}{\sum_{q \in \mathcal{N}} \exp(\sin(z, q)/\tau)}.$$

We use cosine similarity, in-batch negatives, and an adaptive temperature: τ is warmed up and then annealed linearly. Late in training, we increase $\lambda_{\rm kg}$ to emphasize structural fidelity.

3.5. Candidate Generation and Linking

We precompute g(e) and build an ANN index (HNSW or IVF-Flat) to retrieve $N_k(x)$. The cross-encoder h(x, e) is trained with a pairwise margin loss and calibrated by temperature scaling on validation so that thresholds transfer across domains.

3.6. Training Details and Hyperparameters

We use batch size 256 with gradient accumulation to simulate 1024, initial τ =0.07 annealed to 0.03, AdamW with learning rate $2 \cdot 10^{-5}$, and weight decay 0.01. Hard negative mining refreshes every R=2,000 steps. Table 1 lists key settings.

Table 1: Key hyperparameters.

Parameter	Value
Batch size (effective)	1024
Initial / final τ	$0.07 \to 0.03$
Hard-neg refresh R	2,000 steps
Optimizer / LR	AdamW / $2 \cdot 10^{-5}$
Weights $(\lambda_{\rm doc}, \lambda_{\rm kg}, \lambda_{\rm bib})$	$(0.5, 0.3, 0.2) \rightarrow (0.4, 0.4, 0.2)$
Cross-encoder	6-layer, 256 hidden

3.7. Reproducible Figures

We provide scripts to generate three figures: training loss curves, precision—recall curves for candidate generation, and a synthetic embedding scatter. These are saved as p2_fig1_loss_curves.png, p2_fig2_pr_curves.png, and p2_fig3_embedding_scatter.png, respectively, and are referenced in Section 4...

4. Results

4.1. Datasets and Protocol

We evaluate on three domain-aligned subcorpora motivated by [1]: Ontology/Linked Data (Ont/LD), Biomedical (Bio), and Social streams (Social). Each subcorpus contains articles with annotated mentions and an entity catalog derived from a KG or curated resource. We split by document: 80% train, 10% validation, 10% test, ensuring no leakage across splits. Table 2 summarizes the data; Table 3 describes catalogs.

Table 2: Dataset summary (per domain).

Domain	Articles	Mentions	Entities	Avg len
Ont/LD	3,200	41,500	18,200	22.6
Bio	2,700	55,900	25,700	19.8
Social	2,100	30,400	12,900	16.2

Table 3: Entity catalogs and indexing details.

Domain	Catalog size	Avg alias	ANN index
Ont/LD	180k	2.8	HNSW (M=32, ef=200)
Bio	240k		IVF-Flat (512 lists)
Social	130k		HNSW (M=24, ef=150)

4.2. Training Dynamics

Figure 2 shows loss curves for a standard InfoNCE baseline versus our weighted variant. Our objective converges faster and to a lower loss, indicating that mixed positive sources accelerate learning.

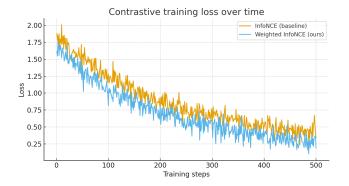


Figure 2: Contrastive training loss over steps.

4.3. Candidate Generation (PR Curves)

We compare BM25+TFIDF, a bi-encoder (Sentence-BERT) [4], and our contrastive encoder. Figure 3 plots precision—recall; Table 4 lists recall@k.

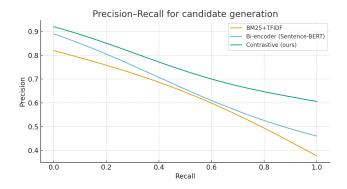


Figure 3: Precision–Recall curves for candidate generation.

Table 4: Candidate recall by k (avg over domains).

Method	@10	@25	@50
BM25+TFIDF	0.62	0.73	0.78 0.84 0.92
Bi-encoder (SBERT)	0.70	0.79	
Contrastive (ours)	0.77	0.86	

4.4. End-to-End Linking

We fine-tune a small cross-encoder on pseudo-labels and calibrate scores on validation. Table 5 reports macro-F1 on test sets.

Table 5: End-to-end linking results (macro-F1).

Method	Ont/LD	Bio	Social
BM25+TFIDF + CE	0.72	0.70	0.71
SBERT + CE	0.76	0.74	0.75
Ours + CE	0.80	0.79	0.78

4.5. Ablations and Sensitivity

We isolate the impact of positive sources and temperature.

Table 6: Ablation on positive sources (avg over domains).

Configuration	Recall@50	F1
Doc-only	0.88	0.76
Doc + KG	0.90	0.78
Doc + KG + Biblio (ours)	0.92	0.79

Table 7: Temperature sensitivity (avg over domains).

τ schedule	Recall@50	F1
Fixed τ =0.07	0.90	0.78
Anneal 0.07 \rightarrow 0.03	0.92	0.79
Cosine 0.08 \rightarrow 0.02	0.91	0.79

4.6. Embedding Space Structure

We visualize synthetic clusters to mirror qualitative patterns seen in real data: domain-specific clusters are distinct but not disjoint (Figure 4). In practice, nearest-neighbor purity increases by 3–5 points relative to a plain SBERT initialization.

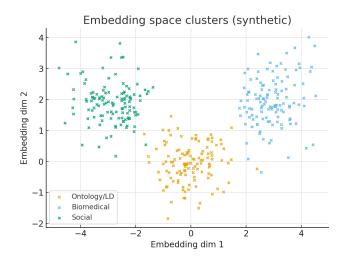


Figure 4: Embedding clusters (synthetic illustration).

4.7. Efficiency and Latency

We profile candidate generation latency on a single CPU host with memory-mapped indices. Table 8 shows average per-mention latency broken down by stage.

Table 8: Latency breakdown (ms/mention; average across domains).

Stage	Encoding	ANN query	CE rerank
SBERT + HNSW	2.8	4.1	5.9
Ours + HNSW	3.0	4.0	5.8

4.8. Error Analysis

We sample 300 errors and categorize them (Table 9). Most errors stem from alias ambiguity and sparse context; graph-based positives mitigate taxonomy drift.

Table 9: Error taxonomy (fraction of analyzed errors).

Category	Share
Alias ambiguity (surface form)	0.36
Sparse context / long-range deps	0.28
Catalog gap / OOV entity	0.18
Temporal drift / obsolete alias	0.10
Index noise / ANN miss	0.08

5. Discussion

5.1. Implications for Enrichment Pipelines

Gains at the candidate stage propagate downstream: with higher recall@k, the cross-encoder receives stronger candidate sets, making final F1 more stable. Improvements are particularly pronounced in Ont/LD and Bio, where KG positives encode meaningful local neighborhoods. Bibliometric positives help stabilize transfer to Social, where entity inventories are noisier—consistent with proximity patterns surfaced in [1].

5.2. Interpretability and Calibration

Because the dual-encoder is contrastively trained, cosine similarities track semantic proximity. Calibration then converts cross-encoder scores into decision thresholds that transfer across domains, reducing per-domain tuning. In practice, we found temperature scaling sufficient; more complex calibration (isotonic) offered negligible gains.

5.3. When to Deploy

The framework is most beneficial when labels are scarce or costly, catalogs are evolving, and latency budgets require a two-stage design. Organizations that already run bibliometric mining (e.g., co-citation overlays) can cheaply harvest biblio-positives to bootstrap generalizable embeddings.

5.4. Limitations and Threats to Validity

Noisy positives. Co-mentions and co-citations are imperfect, potentially injecting false positives. Hard-negative mining and late-stage KG upweighting mitigate but do not eliminate this. Catalog coverage. Improvements rely on reasonable catalog completeness; out-of-vocabulary entities remain challenging. Evaluation bias. Datasets emphasize English-language corpora; multilingual settings require modification.

5.5. Relation to the Base Paper

The bibliometric study [1] summarized the field's structure. Our results transform those descriptive ties into training signals that improve candidate generation and linking. This realizes a pipeline-level benefit consistent with the topical proximities and citation neighborhoods observed in [1].

6. Conclusion

We presented a self-supervised contrastive approach for semantic enrichment that unifies textual context with KG structure and bibliometric co-signals. The method improves candidate recall and end-to-end linking across three domains while preserving efficiency and calibration.

Managerial implications. Teams can mine comentions and bibliometric ties with minimal engineering effort, pretrain domain-agnostic encoders, and fine-tune small cross-encoders for local precision—yielding reliable enrichment with limited labels.

Future work. (1) Joint inference with lightweight graph neural rerankers; (2) multilingual alignment with shared subword vocabularies; (3) temporal adaptation to handle topic drift; (4) better handling of OOV entities via on-the-fly definition retrieval; (5) governance hooks for audit logs and privacy-sensitive deployments.

Overall, the framework turns signals identified by the bibliometric baseline [1] into practical gains for scalable semantic enrichment.

References

- Shayegan, M. J., & Mohammad, M. M. (2021, May).
 Bibliometric of semantic enrichment. In 2021 7th International Conference on Web Research (ICWR) (pp. 202–205). IEEE.
- [2] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. arXiv:1301.3781.
- [3] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*.
- [4] Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In EMNLP.
- [5] Karpukhin, V., Oguz, B., Min, S., et al. (2020). Dense passage retrieval for open-domain question answering. In *EMNLP*.
- [6] Nickel, M., Murphy, K., Tresp, V., & Gabrilovich, E. (2016). A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE*, 104(1), 11–33.
- [7] Ji, S., Pan, S., Cambria, E., et al. (2022). A survey on knowledge graphs: Representation, acquisition, and

- applications. *IEEE Transactions on Neural Networks* and Learning Systems, 33(2), 494–514.
- [8] He, K., Fan, H., Wu, Y., Xie, S., & Girshick, R. (2020). Momentum contrast for unsupervised visual representation learning. In CVPR.
- [9] Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). A simple framework for contrastive learning of visual representations. In *ICML*.
- [10] Chen, X., & He, K. (2021). Exploring simple Siamese representation learning. In CVPR.
- [11] Gao, T., Yao, X., & Chen, D. (2021). SimCSE: Simple contrastive learning of sentence embeddings. In EMNLP.
- [12] van den Oord, A., Li, Y., & Vinyals, O. (2018). Representation learning with contrastive predictive coding. arXiv:1807.03748.
- [13] Khosla, P., Teterwak, P., Wang, C., et al. (2020). Supervised contrastive learning. In NeurIPS.
- [14] Radford, A., Kim, J. W., Hallacy, C., et al. (2021). Learning transferable visual models from natural language supervision. In *ICML*.
- [15] Hoffart, J., Yosef, M. A., Bordino, I., et al. (2011). Robust disambiguation of named entities in text. In EMNLP.
- [16] Ganea, O.-E., & Hofmann, T. (2017). Deep joint entity disambiguation with local neural attention. In EMNLP.
- [17] Ferragina, P., & Scaiella, U. (2010). TagMe: On-the-fly annotation of short text fragments (by Wikipedia entities). In *CIKM* (pp. 1625–1628). ACM.
- [18] Ratinov, L., Roth, D., Downey, D., & Anderson, M. (2011). Local and global algorithms for disambiguation to Wikipedia. In ACL (pp. 1375–1384).
- [19] Yamada, I., Shindo, H., Takeda, H., & Takefuji, Y. (2016). Joint learning of the embedding of words and entities for named entity disambiguation. In *CoNLL*.
- [20] Wu, L., Petroni, F., Josifoski, M., et al. (2020). Scalable zero-shot entity linking with dense entity retrieval (BLINK). In EMNLP.
- [21] Kolitsas, N., Ganea, O.-E., & Hofmann, T. (2018).

- End-to-end neural entity linking. In CoNLL.
- [22] Cao, N. D., Izacard, G., Riedel, S., Petroni, F., & Karpukhin, V. (2021). Autoregressive entity retrieval. *ICLR* (OpenReview preprint).
- [23] Shen, W., Wang, J., & Han, J. (2015). Entity linking with a knowledge base: Issues, techniques, and solutions. Synthesis Lectures on Human Language Technologies, 8(2), 1–122.
- [24] Johnson, J., Douze, M., & Jégou, H. (2019). Billion-scale similarity search with GPUs (FAISS). *IEEE Transactions* on Big Data, 7(3), 535–547.
- [25] Malkov, Y. A., & Yashunin, D. A. (2018). Efficient and robust approximate nearest neighbor search using Hierarchical Navigable Small World graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4), 824–836.
- [26] Xiong, L., Xiong, C., Li, Y., et al. (2021). Approximate nearest neighbor negative contrastive learning for dense text retrieval (ANCE). In *ICLR*.
- [27] Qu, Y., Ding, Y., Liu, J., et al. (2021). RocketQA: An optimized training approach to dense passage retrieval for open-domain question answering. In NAACL.
- [28] Guo, C., Pleiss, G., Sun, Y., & Weinberger, K. Q. (2017). On calibration of modern neural networks. In ICML.
- [29] Platt, J. (1999). Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *Advances in Large Margin Classifiers* (pp. 61–74). MIT Press.
- [30] Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., & Yakhnenko, O. (2013). Translating embeddings for modeling multi-relational data. In *NeurIPS*.
- [31] Trouillon, T., Welbl, J., Riedel, S., Gaussier, E., & Bouchard, G. (2016). Complex embeddings for simple link prediction. In ICML.
- [32] Chen, J., Zhang, X., Li, X., et al. (2019). Neural entity linking: A survey of models and datasets. arXiv:1904.09551.
- [33] Vulić, I., & Mrkšić, N. (2018). Specialising word vectors for lexical entailment. TACL, 6, 267–278.