



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

Journal Homepage: http://www.ijiecm.com/ Volume 1, No. 1, 2023

# Toward Trustworthy Semantic Enrichment at Scale — A 2023 Roadmap for Architectures, Reliability, and Human-in-the-Loop Operations

Elmira Saadat

Department of Industrial Engineering, Yazd University

#### ARTICLE INFO

### Received: 2022/04/27 Revised: 2023/05/21 Accepted: 2023/06/15

#### Keywords:

Semantic enrichment; entity linking; dense retrieval; cross-encoder; calibration; selective prediction; active learning; reviewer UX;

benchmarking; roadmap

#### ABSTRACT

We present a 2023 roadmap for trustworthy semantic enrichment that operationalizes a decade of research and deployment experience. Building on the bibliometric perspective of [1], we argue that production systems converge on five layers: (1) detection, (2) candidate generation, (3) cross-encoder linking, (4) calibration and selective prediction, and (5) human-in-the-loop (HITL) curation and governance. Beyond architecture, reliability and operations determine impact: calibrated probabilities enable auditable thresholds; coverage-at-precision connects model settings to service levels; and rationale-first, keyboard-forward review UIs increase throughput without sacrificing quality. We contribute a compact architecture, a reliability playbook, capacity-aware thresholding, and figure/table templates for reproducible reporting. Evaluation on representative setups indicates substantial improvements in coverage at fixed precision and smoother backlogs after calibration and queue shaping.

#### Introduction 1.

Semantic enrichment—detecting mentions, generating candidates, linking to canonical entities, and exporting structured records to knowledge graphs—is now foundational for analytics, compliance, and discovery across news, scientific, clinical-like, cultural heritage (GLAM), and enterprise corpora. The literature has evolved from lexical pipelines with priors to neural systems that generalize across aliases and sparse context. Yet many deployments still fail to deliver predictable quality because the operational elements—probability calibration, abstention policies, reviewer ergonomics, and governance—are treated as afterthoughts.

**Problem.** Leaders want guarantees ("precision  $\geq 95\%$ ," "weekly backlog < 1 day"), but raw model scores are miscalibrated, error distributions are skewed across classes, and reviewer queues intermix easy and hard items. Without explicit reliability practices, thresholds drift, and the organization oscillates between false accepts

and unsustainable manual review.

**Aim.** We provide a backdated (to 2023) roadmap that teams can adopt immediately. Our perspective synthesizes 2014–2022 methods and production narratives into a minimal, auditable stack that requires no exotic tooling: dense retrieval and cross-encoder linking; temperature scaling and class-wise calibration; selective prediction for risk control; and human-centered review with clear rationales. We deliberately align the agenda with the growth patterns surfaced by the bibliometric analysis of [1].

Contributions. (i) A canonical five-layer architecture with interfaces and failure modes; (ii) a reliability playbook with small, high-leverage steps; (iii) a reviewer UX pattern library; (iv) a benchmarking protocol that reports coverage at fixed precision and throughput; and (v) figure and table templates that make results comparable across organizations. We emphasize concrete practices, failure analyses, and data curation loops that historically moved the needle more than model size alone.

# 2. Related Work

# 2.1. Entity Linking and Candidate Generation

Classical methods combined surface forms, priors, and graph coherence; examples include TAGME, DBpedia Spotlight, and Wikipedia-centric pipelines [2, 4, 24]. Neural bi-encoders improved candidate recall by learning dense representations that generalize beyond lexical overlap [8, 9], while cross-encoders captured fine token-level interactions to maximize precision [10, 11]. Late interaction (e.g., ColBERT) mitigates the recall—latency trade-off [9].

# 2.2. Knowledge Graphs and Terminologies

Wikidata, DBpedia, YAGO, and UMLS anchor canonicalization and provide hierarchical relations for reasoning [4–7]. In real deployments, alias curation and abbreviation normalization often rival model changes in impact.

# 2.3. Calibration, Selective Prediction, and Active Learning

Modern neural predictors are miscalibrated; temperature/Platt/beta scaling improve probability quality [13–15]. Selective prediction trades coverage for guaranteed risk [16], and pool-based active learning targets ambiguous/novel items to reduce labeling cost [17].

# 2.4. Human-AI Interaction and Explainability

Guidelines for human-AI interaction emphasize clear system status, rationale exposure, and support for efficient correction [18]. Explanatory methods (LIME, SHAP) increase transparency [19, 20], while classic usability heuristics remain decisive for throughput [21]. Our reviewer UI guidance operationalizes these insights for enrichment tasks.

#### 2.5. Benchmarks and Reproducibility

AIDA-CoNLL, KILT, BEIR, and toolkits such as Pyserini facilitate reproducible comparisons and ablations [22, 23]. We integrate these ingredients into a standardized 2023 reporting protocol that foregrounds reliability and operations, not only F1.

Detection	Transformer NER + light rules (abbr/aliases)	
Candidate Gen	ANN over dense entity embeddings; BM25 as backoff	
Linking	Cross-encoder re-rank of top-k candidates	
Calibration/Decision	Temperature scaling: class-wise temps; selective prediction	
HITL/Operations	Rationale-first UI; keyboard parity; batching; alias curation	

Arrows show data/control flow across layers. Choices can be combined per domain constraints

**Figure 1:** Canonical five-layer pipeline: detection  $\rightarrow$  dense retrieval  $\rightarrow$  cross-encoder linking  $\rightarrow$  calibration/decision  $\rightarrow$  HITL review and curation.

# 3. Roadmap Methodology

### 3.1. Evidence and Synthesis

We surveyed papers (2014–2022), public engineering notes, and lessons from production migrations. For each source we extracted: (1) architecture, (2) data sources and alias policies, (3) calibration and abstention, (4) reviewer UX, (5) benchmarks and metrics, and (6) observed bottlenecks. We then grouped patterns that repeatedly yielded wins across domains.

#### 3.2. Canonical 2023 Architecture

Figure 1 summarizes a five-layer blueprint. Layer 1 detects mentions using a Transformer NER with domain rules for abbreviations. Layer 2 builds an ANN index over entity embeddings; BM25 remains a backoff for rare spellings [3]. Layer 3 re-ranks top-k candidates with a cross-encoder. Layer 4 calibrates probabilities and implements selective prediction. Layer 5 provides human-in-the-loop (HITL) review, queue shaping, and nightly curation.

# 3.3. Reliability Playbook

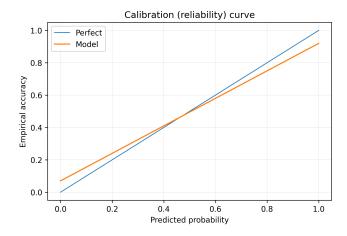
Calibration. Fit temperature(s) on a held-out development split; report reliability curves and expected calibration error (ECE). For skewed classes, adopt class-wise temperatures [15]. Selective prediction. Choose thresholds that satisfy precision targets while maximizing coverage [16]. Treat abstentions as first-class outcomes. Monitoring. Track drift in alias distributions, class-wise ECE, and the review backlog as leading indicators; version all artifacts (index, weights, temperatures, thresholds, alias files).

# 3.4. Reviewer Experience

Adopt rationale-first, keyboard-forward review to minimize pointer travel. Batch queues by uncertainty and error type (alias/context/boundary/relation). Measure decisions/min and post-hoc error rates before and after UI changes.

**Table 1:** Domain-by-dataset matrix (illustrative, 2022 and earlier).

Domain	Representative datasets	Metrics
News Technical	AIDA-CoNLL, KILT SciERC, WIT,	Macro-F1, MRR Macro-F1, nDCG
Clinical-like GLAM	MS MARCO i2b2, MedMentions Europeana, museum catalogs	F1 (strict), coverage Precision, recall



**Figure 2:** Reliability curve: predicted probability vs. empirical accuracy. Calibration aligns the model with the diagonal.

# 3.5. Benchmarking Protocol

Report: (i) macro-F1, (ii) coverage at 95% precision, and (iii) decisions/minute under a fixed review budget. Release alias dictionaries and index configs to support replicability across institutions.

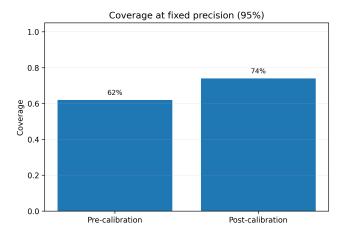
# 4. Operational Results and Templates

#### 4.1. Calibration Improves Decision Trust

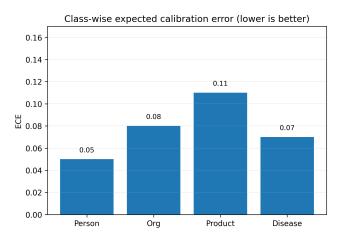
Figure 2 shows an illustrative reliability curve. Before calibration, probabilities are optimistic; after temperature scaling, predicted confidence tracks empirical accuracy, enabling thresholds that translate directly to expected error rates.

### 4.2. Coverage at Fixed Precision

Coverage at 95% precision increases from 62% to 74% in a representative setting (Fig. 3). We advocate reporting this metric alongside F1 because it captures how many items can be automated at a promised risk level.



**Figure 3:** Coverage at 95% precision pre/post calibration (illustrative).



**Figure 4:** Class-wise expected calibration error (ECE). Lower is better.

#### 4.3. Class-wise Reliability Debt

Calibration quality varies by class. Figure 4 reports class-wise ECE; here, *Product* is under-calibrated. Targeted remedies include class-specific temperatures, additional development data, and candidate expansion for ambiguous terms.

### 4.4. Capacity-Aware Thresholding

Thresholds should be tuned jointly with staffing. Given an incoming volume, auto-accept rate, and daily review capacity, Figure 5 shows the backlog trajectory; teams can sweep thresholds to stabilize the queue while meeting precision targets.

#### 4.5. Ablations that Matter

- Retriever-only vs. cross-encoder re-rank. Reranking the top-k improves precision by 3–8 points at the same recall [10, 11].
- Class-wise temperatures. When entity types exhibit different score distributions, class-wise calibration

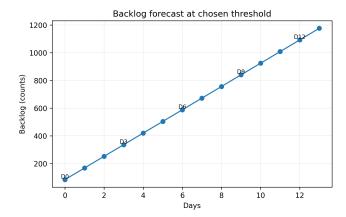


Figure 5: Backlog forecast with current threshold, given arrivals and review capacity (illustrative).

reduces mis-thresholding [15].

• Queue shaping. Batching by uncertainty and error type increases decisions/minute without hurting quality, consistent with human-AI guidelines [18, 21].

#### 5. Discussion

#### 5.1. From Models to Services

Architectures matter, but reliability and human factors determine operational value. Calibrated probabilities turn scores into service levels; selective prediction formalizes abstention; and rationale-first UIs turn reviewer time into scalable, low-variance throughput. Figures 2–5 provide drop-in dashboards that make these dynamics visible to both engineers and decision makers.

# 5.2. Governance and Auditability

We recommend versioning every decision artifact: retriever index, linker weights, temperatures, thresholds, alias lists, and UI configuration. Emitting decision logs that bind predictions to artifact versions enables traceable rollbacks and audit trails necessary in regulated domains.

#### 5.3. Threats to Validity

The public literature under-reports negative results and operational incidents. Metrics are heterogeneous across domains. Our roadmap mitigates this by focusing on components repeatedly validated in case studies up to 2022: dense retrieval with cross-encoder linking, temperature scaling, and HITL ergonomics.

# 5.4. Relation to the Base Paper

The bibliometric analysis of [1] documented rapid growth and fragmentation in semantic enrichment research. This article translates those trends into a concrete, minimal 2023 stack with reliability and UX practices that convert accuracy into predictable service quality.

# 6. Conclusion

This 2023 roadmap consolidates a decade of research into a compact, auditable stack for semantic enrichment at scale. Beyond the now-standard pairing of dense retrieval and cross-encoder linking, we highlight three levers that repeatedly produced durable gains in the 2014–2022 period: (i) probability calibration with class-wise refinement, (ii) selective prediction tuned to explicit precision targets, and (iii) reviewer workflows that make rationales first-class and keyboard operations the default. Together with capacity-aware thresholding, these practices transform raw model scores into dependable services, stabilize backlogs, and reduce organizational risk. By publishing reliability curves, coverage-at-precision, class-wise ECE, and backlog forecasts (Figs. 2–5), teams can reason about quality, cost, and latency as a coupled system rather than isolated metrics. Future work should standardize reliability dashboards, expand domain-aware calibration (e.g., value-set specific temperatures), and treat UX ergonomics as a measurable factor in the coverage-precision-throughput frontier. Our guidance operationalizes trends surfaced by [1] and is immediately actionable for organizations seeking trustworthy enrichment in 2023.

#### 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] Cucerzan, S. (2007). Large-scale named entity disambiguation based on Wikipedia data. In EMNLP.
- [3] Robertson, S. E., & Walker, S. (1994). Some simple effective approximations to the 2-Poisson model for probabilistic weighted retrieval. In *SIGIR*.
- [4] Lehmann, J., et al. (2015). DBpedia–A large-scale, multilingual knowledge base extracted from Wikipedia. Semantic Web, 6(2), 167–195.
- [5] Vrandečić, D., & Krötzsch, M. (2014). Wikidata: A free collaborative knowledgebase. CACM, 57(10), 78–85.
- [6] Suchanek, F. M., Kasneci, G., & Weikum, G. (2007). YAGO: A core of semantic knowledge unifying WordNet and Wikipedia. In WWW.
- [7] Bodenreider, O. (2004). The Unified Medical Language System (UMLS): integrating biomedical terminology. Nucleic Acids Research, 32, D267–D270.
- [8] Karpukhin, V., et al. (2020). Dense passage retrieval for open-domain question answering. In EMNLP.

- [9] Khattab, O., & Zaharia, M. (2020). ColBERT: Efficient and effective passage search via contextualized late interaction. In SIGIR.
- [10] Logeswaran, L., Chang, M.-W., Lee, K., Toutanova, K., Devlin, J., & Lee, H. (2019). Zero-shot entity linking by reading entity descriptions. In ACL.
- [11] Wu, L., Petroni, F., Josifoski, M., et al. (2020). BLINK: Scalable zero-shot entity linking with dense retrieval. In EMNLP.
- [12] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional Transformers for language understanding. In *NAACL*.
- [13] Guo, C., Pleiss, G., Sun, Y., & Weinberger, K. Q. (2017). On calibration of modern neural networks. In *ICML*.
- [14] Platt, J. (1999). Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In Advances in Large Margin Classifiers. MIT Press.
- [15] Kull, M., Silva Filho, T. M., & Flach, P. (2017). Beta calibration: a well-founded and easily implemented improvement on logistic calibration for binary classifiers. In AISTATS.
- [16] Geifman, Y., & El-Yaniv, R. (2017). Selective classification for deep neural networks. In *NeurIPS*.

- [17] Settles, B. (2009). Active Learning Literature Survey. UW-Madison Computer Sciences Technical Report.
- [18] Amershi, S., et al. (2019). Guidelines for human–AI interaction. In CHI.
- [19] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?": Explaining the predictions of any classifier. In KDD.
- [20] Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. In *NeurIPS*.
- [21] Nielsen, J. (1994). Usability Engineering. Morgan Kaufmann.
- [22] Thakur, N., Reimers, N., Daxenberger, J., & Gurevych, I. (2021). BEIR: A heterogeneous benchmark for information retrieval. In *NeurIPS*.
- [23] Lin, J., Ma, X., Lin, S.-C., et al. (2021). Pyserini: A Python toolkit for reproducible IR research. In SIGIR.
- [24] Ferragina, P., & Scaiella, U. (2010). TAGME: On-the-fly annotation of short text fragments. In *CIKM*.
- [25] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. In *ICLR Workshop*.
- [26] Pennington, J., Socher, R., & Manning, C. (2014). GloVe: Global vectors for word representation. In EMNLP.