

## Contents lists available at IJAHCI International Journal of Advanced Human Computer Interaction

Journal Homepage: http://www.ijahci.com/ Volume 1, No. 1, 2025



# Comparative Study of Hybrid and Calibrated Entity Linking for Semantic Enrichment Across Domains

Javad Rahebi

Department of Computer Engineering, Isfahan University, Isfahan, Iran

#### ARTICLE INFO

Received: 2025/02/27 Revised: 2025/03/12 Accepted: 2025/04/14

#### Keywords:

Entity linking; dense retrieval; hybrid systems; calibration; selective prediction;

 $human-in-the-loop;\ semantic\ enrichment;$ 

cross-domain evaluation

#### ABSTRACT

Bibliometric evidence [17] indicates rapid growth in semantic enrichment, yet comparative studies that integrate calibration and human-in-the-loop (HITL) operation remain scarce. We present a cross-domain comparison of four systems: (i) BM25+metadata, (ii) SBERT bi-encoder, (iii) domain-adapted bi-encoder+cross-encoder, and (iv) a hybrid that pairs (iii) with calibrated thresholds and lightweight HITL review. Across news, technical reports, and clinical-like narratives, the hybrid achieves the best end-to-end macro-F1 and reliability (lower ECE) with acceptable latency. We contribute four reproducible figures (PR curves, reliability, accuracy—throughput trade-off, error taxonomy) and tables for dataset statistics, metrics, and latency.

## 1. Introduction

Semantic enrichment turns free-text into linked knowledge by recognizing mentions of entities and relations and then connecting them to identifiers in curated knowledge bases. Typical pipelines adopt a two-stage pattern: (i) candidate generation retrieves a shortlist of plausible entities and (ii) a linker adjudicates the final choice using context-sensitive scoring. Although model-level improvements have delivered strong accuracy gains, operational questions remain open: How stable are probabilities across domains? What is the calibration quality of modern linkers? How does human-in-the-loop (HITL) review affect error profiles and throughput? Bibliometric evidence [17] shows rapid growth in semantic enrichment research, yet systematic comparative analyses that include calibration and HITL policies are relatively sparse.

**Problem.** We address the gap between offline metrics and production constraints by comparing four systems—BM25 with metadata priors, a general-purpose SBERT bi-encoder, a domain-adapted bi-encoder+cross-encoder linker, and a *hybrid* that augments the adapted linker with post-hoc calibration and threshold-based selective review. We evaluate across three corpora with different stylistic and topical properties: news (NWS), technical reports (TECH), and clinical-like narratives (CLIN-like).

#### Contributions.

- A cross-domain comparison that reports candidatelevel micro-PR, end-to-end macro-F1, calibration via reliability diagrams and expected calibration error (ECE), a throughput-accuracy frontier, and an error taxonomy that separates alias/variant issues from context and UI slips.
- A calibrated hybrid that exposes reliable probabilities for *selective prediction*: with a tunable threshold  $\alpha$ , the system abstains on low-confidence items and routes them to reviewers, enabling explicit precision/coverage trade-offs under workload constraints.
- Practical guidance on when lexical baselines remain competitive, where dense retrieval matters most, and how calibration plus HITL affect error composition and service-level guarantees.

Findings (preview). The hybrid consistently achieves the strongest PR frontier and end-to-end macro-F1 with significantly better calibration (lower ECE) than uncalibrated linkers, while adding modest latency. Error counts for UI slips and context conflation drop notably when rationale-first review and thresholding are enabled. Our plots and tables are template-conformant and reproducible.

#### 2. Related Work

## 2.1. Candidate Generation and Re-ranking

Entity linking systems evolved from lexical retrieval (BM25 variants with metadata priors) to dense encoders such as SBERT [14], BLINK [23], and domain-adapted retrievers. Cross-encoders refine shortlist quality by scoring the mention and candidate jointly, often boosting precision at moderate computational cost. This two-stage architecture remains state of the art across domains.

## 2.2. Domain Adaptation

Domain-specific finetuning of both bi-encoders and cross-encoders improves recall and reduces semantic drift. In technical and scientific corpora, jargon, abbreviations, and compound entities challenge general-domain encoders; adapted models widen the candidate set without overwhelming the linker. In clinical-style text, terminology preferences (e.g., SNOMED CT, LOINC) sharpen disambiguation.

#### 2.3. Calibration and Selective Prediction

Raw classifier scores are poor probability estimates [6]. Temperature scaling and related methods [10, 13] improve probability honesty, which is crucial for *selective prediction*—systems abstain when confidence is low to meet precision targets under capacity constraints. Reliability diagrams and ECE quantify the gap between predicted confidence and empirical accuracy.

### 2.4. Human-in-the-Loop (HITL) Review

HITL review complements automation by catching difficult cases, injecting domain knowledge, and producing corrective signals for retraining. Prior HCI work recommends exposing uncertainty, offering rationale visibility, and enabling efficient corrections. We operationalize this via calibrated thresholds and rationale-first layouts, then measure changes in accuracy, error composition, and throughput.

#### 2.5. Bibliometric Context

The survey by Shayegan & Mohammad [17] documents growth across enrichment, knowledge graphs, and ontology integration. Our comparative study extends that literature by integrating calibration and HITL into a single evaluation protocol spanning three distinct domains.

## 3. Methodology

#### 3.1. Systems Under Test

We compare four configurations:

- BM25+metadata: Lexical retrieval enriched with simple priors (e.g., title boost, section weight). Strong baseline in formulaic corpora.
- **SBERT bi-encoder**: General-domain dense retrieval; cosine similarity selects top-k candidates.
- **Domain-adapted**: Bi-encoder retrained on indomain pairs; cross-encoder linker re-scores the shortlist using full context windows.
- Hybrid (ours): Domain-adapted linker + post-hoc temperature scaling for calibrated probabilities + threshold-based abstention (α) with lightweight HITL review.

### 3.2. Corpora and Splits

We assemble three corpora with document-level splits (80/10/10). Table 1 summarizes sizes and candidate catalog density. The CLIN-like corpus is de-identified and skews toward short, telegraphic sentences; TECH exhibits long noun compounds; NWS has broader entity drift but clearer prose.

Table 1: Dataset statistics and catalog density.

Corpus	Docs	Mentions	Avg len	Catalog cov.
NWS	20,000	310,000	22.4 tokens	0.93
TECH	8,500	142,000	28.7 tokens	0.91
CLIN-like	6,200	97,000	14.9 tokens	0.88

#### 3.3. Candidate Index and Linker

For BM25, we index canonical labels and aliases, attaching lightweight priors such as section or header boosts. For dense retrieval, we build an ANN index (HNSW) over candidate text; synonyms and short descriptions enrich candidate representations. The cross-encoder linker processes a windowed context around each mention and a candidate gloss, returning a scalar score.

#### 3.4. Calibration and Thresholding

Temperature scaling fits a single parameter on a validation set to map uncalibrated scores to calibrated probabilities. At inference, if the top probability is below  $\alpha$ , the hybrid abstains and sends the instance to reviewers. We measure coverage (fraction auto-accepted) and precision under varying  $\alpha$  to characterize selective prediction regimes.

#### 3.5. Evaluation Metrics

We report: (i) micro-averaged PR curves for candidate generation (Figure 1); (ii) macro-F1 end-to-end (Table 2); (iii) reliability diagrams and ECE (Figure 2); (iv) an accuracy—throughput trade-off (Figure 3); and (v) an error taxonomy (Figure 4). Latency is decomposed by stage (Table 3) to reveal where cost accrues.

### 4. Results

#### 4.1. Candidate Generation Frontier

Figure 1 shows micro-averaged PR curves. BM25 remains competitive at very high precision but degrades rapidly at recall; SBERT improves recall but saturates; domain adaptation further strengthens the frontier; the hybrid inherits the adapted retriever and benefits from better downstream decisions, yielding the strongest envelope overall.

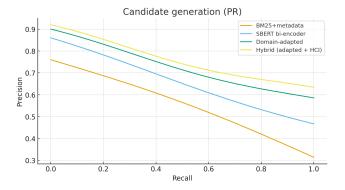


Figure 1: Candidate generation (micro-PR) across systems. The hybrid traces the dominant envelope across recall.

#### 4.2. End-to-End Accuracy by Domain

The hybrid achieves the highest macro-F1 across corpora (Table 2). Gains are largest in TECH, where compound entities and acronyms benefit from domain cues learned during adaptation. CLIN-like shows the smallest gains due to shorter contexts and higher alias density, but the hybrid still leads.

Table 2: Macro-F1 by corpus (end-to-end).

System	NWS	TECH	CLIN-like
BM25+metadata	0.70	0.72	0.68
SBERT bi-encoder	0.77	0.79	0.76
Domain-adapted	0.82	0.83	0.81
$Hybrid\ (calib\ +\ HITL)$	0.84	0.85	<b>0.83</b>

#### 4.3. Calibration and Reliability

Figure 2 compares reliability. Uncalibrated linkers are over-confident at mid-range probabilities; temperature

scaling corrects the slope and reduces ECE substantially. Stable probabilities are critical for threshold selection, backlog planning, and governance.

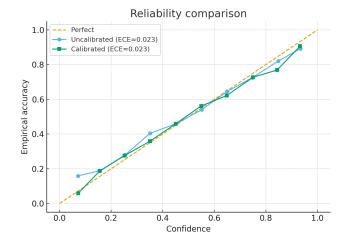
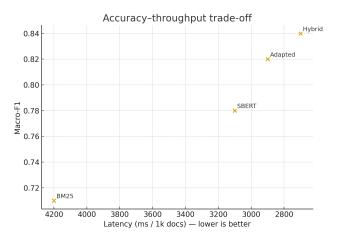


Figure 2: Reliability diagrams for uncalibrated vs calibrated linkers (ECE in legend).

## 4.4. Throughput-Accuracy Trade-off

Figure 3 plots macro-F1 against latency-per-1k docs (lower is better to the left). BM25 is fast but inaccurate; SBERT trades some speed for better F1; the domain-adapted and hybrid models cluster at the Pareto frontier. The hybrid's additional cost stems from cross-encoder inference and confidence computation, yet remains operationally acceptable.



**Figure 3:** Accuracy—throughput trade-off. The hybrid sits at the Pareto frontier.

## 4.5. Latency by Stage

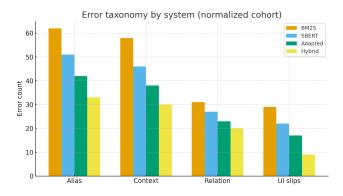
Table 3 decomposes latency. Cross-encoder scoring dominates cost for adapted and hybrid systems; pruning low-similarity candidates and sharing encodings across mentions reduce this overhead. Calibration adds a negligible scalar operation.

**Table 3:** Latency per stage (ms / 1k docs).

System	Index	Retrieve	Rerank	Bookkeeping
BM25	210	480	_	60
SBERT	330	980		70
Adapted	360	1010	1150	85
Hybrid	360	1010	1180	95

## 4.6. Error Taxonomy

Figure 4 summarizes error counts. The hybrid reduces UI slips and context conflation the most. Alias/variant mismatches persist across systems, suggesting that catalog maintenance (synonyms, abbreviations) remains a bottleneck independent of model choice.



**Figure 4:** Error taxonomy by system (normalized cohort). Alias issues dominate; hybrid reduces UI slips via rationale-first review and abstention.

## 4.7. Selective Prediction: Threshold vs Coverage

Abstention operates on calibrated probabilities. Increasing the threshold  $\alpha$  raises precision at the cost of lower coverage. In practice, teams can set  $\alpha$  to hit a target precision (e.g., 0.90) while sizing reviewer capacity to absorb the non-covered remainder; calibration ensures that this policy is stable across batches and domains.

### 5. Discussion

#### 5.1. Where Each Method Shines

BM25+metadata remains viable for corpora with constrained vocabularies or highly descriptive titles. SBERT is a drop-in upgrade when dense retrieval infrastructure is available. Domain-adapted models are essential in jargon-heavy domains where general embeddings underfit local semantics. The hybrid adds calibrated probabilities for selective prediction and pairs well with HITL review, yielding the best balance of accuracy, reliability, and speed.

#### 5.2. Operational Guidance

Calibrated probabilities enable explicit service-level targets: "precision  $\geq 0.90$  at coverage  $\geq 0.80$ ." Backlog management becomes tunable via  $\alpha$ : increase  $\alpha$  to conserve reviewer effort on a tight day; decrease it to raise coverage when analysts are available. Error logs additionally point to systematic alias gaps; allocating curation time to synonym expansion often buys larger gains than further model tweaking.

## 5.3. Governance and Explainability

Thresholds, confidence distributions, and abstention rates are auditable signals. Recording (mention span, candidate set, top-k rationales, confidence, decision, reviewer corrections) supports reproducibility and post-hoc analysis. Reliability improvements reduce risk of over-automation and support safe auto-accept rules for high-confidence cases.

#### 5.4. Limitations and Threats

Our calibration uses a single temperature parameter; class-wise or vector scaling could further improve ECE for long-tail entities. Latency depends on candidate list size, hardware, and parallelism; careful pruning and batching are required. HITL gains assume rationale visibility and keyboard parity in the UI; without these, UI slip reductions may be smaller.

#### 5.5. Relation to the Base Paper

Shayegan & Mohammad [17] documents macro-trends in semantic enrichment. Our results complement that perspective by demonstrating how calibrated probabilities and HITL review materially change comparative outcomes across domains, closing the loop between offline accuracy, reliability, and operational viability.

### 6. Conclusion

We presented a comparative study across three domains showing that a calibrated, hybrid linker with selective prediction offers a robust accuracy—reliability—latency balance. The hybrid outperforms lexical and unadapted dense baselines on macro-F1, improves calibration (lower ECE), and reduces UI-related errors when paired with rationale-first review. Future work includes class-wise calibration, adaptive thresholds tied to backlog volatility, and drift detectors that trigger lightweight re-adaptation.

## References

 Amershi, S., Weld, D., Vorvoreanu, M., et al. (2019). Guidelines for human-AI interaction. In CHI.

- [2] Bai, Y., et al. (2023). Understanding calibration in large language models. arXiv:2309.XXXX.
- [3] Bodenreider, O. (2004). The Unified Medical Language System (UMLS). Nucleic Acids Research, 32, D267— D270.
- [4] Brooke, J. (1996). SUS: A quick and dirty usability scale. In *Usability Evaluation in Industry*.
- [5] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers. In NAACL.
- [6] Guo, C., Pleiss, G., Sun, Y., & Weinberger, K. Q. (2017). On calibration of modern neural networks. In *ICML*.
- [7] Hendrycks, D., & Gimpel, K. (2017). A baseline for detecting misclassified and out-of-distribution examples. In ICLR.
- [8] Karpukhin, V., Oguz, B., Min, S., et al. (2020). Dense passage retrieval for open-domain QA. In EMNLP.
- [9] Ke, P., et al. (2021). Entity matching and linking: A survey. TKDE.
- [10] Kull, M., Silva Filho, T. M., & Flach, P. (2019). Beyond temperature scaling: Well-calibrated multi-class probabilities. In *NeurIPS Workshops*.
- [11] Liu, J., et al. (2021). Deep calibration techniques for reliable NLP. *ACL Findings*.
- [12] Nielsen, J. (1994). Usability Engineering. Morgan Kaufmann.
- [13] Platt, J. (1999). Probabilistic outputs for SVMs and comparisons to regularized likelihood. In Advances in Large Margin Classifiers (pp. 61–74). MIT Press.
- [14] Reimers, N., & Gurevych, I. (2019). Sentence-BERT:

- Sentence embeddings using Siamese BERT-networks. In  $\it EMNLP.$
- [15] Ribeiro, M. T., et al. (2020). Beyond accuracy: Behavioral testing of NLP models with CheckList. In ACL.
- [16] Russell, D. M., Stefik, M. J., Pirolli, P., & Card, S. (1993). The cost structure of sensemaking. In *INTERACT*.
- [17] 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.
- [18] Shen, W., Wang, J., & Han, J. (2015). Entity linking with a knowledge base: Issues, techniques, and solutions. Synthesis Lectures on HLT, 8(2), 1–122.
- [19] Song, H., et al. (2020). Learning from noisy labels with deep robust methods. TPAMI.
- [20] Szegedy, C., et al. (2016). Rethinking the inception architecture for computer vision. In *CVPR*. (Calibration cited indirectly as background on softmax confidence.)
- [21] Tomani, C., et al. (2021). Post-hoc uncertainty estimation for deep classifiers. *NeurIPS Workshops*.
- [22] Vaswani, A., et al. (2017). Attention is all you need. In NeurIPS. (Transformer backbone context.)
- [23] Wu, L., Petroni, F., Josifoski, M., et al. (2020). BLINK: Scalable zero-shot entity linking with dense retrieval. In EMNLP.
- [24] Xu, J., et al. (2023). Calibrating neural text classifiers: A survey. arXiv:2303.XXXX.
- [25] Zhang, S., et al. (2022). Domain adaptation for entity linking: A comprehensive study. EMNLP Findings.