
Hybrid End-to-End Knowledge Graph Construction and Validation: A Cross-Domain Study with LLM-as-a-Judge

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The automated construction of knowledge graphs (KGs) from unstructured text
2 remains a central challenge in information management and artificial intelligence.
3 This paper introduces a hybrid framework that combines the conceptual reasoning
4 of large language models (LLMs) with the efficiency of scalable, rule-based
5 methods to deliver an end-to-end pipeline for KG construction and validation. The
6 framework begins with ontology induction using an LLM to define domain-specific
7 entity and relation types, followed by large-scale rule-based information extraction,
8 entity resolution, and graph assembly. A novel extrinsic evaluation method,
9 *LLM-as-a-Judge*, is employed to assess the semantic quality of the resulting graphs.
10 We evaluate the pipeline across three diverse benchmarks. In the financial domain,
11 the FiQA dataset (5,500+ documents) yielded a graph with 475 nodes and 36
12 edges, achieving an overall quality score of 2.97/5 at a total cost of \$2.63. In the
13 document-level relation extraction setting, the DocRED dataset (100 annotated
14 documents) produced 5,000 nodes and 389 edges, with a lower quality score of
15 2.68/5, primarily due to systematic entity type misclassification. In the biomedical
16 domain, the CDR dataset (100 sampled abstracts) generated 966 nodes and 13 edges,
17 but achieved the highest semantic precision, with an overall quality of 3.91/5 at a
18 cost of \$0.65. Across all datasets, the pipeline demonstrated efficiency, with end-to-
19 end processing times under one hour, and highlighted complementary strengths and
20 weaknesses: FiQA emphasized scale but sparse connectivity, DocRED revealed
21 classification challenges, and CDR achieved high entity-level precision despite
22 graph fragmentation. These results validate the effectiveness of hybrid architectures
23 for KG construction: LLMs provide strong conceptual modeling, while rule-based
24 systems ensure scalability and cost-efficiency. The *LLM-as-a-Judge* framework
25 further supplies actionable feedback, exposing domain-specific error modes and
26 guiding refinement. Our work establishes a cost-effective, modular, and adaptable
27 methodology for automated KG construction, offering a foundation for future
28 research on improving connectivity, refining extraction accuracy, and extending to
29 new domains.

30 **1 Background and Related Work**

31 Automating the construction of knowledge graphs (KGs) from unstructured text remains a central
32 challenge in information extraction and artificial intelligence [16, 2]. KGs provide structured,
33 machine-readable representations that power semantic search, analytics, and question answering, yet
34 end-to-end systems must discover relevant concepts, extract factual assertions, and resolve entity
35 ambiguity at scale and under strict cost constraints [7]. We present a hybrid, modular framework
36 that couples the conceptual strength of a Large Language Model (LLM) with the efficiency of

37 deterministic, rule-based components. In our design, the LLM is used sparingly for ontology induction
38 and extrinsic evaluation (*LLM-as-a-Judge*), while information extraction (IE) and entity resolution
39 (ER) are implemented as scalable, low-cost rule-based stages. This separation preserves semantic
40 guidance where it matters most while keeping the bulk processing economical and reproducible.

41 Beyond the financial domain that motivated our initial system, we evaluate the framework across
42 three distinct corpora to assess portability and robustness: finance (FiQA), general document-level
43 relation extraction (DocRED [15]), and biomedicine (BioCreative V CDR). In the FiQA setting,
44 we induce a financial ontology with an LLM, run rule-based IE and ER over 5,500+ documents,
45 assemble the KG, and validate with an LLM judge. The resulting graph comprises 475 nodes and 36
46 edges, with an overall quality score of 2.97/5 at a total cost of approximately \$2.63. On DocRED,
47 where the schema is implicit in the annotations, we demonstrate the pipeline’s modularity by omitting
48 ontology induction and running a three-stage variant (IE → KG construction → judging). This yields
49 a 5,000-node/389-edge graph with an overall quality of 2.68/5 at roughly \$0.15 and exposes a salient
50 error mode: pervasive entity-typing drift (e.g., over-assignment to PERSON) [13, 8]. In the biomedical
51 CDR corpus, we reinstate LLM-based ontology induction to define a domain schema, apply rule-
52 based extraction and ER, construct the KG, and validate. The outcome is a 966-node/13-edge graph
53 with an overall quality of 3.91/5 at approximately \$0.65, illustrating improved entity-level fidelity in
54 a specialized domain but also highlighting sparse cross-document connectivity [12].

55 These experiments support three claims. First, careful placement of LLM capability—restricted to
56 concept modeling and lightweight extrinsic judging—can substantially reduce costs while retaining
57 semantic leverage [1]. Second, rule-based IE and ER scale across domains with only configuration
58 changes, enabling rapid transfer without model retraining [9, 11]. Third, a uniform LLM-as-a-Judge
59 protocol provides consistent, actionable feedback at node, edge, and graph levels, revealing domain-
60 specific failure modes: DocRED stresses entity typing, whereas FiQA and CDR emphasize relation
61 sparsity and KG connectivity [3, 5]. Collectively, the results argue for a practical path to cost-aware,
62 domain-adaptable KG construction: employ the LLM where its conceptual advantage is highest,
63 keep high-throughput stages deterministic, and close the loop with a diagnostic LLM-based evaluator
64 [7, 14].

65 Our contributions are fourfold. We introduce an end-to-end, modular KG pipeline that cleanly
66 separates LLM-centric reasoning (ontology induction and judging) from deterministic extraction and
67 resolution. We demonstrate cross-domain portability on FiQA, DocRED, and CDR with sub-\$3 total
68 cost per corpus and minute-scale runtimes. We formalize a stratified LLM-as-a-Judge evaluation
69 method that yields interpretable node, edge, and graph scores alongside error diagnostics. Finally, we
70 provide empirical evidence that different domains stress different components, motivating targeted
71 refinements such as typed NER for DocRED [4, 6], tighter IE–ER coupling to reduce relation loss in
72 FiQA and CDR [17, 19], and connectivity-enhancing strategies for sparsely linked graphs [10, 18].

73 **2 A Framework for Automated Knowledge Graph Construction and 74 Validation**

75 Our framework implements a multi-stage pipeline for the end-to-end construction and validation of a
76 knowledge graph from unstructured text corpora. The system follows a hybrid approach, combining
77 a Large Language Model (LLM) for initial conceptual modeling with scalable, rule-based methods
78 for large-scale data processing. This design balances the nuanced understanding of LLMs with the
79 efficiency and cost-effectiveness of deterministic algorithms.

80 Depending on the dataset and domain, the pipeline can operate in either three, four, or five stages.
81 In the most complete form (as used for FiQA and CDR), the pipeline includes Ontology Induction,
82 Information Extraction (IE), Entity Resolution (ER), Knowledge Graph Population, and LLM-based
83 Validation. For DocRED, where the ontology is pre-defined, the Ontology Induction stage is skipped,
84 and the pipeline runs with three stages (IE, graph construction, and validation). This flexibility
85 demonstrates the modular design of the system.

86 **2.1 System Pipeline Overview**

87 The pipeline begins by establishing a domain-specific schema (when required) and then uses that
88 schema to extract, resolve, and structure information from the source documents. Ontology Induction

89 uses an LLM to create a conceptual framework of entity and relation types from a small representative
90 sample. Guided by this ontology, the Information Extraction stage systematically processes the entire
91 document set to identify specific instances of these entities and relations using a rule-based engine.
92 Next, the Entity Resolution stage identifies and merges duplicate entity mentions into canonical
93 forms through similarity clustering. Finally, the Knowledge Graph Population stage assembles these
94 resolved entities and relations into a formal graph structure. Each component is designed to be
95 modular, allowing for independent operation and refinement across domains.

96 **2.2 Component 1: Ontology Induction**

97 When required (e.g., FiQA and CDR), the initial stage of the pipeline creates a conceptual ontology
98 that defines the semantic schema for the knowledge graph. The objective is to produce a formal
99 definition of relevant entity types and relation types specific to the target domain. To achieve this, the
100 framework utilizes a large-scale LLM (e.g., Llama 3.3 70B). Rather than processing the entire corpus,
101 which would be computationally expensive, an intelligent sampling strategy is employed where the
102 LLM analyzes a small subset of documents. The output is a structured ontology file that serves as
103 the guiding schema for subsequent extraction tasks. In datasets like DocRED, where the ontology is
104 provided with the annotations, this stage is skipped.

105 **2.3 Component 2: Information Extraction**

106 Following ontology induction (or using a pre-defined schema), the Information Extraction (IE)
107 component populates the schema with specific instances from the full text corpus. To ensure
108 scalability and cost efficiency, this stage is implemented as a completely rule-based system that avoids
109 LLM calls. The IE component processes the entire dataset, applying a comprehensive library of
110 extraction patterns for both entity recognition and relation extraction. Flexible matching techniques,
111 such as exact, partial, and word-overlap matching, are employed to maximize coverage.

112 **2.4 Component 3: Entity Resolution**

113 The raw output of the IE stage contains many duplicate or variant entity mentions. The Entity
114 Resolution (ER) stage disambiguates and consolidates these mentions into canonical forms. The
115 approach is based on rule-based similarity clustering, using metrics such as Jaccard, Levenshtein, and
116 Cosine similarity. The system applies type-aware thresholds to reflect the typical naming conventions
117 of different entity types. Mentions that are sufficiently similar are grouped into clusters, and a
118 canonical form is chosen for each.

119 **2.5 Component 4: Knowledge Graph Population**

120 The refined data is then assembled into a formal graph structure. This component takes the canonical
121 entities from ER as input for graph nodes and the extracted relations as input for edges. Using the
122 NetworkX library, the system instantiates the final graph and aggregates node/edge properties such as
123 document identifiers. The resulting KG is exported into multiple formats (JSON, GraphML, GEXF),
124 making it available for analysis, visualization, and downstream tasks.

125 **2.6 Component 5: LLM-as-a-Judge Validation**

126 The final stage assesses the semantic quality of the constructed knowledge graph. An LLM-based
127 judging framework is employed to evaluate sampled nodes and edges according to a multi-criteria
128 rubric. Stratified sampling ensures coverage across entity and relation types. This extrinsic evaluation
129 produces interpretable scores for nodes, edges, and the graph as a whole, while also identifying
130 common error modes (e.g., misclassification, vague labels, relation inconsistencies). In all datasets,
131 this stage provides actionable feedback for improving earlier pipeline components.

132 **3 Experimental Setup**

133 This section details the experimental design for constructing and evaluating knowledge graphs
134 across three benchmark datasets. We describe the corpora used, the implementation of our pipeline
135 components, and the metrics employed for both intrinsic and extrinsic evaluation.

136 **3.1 Text Corpora and Datasets**

137 We evaluate the framework on three diverse corpora spanning finance, general document-level relation
138 extraction, and biomedicine. The FiQA dataset consists of over 5,500 financial documents from
139 Hugging Face, processed in full for information extraction, entity resolution, and graph construction,
140 with a sample of 50 documents used for ontology induction. The DocRED dataset contains 100
141 human-annotated documents with entity and relation spans; because the schema is already provided,
142 the ontology induction stage is omitted in this setting. The BioCreative V CDR dataset consists of
143 1,500 PubTator abstracts annotated with chemical and disease entities and their relations; we sample
144 100 documents for processing and employ LLM-based ontology induction to generate a biomedical
145 schema. Together, these datasets stress different aspects of the pipeline: large-scale extraction in
146 FiQA, dense annotations and type imbalance in DocRED, and specialized biomedical terminology in
147 CDR.

148 **3.2 Pipeline Implementation**

149 The pipeline is implemented as a modular system in which each stage can be configured and executed
150 independently. It follows a hybrid model that combines an LLM for conceptual modeling with
151 rule-based systems for scalable processing. Depending on the dataset, the pipeline consists of the
152 following components:

- 153 **1. Ontology Induction:** For FiQA and CDR, the Llama 3.3 70B model is used via the
154 OpenRouter API to generate an ontology from a small sample of documents. In DocRED,
155 this stage is skipped since entity and relation types are provided.
- 156 **2. Information Extraction:** A rule-based system processes all documents in each dataset,
157 applying regular expression patterns and flexible string-matching strategies. FiQA uses
158 patterns for 10 entity and 15 relation types, DocRED for 6 entity and 7 relation types, and
159 CDR for a biomedical-specific schema including chemicals, diseases, and treatments.
- 160 **3. Entity Resolution:** For FiQA and CDR, entity mentions are consolidated into canonical
161 forms using similarity clustering with Jaccard, Levenshtein, and Cosine similarity metrics
162 under type-aware thresholds. DocRED does not require this step because entities are
163 pre-disambiguated in the dataset.
- 164 **4. Knowledge Graph Population:** Across all datasets, the NetworkX library is used to
165 construct the graph from canonical entities and extracted relations, and graphs are exported
166 into JSON, GraphML, and GEXF formats.
- 167 **5. LLM-as-a-Judge Validation:** For all datasets, the Llama 3.3 70B model evaluates sampled
168 nodes and edges. Stratified sampling ensures balanced coverage across entity and relation
169 types: 100 items (50 nodes and 50 edges) for FiQA, 50 items for DocRED, and 50 items for
170 CDR.

171 This modular design allows for consistent comparison between LLM-based conceptual modeling
172 and rule-based extraction, while permitting dataset-specific adjustments such as skipping ontology
173 induction in DocRED.

174 **3.3 Evaluation Metrics**

175 We adopt a two-part evaluation strategy that combines intrinsic and extrinsic measures. Intrinsic
176 evaluation captures structural and statistical properties of intermediate and final outputs. For FiQA
177 and CDR, entity resolution quality is measured by the resolution rate, defined as the proportion of raw
178 mentions successfully grouped into canonical clusters. For all datasets, graph-theoretic properties
179 such as the number of nodes, number of edges, density, the number of connected components, and

180 the size of the largest component are computed. Processing efficiency is also tracked in terms of
181 runtime and monetary cost.

182 Extrinsic evaluation is conducted with the LLM-as-a-Judge framework. Stratified samples of nodes
183 and edges are rated on a 1–5 scale according to correctness of entity types, clarity of labels, semantic
184 validity of relations, and contextual alignment. The outputs are aggregated into node, edge, and
185 graph-level quality scores, supplemented with diagnostic feedback on systematic errors such as
186 misclassification, vague labels, or relation inconsistencies. This combination of intrinsic and extrinsic
187 metrics provides a comprehensive view of both pipeline performance and semantic quality. All code
188 can be found at [here](#).

189 4 Results and Analysis

190 We now present the results obtained by executing our pipeline across the three datasets: FiQA,
191 DocRED, and CDR. The analysis covers ontology induction, information extraction, entity resolution,
192 knowledge graph construction, and final validation with an LLM-as-a-Judge. Results are reported in
193 terms of structural statistics, processing efficiency, cost, and semantic quality.

194 4.1 FiQA Results

195 In the FiQA financial domain, ontology induction using the Llama 3.3 70B model successfully
196 produced a conceptual schema from a 50-document sample, yielding 42 entity concept types and
197 120 relation types. This stage required approximately five minutes of processing at a cost of \$2.50.
198 The subsequent rule-based information extraction stage processed the full corpus of more than
199 5,500 documents, identifying over 15,000 entity mentions and 10,276 candidate relations. Entity
200 resolution consolidated these mentions into 475 canonical clusters, achieving a resolution rate of
201 95.9%. Knowledge graph population resulted in a sparse graph of 475 nodes and 36 edges, with
202 a density of 0.0002 and 448 connected components. The largest component contained 27 nodes,
203 reflecting the topical clustering typical of financial documents. Extrinsic validation was conducted
204 over 100 sampled nodes and edges, all of which were successfully evaluated. The LLM judge
205 assigned an overall quality score of 2.97/5, with node quality rated slightly higher at 3.03/5 and edge
206 quality slightly lower at 2.90/5. The primary weaknesses identified were entity type misclassification
207 and vague entity labels, while relation quality showed inconsistency across types.

208 4.2 DocRED Results

209 For DocRED, the pipeline operated without ontology induction, relying instead on the dataset’s
210 annotated schema. Information extraction from 100 documents yielded approximately 5,000 entity
211 mentions and 389 relations. Knowledge graph construction produced a graph with 5,000 nodes and
212 389 edges, characterized by extreme sparsity: the density was only 0.00003, with 4,611 connected
213 components. The largest component contained 389 nodes, corresponding to a subset of tightly
214 interlinked documents. Validation was carried out on 50 samples, consisting of 25 entities and 25
215 relations, with a 96% success rate. The overall quality score was lower than FiQA, at 2.68/5. Entity
216 quality was the weakest dimension, averaging 2.34/5, largely due to severe type misclassification,
217 with 98.6% of entities labeled as PERSON. Relation quality was somewhat stronger at 2.98/5, and
218 graph quality averaged 2.50/5. These results demonstrate that while the pipeline can process DocRED
219 effectively, it exposes a clear failure mode in entity classification when the schema is imbalanced or
220 difficult to capture through rules alone.

221 4.3 CDR Results

222 In the biomedical domain, ontology induction was again performed using the Llama 3.3 70B model,
223 which defined a schema of 15 concept types and 25 relation types from a sample of 100 docu-
224 ments. Rule-based extraction then identified 966 unique entities and a set of relations centered on
225 chemical–disease interactions. Entity resolution grouped mentions into canonical forms, supporting
226 coherent biomedical concepts. The assembled knowledge graph contained 966 nodes and 13 validated
227 relations, reflecting the narrow but highly specialized focus of the corpus. Graph density was 0.0002,
228 and although the graph was fragmented, it provided meaningful biomedical substructures. Validation
229 on 50 samples achieved a high overall quality of 3.91/5, with entity quality at 4.03/5, relation quality

230 at 3.85/5, and graph quality somewhat weaker at 2.20/5 due to sparsity and limited connectivity.
231 Processing time was approximately 6.5 minutes at a total cost of \$0.65. The biomedical domain thus
232 produced the most semantically precise results, though at the expense of global connectivity.

233 **4.4 Cross-Domain Comparison**

234 Across the three datasets, the pipeline demonstrated strong portability and consistent processing
235 efficiency. FiQA and CDR benefited from LLM-driven ontology induction, while DocRED relied
236 on its built-in schema. Rule-based information extraction proved scalable and cost-effective, with
237 all datasets processed for under \$3. FiQA emphasized scale, producing thousands of mentions from
238 a large corpus; DocRED highlighted structural challenges and type misclassification; and CDR
239 demonstrated domain adaptation with high-quality entities but sparse relations. The LLM-as-a-Judge
240 framework consistently provided actionable insights, exposing different error modes in each domain.
241 The financial graph achieved “fair” quality at 2.97/5, the DocRED graph slightly lower at 2.68/5, and
242 the biomedical graph the strongest at 3.91/5. These results validate the modularity of the framework
243 and underscore the importance of balancing rule-based precision, schema quality, and semantic
244 validation in cross-domain KG construction.

245 **4.5 Comparative Summary**

246 To illustrate the differences across datasets, Table 1 summarizes the key outcomes of our experiments.

Dataset	Docs	Nodes	Edges	Density	Cost (\$)	Quality
FiQA	5,500+	475	36	0.0002	2.63	2.97/5
DocRED	100	5,000	389	0.00003	0.15	2.68/5
CDR	100 (of 1,500)	966	13	0.0002	0.65	3.91/5

Table 1: Comparative results of KG construction and validation across FiQA, DocRED, and CDR. Costs reflect total pipeline execution, including LLM-based ontology induction (if used) and LLM-as-a-Judge validation.

247 **5 Discussion**

248 The results across FiQA, DocRED, and CDR demonstrate the robustness and adaptability of our
249 hybrid pipeline for knowledge graph construction. By combining LLM-driven ontology induction
250 with scalable rule-based extraction and entity resolution, the system consistently produced structured
251 knowledge graphs across domains as different as financial text, general document-level relations, and
252 biomedical literature. At the same time, the experiments reveal distinct strengths and weaknesses that
253 highlight where future refinements are most needed.

254 In the financial domain, the FiQA experiments showed that LLM-guided ontology induction pro-
255 vides a rich conceptual schema at low cost, enabling subsequent rule-based components to process
256 thousands of documents efficiently. The high resolution rate of 95.9% illustrates the strength of
257 type-aware similarity clustering for entity canonicalization. However, the final graph was sparse,
258 with only 36 edges among 475 nodes, and the quality assessment pointed to recurring issues of entity
259 type misclassification and vague labels. These findings suggest that while financial concepts are
260 well captured at the ontology level, relation extraction patterns require refinement to ensure more
261 meaningful connectivity in the final graph.

262 The DocRED experiments revealed a different limitation. Although the pipeline produced a graph
263 with 5,000 nodes and 389 edges, validation exposed a systemic failure in entity classification: nearly
264 all entities were labeled as PERSON. This misclassification drastically reduced entity-level quality
265 scores, bringing the overall quality down to 2.68/5 despite moderately better relation scores. Unlike
266 FiQA and CDR, where ontology induction established a balanced schema, DocRED relied on its
267 predefined annotation types, which proved difficult to capture with purely rule-based extraction.
268 This points to a weakness of relying exclusively on handcrafted patterns in contexts where entity
269 diversity is high and annotations are dense. The pipeline, however, still demonstrated efficiency and
270 reproducibility, with extraction and validation completed in under 22 minutes at a cost of only \$0.15.

271 The biomedical CDR experiments highlight the pipeline’s adaptability to specialized domains. Ontology induction successfully produced a rich biomedical schema with 15 entity types and 25 relation
272 types, supporting accurate extraction of chemicals, diseases, and treatments. The resulting graph
273 was small, with 966 nodes and only 13 validated edges, but the semantic quality was notably higher
274 than in FiQA and DocRED. Validation produced an overall quality score of 3.91/5, with particularly
275 strong performance in entity recognition (4.03/5) and relation identification (3.85/5). The relatively
276 low graph-level score of 2.20/5 reflected sparsity and fragmentation, but the pipeline nevertheless
277 produced meaningful biomedical subgraphs. These results indicate that in highly technical domains,
278 LLM-driven ontology induction coupled with domain-specific extraction rules can achieve high
279 semantic precision, even if global connectivity remains weak.
280

281 Taken together, these experiments demonstrate that the hybrid design achieves a balance between
282 cost-effectiveness and semantic depth. In all three cases, the pipeline operated at minimal cost —
283 under \$3 for FiQA, \$0.15 for DocRED, and \$0.65 for CDR — while processing times remained
284 in the range of minutes rather than hours. The modularity of the framework allowed it to adapt
285 seamlessly across domains, skipping ontology induction where unnecessary, and tailoring rule-based
286 extraction to dataset-specific schemas. Most importantly, the LLM-as-a-Judge validation framework
287 consistently exposed systematic weaknesses, whether entity type misclassification in DocRED, sparse
288 relation connectivity in FiQA, or graph-level fragmentation in CDR. This validates the use of an
289 external LLM evaluator not only for scoring quality but also for providing actionable insights that
290 guide iterative refinement.

291 The broader implication is that hybrid architectures are particularly well suited to automated knowl-
292 edge graph construction. Purely LLM-driven pipelines may be prohibitively expensive at scale, while
293 purely rule-based systems lack the conceptual flexibility to adapt across domains. By combining the
294 two, our approach demonstrates portability, reproducibility, and efficiency, while producing quality
295 scores that reveal a clear trajectory for improvement. In financial and biomedical domains, better
296 relation extraction and tighter integration with entity resolution could increase graph connectiv-
297 ity. In document-level tasks such as DocRED, augmenting rule-based extraction with lightweight
298 LLM-based classification may correct systematic misclassification errors. In all cases, the feedback
299 loop provided by the LLM-as-a-Judge can serve as a foundation for semi-automatic refinement of
300 extraction rules.

301 Overall, the experiments confirm that end-to-end knowledge graph construction and validation can
302 be achieved in a cost-effective and scalable manner across domains. The framework succeeds in
303 producing domain-specific graphs with varying degrees of precision and connectivity, while the
304 validation stage ensures that errors are not only measured but also interpreted. This positions the
305 pipeline as both a practical system for applied KG construction and a methodological contribution for
306 research in hybrid approaches to information extraction.

307 6 Conclusion and Future Work

308 This paper introduced a hybrid framework for the automated construction and validation of knowledge
309 graphs, combining LLM-driven ontology induction with scalable rule-based methods for extraction,
310 resolution, and assembly. Evaluations across three benchmark datasets — FiQA in the financial
311 domain, DocRED for document-level relation extraction, and CDR in the biomedical domain —
312 demonstrate both the versatility and the limitations of this approach. In all cases, the pipeline
313 produced structured knowledge graphs at minimal cost, with total expenses ranging from only \$0.15
314 for DocRED to under \$3 for FiQA. Processing times were consistently short, measured in minutes
315 rather than hours, confirming the efficiency of the system.

316 The experiments revealed complementary strengths and weaknesses across domains. In FiQA,
317 ontology induction successfully captured a rich conceptual schema, but relation connectivity was
318 sparse and entity classification errors persisted. In DocRED, dense annotations exposed the limits
319 of purely rule-based extraction, with entity type misclassification dominating the error profile. In
320 contrast, the biomedical CDR dataset showed that domain-specific ontology induction and tailored
321 extraction rules could yield high semantic quality, even if graph-level sparsity remained a challenge.
322 The LLM-as-a-Judge validation framework proved invaluable across all settings, not only quantifying
323 node, edge, and graph quality, but also surfacing systematic issues such as vague entity labels,

324 semantic inconsistencies, or structural fragmentation. These results validate the design choice of
325 combining LLM reasoning with rule-based scalability, producing actionable insights at low cost.
326 Future work will build on these findings in three directions. First, improvements in entity classification
327 and relation extraction are necessary to enhance precision and reduce sparsity, particularly in financial
328 and general knowledge corpora. This may include expanding the pattern libraries, integrating
329 contextual cues, and leveraging lightweight LLM components for disambiguation. Second, tighter
330 integration between entity resolution and relation extraction should preserve more valid connections,
331 increasing graph connectivity without inflating noise. Third, we plan to extend the role of LLMs
332 beyond ontology induction and validation to include semi-automatic refinement of extraction rules
333 and support for entity canonicalization. Additional avenues include multi-model validation, cross-
334 domain adaptation beyond the three benchmarks studied here, and the development of new quality
335 metrics that capture not only correctness but also coverage and coherence.
336 In conclusion, the presented framework establishes a cost-effective, modular, and adaptable pipeline
337 for knowledge graph construction. By demonstrating its applicability across financial, general,
338 and biomedical domains, and by coupling construction with explicit validation, we provide both a
339 practical methodology and a foundation for future research. The system offers a clear path toward
340 richer, more connected, and more accurate knowledge graphs, while remaining accessible in terms of
341 computational cost and reproducibility.

342 References

- 343 [1] J. Devlin, M. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional
344 transformers for language understanding. In *Proceedings of NAACL-HLT (Volume 1: Long and
345 Short Papers)*, pages 4171–4186, 2019.
- 346 [2] O. Etzioni, M. J. Cafarella, D. Downey, S. Kok, A. Popescu, T. Shaked, S. Soderland, D. S.
347 Weld, and A. Yates. Web-scale information extraction in knowitall: (preliminary results). In
348 *Proceedings of the International World Wide Web Conference (WWW)*, pages 100–110, 2004.
- 349 [3] Z. Han, P. Chen, Y. Ma, and V. Tresp. Explainable subgraph reasoning for forecasting on
350 temporal knowledge graphs. In *ICLR*, 2021.
- 351 [4] Z. Huang, W. Xu, and K. Yu. Bidirectional lstm-crf models for sequence tagging. *CoRR*,
352 abs/1508.01991, 2015.
- 353 [5] T. Jiang, T. Zhao, B. Qin, T. Liu, N. V. Chawla, and M. Jiang. Multi-input multi-output sequence
354 labeling for joint extraction of fact and condition tuples from scientific text. In *Proceedings of
355 EMNLP-IJCNLP*, pages 302–312, 2019.
- 356 [6] X. Ma and E. H. Hovy. End-to-end sequence labeling via bi-directional lstm-cnns-crf. In
357 *Proceedings of ACL (Volume 1: Long Papers)*, 2016.
- 358 [7] H. Paulheim. Knowledge graph refinement: A survey of approaches and evaluation methods.
359 *Semantic Web*, 8(3):489–508, 2017.
- 360 [8] X. Ren, W. He, M. Qu, C. R. Voss, H. Ji, and J. Han. Label noise reduction in entity typing by
361 heterogeneous partial-label embedding. In *Proceedings of KDD*, pages 1825–1834, 2016.
- 362 [9] S. Riedel, L. Yao, and A. McCallum. Modeling relations and their mentions without labeled
363 text. In *ECML PKDD 2010, Proceedings, Part III*, volume 6323, pages 148–163, 2010.
- 364 [10] J. Tang, D. Zhang, and L. Yao. Social network extraction of academic researchers. In *Proceed-
365 ings of ICDM*, pages 292–301, 2007.
- 366 [11] X. Wang, X. Han, Y. Lin, Z. Liu, and M. Sun. Adversarial multi-lingual neural relation
367 extraction. In *Proceedings of COLING*, pages 1156–1166, 2018.
- 368 [12] W. Xin-Dong, S. Shao-Jing, J. Ting-Ting, B. Chen-Yang, and W. Ming-Hui. Huapu-cp: from
369 knowledge graphs to a data central-platform. *Acta Automatica Sinica*, 46(10):2045–2059, 2020.
- 370 [13] P. Xu and D. Barbosa. Neural fine-grained entity type classification with hierarchy-aware loss.
371 In *Proceedings of NAACL-HLT (Long Papers)*, pages 16–25, 2018.

- 372 [14] J. Yan, C. Wang, W. Cheng, M. Gao, and A. Zhou. A retrospective of knowledge graphs.
373 *Frontiers of Computer Science*, 12(1):55–74, 2018.
- 374 [15] Y. Yao, D. Ye, P. Li, X. Han, Y. Lin, Z. Liu, Z. Liu, L. Huang, J. Zhou, and M. Sun. Docred: A
375 large-scale document-level relation extraction dataset. In *Proceedings of ACL (Volume 1: Long
376 Papers)*, pages 764–777, 2019.
- 377 [16] A. Yates, M. Banko, M. Broadhead, M. J. Cafarella, O. Etzioni, and S. Soderland. Textrunner:
378 Open information extraction on the web. In *HLT-NAACL, Proceedings*, pages 25–26, 2007.
- 379 [17] D. Zeng, K. Liu, Y. Chen, and J. Zhao. Distant supervision for relation extraction via piecewise
380 convolutional neural networks. In *Proceedings of EMNLP*, pages 1753–1762, 2015.
- 381 [18] F. Zhang, X. Liu, J. Tang, Y. Dong, P. Yao, J. Zhang, X. Gu, Y. Wang, B. Shao, R. Li, and
382 K. Wang. Oag: Toward linking large-scale heterogeneous entity graphs. In *Proceedings of
383 KDD*, pages 2585–2595, 2019.
- 384 [19] P. Zhou, W. Shi, J. Tian, Z. Qi, B. Li, H. Hao, and B. Xu. Attention-based bidirectional long
385 short-term memory networks for relation classification. In *Proceedings of ACL (Volume 2:
386 Short Papers)*, 2016.

387 **Agents4Science AI Involvement Checklist**

- 388 1. **Hypothesis development:** Hypothesis development includes the process by which you
389 came to explore this research topic and research question. This can involve the background
390 research performed by either researchers or by AI. This can also involve whether the idea
391 was proposed by researchers or by AI.

392 Answer: **[A]**

393 Explanation: The ideas and overall process is generated by human researchers

- 394 2. **Experimental design and implementation:** This category includes design of experiments
395 that are used to test the hypotheses, coding and implementation of computational methods,
396 and the execution of these experiments.

397 Answer: **[D]**

398 Explanation: use Cursor.ai for all the code

- 399 3. **Analysis of data and interpretation of results:** This category encompasses any process to
400 organize and process data for the experiments in the paper. It also includes interpretations of
401 the results of the study.

402 Answer: **[D]**

403 Explanation: use Cursor.ai for all the evaluation code and innovai.pro AI-copilot for writing
404 the analysis

- 405 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final
406 paper form. This can involve not only writing of the main text but also figure-making,
407 improving layout of the manuscript, and formulation of narrative.

408 Answer: **[D]**

409 Explanation: all used innovai.pro AI-copilot for writing

- 410 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or
411 lead author?

412 Description: It needs very fine-grained control, otherwise very loose work overall

413 **Agents4Science Paper Checklist**

414 **1. Claims**

415 Question: Do the main claims made in the abstract and introduction accurately reflect the
416 paper's contributions and scope?

417 Answer: [Yes]

418 Justification: see introduction

419 Guidelines:

- 420 • The answer NA means that the abstract and introduction do not include the claims
421 made in the paper.
- 422 • The abstract and/or introduction should clearly state the claims made, including the
423 contributions made in the paper and important assumptions and limitations. A No or
424 NA answer to this question will not be perceived well by the reviewers.
- 425 • The claims made should match theoretical and experimental results, and reflect how
426 much the results can be expected to generalize to other settings.
- 427 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
428 are not attained by the paper.

429 **2. Limitations**

430 Question: Does the paper discuss the limitations of the work performed by the authors?

431 Answer: [Yes]

432 Justification: see the discussion section

433 Guidelines:

- 434 • The answer NA means that the paper has no limitation while the answer No means that
435 the paper has limitations, but those are not discussed in the paper.
- 436 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 437 • The paper should point out any strong assumptions and how robust the results are to
438 violations of these assumptions (e.g., independence assumptions, noiseless settings,
439 model well-specification, asymptotic approximations only holding locally). The authors
440 should reflect on how these assumptions might be violated in practice and what the
441 implications would be.
- 442 • The authors should reflect on the scope of the claims made, e.g., if the approach was
443 only tested on a few datasets or with a few runs. In general, empirical results often
444 depend on implicit assumptions, which should be articulated.
- 445 • The authors should reflect on the factors that influence the performance of the approach.
446 For example, a facial recognition algorithm may perform poorly when image resolution
447 is low or images are taken in low lighting.
- 448 • The authors should discuss the computational efficiency of the proposed algorithms
449 and how they scale with dataset size.
- 450 • If applicable, the authors should discuss possible limitations of their approach to
451 address problems of privacy and fairness.
- 452 • While the authors might fear that complete honesty about limitations might be used by
453 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
454 limitations that aren't acknowledged in the paper. Reviewers will be specifically
455 instructed to not penalize honesty concerning limitations.

456 **3. Theory assumptions and proofs**

457 Question: For each theoretical result, does the paper provide the full set of assumptions and
458 a complete (and correct) proof?

459 Answer: [NA]

460 Justification: [NA]

461 Guidelines:

- 462 • The answer NA means that the paper does not include theoretical results.

- 463 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
464 referenced.
465 • All assumptions should be clearly stated or referenced in the statement of any theorems.
466 • The proofs can either appear in the main paper or the supplemental material, but if
467 they appear in the supplemental material, the authors are encouraged to provide a short
468 proof sketch to provide intuition.

469 **4. Experimental result reproducibility**

470 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
471 perimental results of the paper to the extent that it affects the main claims and/or conclusions
472 of the paper (regardless of whether the code and data are provided or not)?

473 Answer: [Yes]

474 Justification: see github link

475 Guidelines:

- 476 • The answer NA means that the paper does not include experiments.
477 • If the paper includes experiments, a No answer to this question will not be perceived
478 well by the reviewers: Making the paper reproducible is important.
479 • If the contribution is a dataset and/or model, the authors should describe the steps taken
480 to make their results reproducible or verifiable.
481 • We recognize that reproducibility may be tricky in some cases, in which case authors
482 are welcome to describe the particular way they provide for reproducibility. In the case
483 of closed-source models, it may be that access to the model is limited in some way
484 (e.g., to registered users), but it should be possible for other researchers to have some
485 path to reproducing or verifying the results.

486 **5. Open access to data and code**

487 Question: Does the paper provide open access to the data and code, with sufficient instruc-
488 tions to faithfully reproduce the main experimental results, as described in supplemental
489 material?

490 Answer: [Yes]

491 Justification: see github link

492 Guidelines:

- 493 • The answer NA means that paper does not include experiments requiring code.
494 • Please see the Agents4Science code and data submission guidelines on the conference
495 website for more details.
496 • While we encourage the release of code and data, we understand that this might not be
497 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
498 including code, unless this is central to the contribution (e.g., for a new open-source
499 benchmark).
500 • The instructions should contain the exact command and environment needed to run to
501 reproduce the results.
502 • At submission time, to preserve anonymity, the authors should release anonymized
503 versions (if applicable).

504 **6. Experimental setting/details**

505 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
506 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
507 results?

508 Answer: [NA]

509 Justification: does not apply

510 Guidelines:

- 511 • The answer NA means that the paper does not include experiments.
512 • The experimental setting should be presented in the core of the paper to a level of detail
513 that is necessary to appreciate the results and make sense of them.

- 514 • The full details can be provided either with the code, in appendix, or as supplemental
515 material.

516 **7. Experiment statistical significance**

517 Question: Does the paper report error bars suitably and correctly defined or other appropriate
518 information about the statistical significance of the experiments?

519 Answer: **[No]**

520 Justification: all done by llm

521 Guidelines:

- 522 • The answer NA means that the paper does not include experiments.
523 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
524 dence intervals, or statistical significance tests, at least for the experiments that support
525 the main claims of the paper.
526 • The factors of variability that the error bars are capturing should be clearly stated
527 (for example, train/test split, initialization, or overall run with given experimental
528 conditions).

529 **8. Experiments compute resources**

530 Question: For each experiment, does the paper provide sufficient information on the com-
531 puter resources (type of compute workers, memory, time of execution) needed to reproduce
532 the experiments?

533 Answer: **[Yes]**

534 Justification: used Cursor.ai and innovai.pro

535 Guidelines:

- 536 • The answer NA means that the paper does not include experiments.
537 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
538 or cloud provider, including relevant memory and storage.
539 • The paper should provide the amount of compute required for each of the individual
540 experimental runs as well as estimate the total compute.

541 **9. Code of ethics**

542 Question: Does the research conducted in the paper conform, in every respect, with the
543 Agents4Science Code of Ethics (see conference website)?

544 Answer: **[Yes]**

545 Justification:

546 Guidelines:

- 547 • The answer NA means that the authors have not reviewed the Agents4Science Code of
548 Ethics.
549 • If the authors answer No, they should explain the special circumstances that require a
550 deviation from the Code of Ethics.

551 **10. Broader impacts**

552 Question: Does the paper discuss both potential positive societal impacts and negative
553 societal impacts of the work performed?

554 Answer: **[Yes]**

555 Justification: see discussion and conclusion section

556 Guidelines:

- 557 • The answer NA means that there is no societal impact of the work performed.
558 • If the authors answer NA or No, they should explain why their work has no societal
559 impact or why the paper does not address societal impact.
560 • Examples of negative societal impacts include potential malicious or unintended uses
561 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations,
562 privacy considerations, and security considerations.
563 • If there are negative societal impacts, the authors could also discuss possible mitigation
564 strategies.