

SocraticKG: Knowledge Graph Construction via QA-Driven Fact Extraction

Anonymous ACL submission

Abstract

Constructing Knowledge Graphs (KGs) from unstructured text provides a structured framework for knowledge representation and reasoning, yet current LLM-based approaches struggle with a fundamental trade-off: factual coverage often leads to relational fragmentation, while premature consolidation causes information loss. To address this, we propose SocraticKG, an automated KG construction method that introduces question-answer pairs as a structured intermediate representation to systematically unfold document-level semantics prior to triple extraction. By employing 5W1H-guided QA expansion, SocraticKG captures contextual dependencies and implicit relational links typically lost in direct KG extraction pipelines, providing explicit grounding in the source document that helps mitigate implicit reasoning errors. Evaluation on the MINE benchmark demonstrates that our approach effectively addresses the coverage-connectivity trade-off, achieving superior factual retention while maintaining high structural cohesion even as extracted knowledge volume substantially expands. These results highlight that QA-mediated semantic scaffolding plays a critical role in structuring semantics prior to KG extraction, enabling more coherent and reliable graph construction in subsequent stages.

1 Introduction

As large language models (LLMs) are widely used in knowledge-intensive applications, concerns surrounding factual reliability, interpretability, and grounding have become more pronounced (Ji et al., 2023; Huang et al., 2025). While Retrieval-Augmented Generation (RAG) addresses these concerns by anchoring models to external sources, it often struggles with fragmented contexts and shallow integration of complex facts (Lewis et al., 2020; Gao et al., 2023). In response, Knowledge Graphs (KGs) have re-emerged as a complementary solution, providing a structured and verifiable backbone

for explicit knowledge representation and reasoning (Pan et al., 2023; Rajabi and Etmnani, 2024). However, the reliance on manual curation has historically limited the availability of domain-specific KGs, thereby motivating growing interest in automated construction methods that can scale to diverse and large-scale text sources (Ren et al., 2024).

Recent advances in LLMs have enabled more semantically grounded approaches to knowledge graph construction, moving beyond rule-based pattern matching toward methods that leverage neural reasoning to interpret unstructured text (Zhu et al., 2024). Current approaches address the construction challenge through different strategies. Some methods focus on capturing explicit factual mentions in a single pass, extracting triples directly from text (Cabot and Navigli, 2021; Shang et al., 2022; Zhang and Soh, 2024). Others adopt consolidation-centric strategies, organizing extracted facts around pre-identified entity structures to improve graph coherence (Zhong and Chen, 2021; Ye et al., 2022a; Wei et al., 2023; Mo et al., 2025).

However, these approaches face a persistent challenge: fully externalizing the narrative logic of source documents into structured graphs. The resulting knowledge graphs often struggle with a fundamental tension between factual coverage and structural coherence. Graphs may contain many facts but remain fragmented with weak semantic connectivity, or they may be well-organized yet incomplete, having filtered out contextual nuances that do not conform to predefined structures. At the core of this challenge lies the difficulty of balancing comprehensive information extraction with meaningful connectivity across the graph.

To address this limitation, we draw inspiration from how humans naturally process and organize information from text. Rather than attempting to extract structured knowledge in a single step, human comprehension is fundamentally interrogative: readers construct understanding by progressively

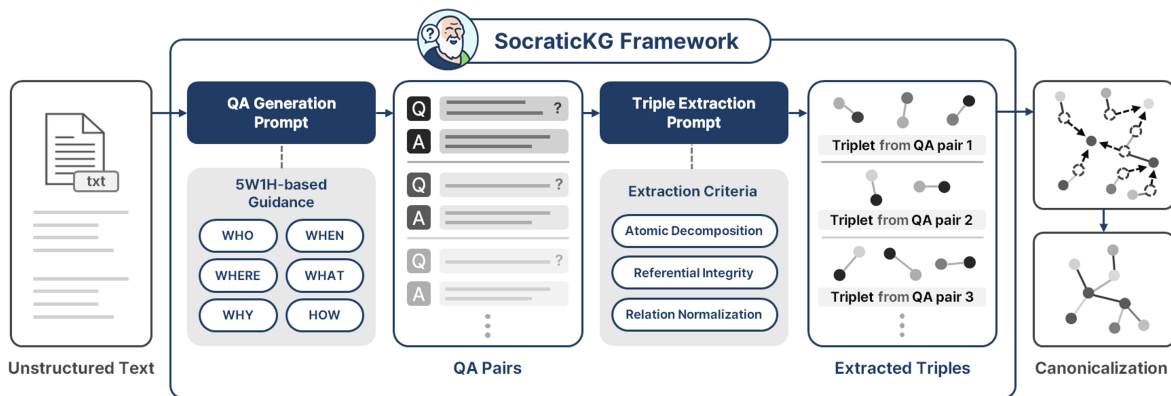


Figure 1: The overall architecture of the SocraticKG framework. Given unstructured text, the method first generates atomic QA pairs through 5W1H-guided questioning, then extracts triples from these QA pairs, and finally canonicalizes the triples to produce a cohesive knowledge graph.

clarifying salient concepts through active inquiry (Graesser and Person, 1994; Ambrose et al., 2010). This process of interrogative learning serves as a natural scaffold for organizing complex information. Question-Answering (QA), in particular, facilitates focused attention and explicit articulation of relationships that might otherwise remain implicit in direct extraction (Wu et al., 2020).

Building on this insight, we propose SocraticKG (SoKG), a method that treating QA not merely as a retrieval mechanism, but as a structured intermediate representation that systematically unfolds document-level semantics prior to graph construction (FitzGerald et al., 2018; Cohen et al., 2023). SoKG employs a structured interrogative framework based on the 5W1H framework (*who, what, when, where, why, and how*) to generate document-grounded QA pairs that capture key concepts, relationships, and contextual dependencies. This QA-mediated expansion articulates implicit connections and contextual nuances in explicit natural language format. The resulting intermediate representation facilitates more consistent and complete triple extraction by providing well-defined semantic units rather than requiring simultaneous resolution of semantics and structure. These extracted triples are then unified through a canonicalization process (Mo et al., 2025) that resolves surface-form variations and consolidates the graph into a coherent structure.

We evaluate our proposed method on the MINE (Measure of Information in Nodes and Edges) benchmark (Mo et al., 2025), a recently proposed benchmark designed to measure factual recoverability from automatically constructed knowledge graphs. Our results demonstrate that SocraticKG

consistently outperforms state-of-the-art counterparts across multiple LLM backbones, achieving superior factual retention while producing more densely connected and less fragmented graphs. By reconciling factual coverage with structural coherence, SoKG provides a scalable approach for high-fidelity KG construction and more reliable structured reasoning.

In summary, we make the following contributions in this work:

- We propose SocraticKG, a QA-mediated method for knowledge graph construction that formalizes question-answering as a semantic scaffold for unfolding document narratives and explicitly articulating implicit connections prior to structural extraction.
- We introduce 5W1H-guided QA expansion as a systematic approach for surfacing latent dependencies typically overlooked in direct extraction, thereby improving factual coverage while reducing implicit reasoning errors.
- We demonstrate that our approach mitigates structural fragmentation and information loss, achieving superior factual retention and recoverability across various LLMs.

2 Related Work

2.1 Direct Triple Extraction

Knowledge Graph (KG) construction has evolved from conventional Open Information Extraction (OpenIE) (Etzioni et al., 2008; Fader et al., 2011) to modern approaches that extract triples directly via LLMs (Cabot and Navigli, 2021; Bi et al., 2024; Zhang and Soh, 2024). While OpenIE is

154	constrained by surface linguistic patterns (Niklaus et al., 2018), such direct extraction methods leverage LLM reasoning capabilities to bridge semantic gaps without explicit intermediate representations.	204
155		205
156		206
157		207
158	However, this direct extraction approach often limits the model to capturing surface-level, explicit mentions while overlooking the latent logical ties that bind them. As noted by Zhu et al. (2024); Meher et al. (2025), this approach often yields shallow factual coverage, often producing fragmented subgraphs that lack the connectivity required for effective graph-based reasoning.	208
159		209
160		210
161		211
162		212
163		213
164		214
165		215
166	2.2 Consolidation-Centric Strategies	216
167	To address the fragmentation issues, various pipelines emphasize structural coherence through post-extraction consolidation. These entity-first approaches organize extracted facts by first identifying key entities, then structuring relations around this pre-established entity framework. GraphRAG (Edge et al., 2024) builds a global index of entities and relationships partitioned into hierarchical communities for query-focused summarization, whereas KGGen (Mo et al., 2025) emphasizes clustering-based canonicalization of entities and relations to produce compact and reusable knowledge graphs. Similarly, CLARE (Henry and Gong, 2025) anchors its relational extraction on initial entity identification to ensure semantic precision within consolidated text.	217
168		218
169		219
170		220
171		221
172		222
173		223
174		224
175		225
176		226
177		227
178		228
179		229
180		230
181		231
182		232
183	Despite their effectiveness in organizing triples, these consolidation-focused strategies can act as a representational bottleneck (Ye et al., 2022b). When entity sets are fixed early in the pipeline, relations or contextual dependencies that do not conform to the initial entity structure may be excluded. This sequencing effectively prioritizes structural utility over factual density, potentially under-representing the document’s latent relations.	233
184		234
185		235
186		236
187		237
188		238
189		239
190		240
191		241
192	2.3 Transform-Then-Extract Approaches	242
193	To reduce extraction complexity, various approaches employ a two-stage process: first transforming raw text through intermediate representations, then extracting triples from the transformed output. Common transformation strategies include coreference resolution to handle referential expressions (Manning et al., 2014; Cetto et al., 2018) and syntactic sentence decomposition to simplify complex structures (Niklaus et al., 2019; Niklaus, 2022). CoDe-KG (Anuyah et al., 2025), for instance, leverages human-guided prompt intervention to incorpo-	243
194		244
195		245
196		246
197		247
198		248
199		249
200		250
201		251
202		252
203		
	rate these transformation tasks, ensuring structural clarity prior to extraction.	
	These transformation-based approaches operate primarily at the sentence level, focusing on local syntactic normalization rather than document-level semantic organization. While effective for resolving surface-level ambiguities within individual sentences, they do not systematically capture cross-sentence dependencies or contextual relationships that span the document. This limits their ability to externalize the broader narrative structure and global semantics required for comprehensive knowledge graph construction.	
	2.4 QA for Knowledge Extraction	
	Question-answering has been widely used to elicit structured information from text (Levy et al., 2017; Li et al., 2019; Du and Cardie, 2020), by leveraging the cognitive process of interrogative inquiry, which facilitates the construction of situation models (Chi et al., 1989; Graesser and Person, 1994). Recent extraction methods, such as StoryNet (Nagireddy, 2021) and ChatIE (Wei et al., 2023), incorporate QA-driven prompting as a core component of their extraction pipelines.	
	However, these approaches treat QA pairs as transient artifacts, generating and consuming them within a single extraction pass, without formalizing them as an intermediate representation for organizing document-level semantics. As a result, they lack systematic question generation and struggle to surface implicit relational and contextual dependencies prior to triple extraction.	
	While recent work has explored QA as an intermediate step for interpretable knowledge construction (Aneja et al., 2025), it primarily emphasizes retrieval utility through factual restatement, rather than semantic organization. Collectively, these gaps suggest that formalizing QA as a structured intermediate representation provides a more robust foundation for construction, particularly when systematic inquiry is used to proactively externalize latent relational and causal dependencies.	
	3 Methods	
	SoKG introduces QA pairs as a structured intermediate representation for LLM-based KG construction. Rather than prompting LLMs to extract triples directly from raw text, our approach first decomposes the document into explicit QA pairs that resolve contextual dependencies and referential am-	

253	biguities in natural language. These QA pairs are	Simplified Relations Predicates are distilled into	299
254	then mapped to atomic triples and unified through	concise verb phrases to reduce surface-form varia-	300
255	canonicalization to produce the final KG.	tions, facilitating subsequent canonicalization. The	301
		model is instructed to skip extraction if the relation-	302
256	3.1 5W1H-Guided QA Generation	ship remains ambiguous.	303
257	This stage transforms document into a collection of	3.3 Graph Construction from Triples	304
258	discrete, self-contained QA pairs. To ensure com-	The final stage unifies discrete triples into a cohe-	305
259	prehensive coverage of the factual content in the	sive graph structure. Since extraction occurs across	306
260	text, we design a prompt strategy based on two core	independent QA units, the raw set often contains re-	307
261	principles: systematic questioning and contextual	dundant or synonymous mentions for the same con-	308
262	independence (detailed prompt in Appendix B.1).	cept. To resolve these redundancies, we adopt the	309
263	Detailed Questioning via 5W1H We leverage	canonicalization procedure from Mo et al. (2025),	310
264	the 5W1H framework to guide systematic ques-	which combines embedding-based clustering with	311
265	tion formulation. The LLM generates multiple	LLM-based refinement.	312
266	questions spanning all six categories and diverse	The canonicalization process is performed inde-	313
267	aspects of the document. As a result, the resulting	pendently on entities and relations through a	314
268	QA pairs capture both surface-level entities and	cluster-then-refine process. First, semantic embed-	315
269	complex dependencies, including causal rationales	dings are generated for all unique entities and rela-	316
270	(<i>why</i>) and procedural details (<i>how</i>).	tions using a text embedding model. To narrow the	317
271	Contextual Independence To ensure each QA	search space, these embeddings are partitioned into	318
272	pair functions as a standalone unit, we instruct the	clusters of a manageable size for entities and rela-	319
273	LLM to generate answers that are fully understand-	tions respectively via K-means clustering. Within	320
274	able without referencing the original source text.	each cluster, the top- <i>k</i> potential matches for each	321
275	Specifically, the model is required to replace pro-	anchor are identified by balancing dense seman-	322
276	nouns (<i>e.g., it, they</i>) with their explicit entity names,	tic similarity with sparse lexical overlap (BM25).	323
277	resolving referential ambiguities. This constraint	Finally, synonyms and abbreviations are resolved	324
278	prevents information loss when each QA pair is pro-	by an LLM, which maps these variants to a sin-	325
279	cessed individually in the triple extraction phase.	gle representative entity or relation to consolidate	326
280	3.2 Triple Extraction from QA	fragmented triples into a cohesive, canonicalized	327
281	This stage transforms the QA pairs into structured	graph.	328
282	triples by treating each pair as an independent ex-	4 Experiments	329
283	traction unit. By operating on these logically self-	To validate the effectiveness of SoKG, we utilized	330
284	contained units allows the extraction process to fo-	the MINE benchmark (Mo et al., 2025). MINE was	331
285	cus on well-defined semantic boundaries, reducing	designed to quantify the information gap between	332
286	errors common in direct extraction from long, com-	raw text and its graph representation by measuring	333
287	plex texts. To achieve this, the LLM is instructed to	how much source information is recoverable. The	334
288	follow three specific constraints (detailed prompt	benchmark comprises 100 diverse articles, each	335
289	in Appendix C.1).	paired with 15 verified atomic facts, providing a	336
290	Atomic Decomposition The model decomposes	rigorous evaluation framework across 1,500 inde-	337
291	each QA pair into separate, atomic triples, captur-	pendent factual instances. Following the bench-	338
292	ing fine-grained facts from both the inquiry and the	mark protocol, we constructed one KG per article	339
293	response to maximize factual richness.	and evaluated each graph from two complemen-	340
294	Entity Clarity All entities are expressed as spe-	tary perspectives: factual retention and structural	341
295	cific noun phrases, and any triple containing am-	characteristics.	342
296	biguous pronouns is discarded. This ensures that	4.1 Evaluation Metrics	343
297	every extracted fact is self-contained and grounded	Factual Retention Score As the primary metric,	344
298	in clear evidence.	we measured the proportion of ground-truth facts	345
		successfully recovered from the constructed KGs.	346

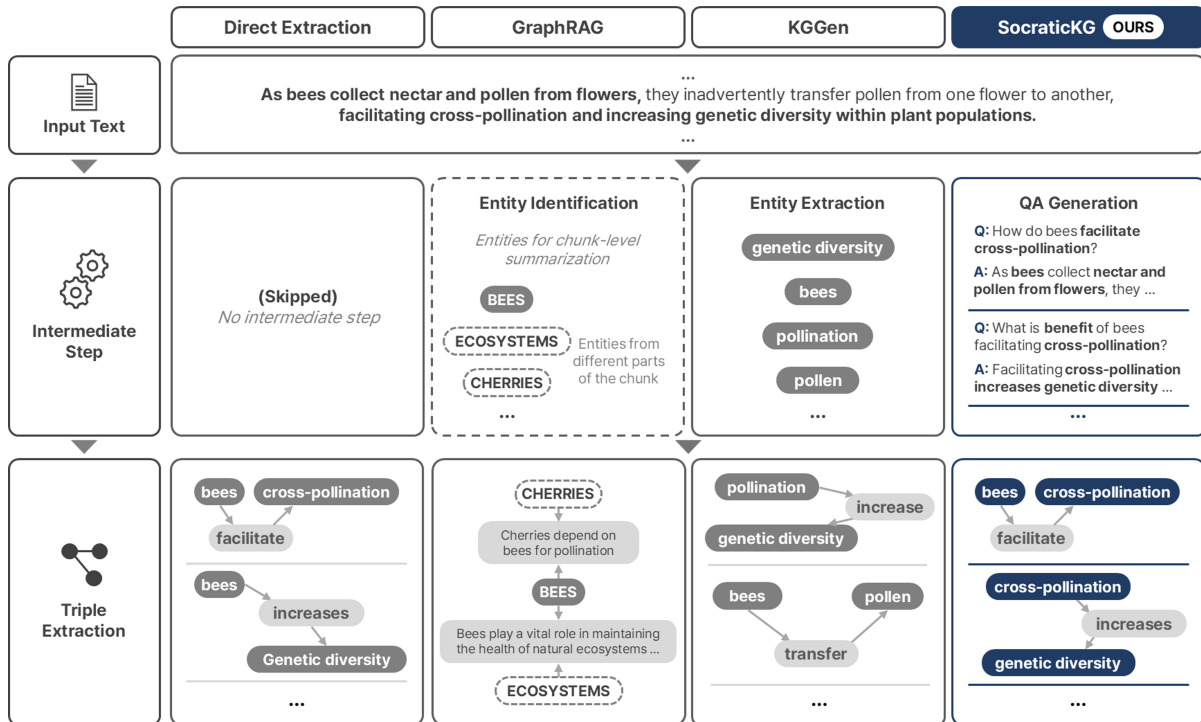


Figure 2: Comparison of extraction pipelines using an example output from Gemini-2.5-flash-lite. While baseline pipelines often miss the syntactic connection in complex sentences, failing to recover the causal link between bees and genetic diversity, SoKG leverages QA-driven reasoning to explicitly reconstruct the intermediate concept. As a result, SoKG successfully recovers the complete causal chain (bees → cross-pollination → genetic diversity), whereas baselines tend to simplify or fragment this relationship.

347 Following the MINE benchmark protocol, we re- 369
 348 trieved a local subgraph for each fact, consisting 370
 349 of the top-8 nodes most semantically similar to 371
 350 the target statement and their 2-hop neighbors. An 372
 351 LLM-judge then determined whether the fact was 373
 352 logically supported by the retrieved subgraph con- 374
 353 text. The score represents the percentage of veri- 375
 354 fiable facts, reflecting how well the graph preserves 376
 355 information from the source text for downstream 377
 356 tasks such as retrieval and reasoning.

357 **Structural Cohesion and Density** To analyze 378
 358 the organization and coherence of the KGs, we 379
 359 investigated:

- 360 • **Average Degree (Deg):** The average number 380
 361 of unique neighboring nodes per node, cap- 381
 362 turing the local connectivity density of the 382
 363 graph (Barabási, 2013). It reflects how many 383
 364 distinct entities a node is connected to, irre- 384
 365 spective of relation direction. Following stan- 385
 366 dard practice for undirected graphs, we com-
 367 pute the average degree as

$$368 \text{Deg} = \frac{2E}{N},$$

where N denotes the number of nodes and E 369
 the number of edges. 370

- **Triple Count (#Tri):** The total number of 371
 atomic facts externalized in the graph. 372
- **Normalized Fragmentation Index (NFI):** 373
 Motivated by the notion of graph fragmen- 374
 tation as the decomposition of a network into 375
 disconnected components (Borgatti, 2003), 376
 we define a component-based metric as: 377

$$378 \text{NFI} = \frac{C - 1}{N - 1},$$

379 where C denotes the number of connected 380
 380 components and N is the total number of 381
 381 nodes ($N \geq 2$). This formulation normal- 382
 382 izes fragmentation to the unit interval $[0, 1]$, 383
 383 where 0 corresponds to a fully connected 384
 384 graph ($C = 1$) and 1 indicates a completely 385
 385 fragmented network ($C = N$).

386 4.2 Comparative Analysis Design

387 **KG Construction Methods** We compared 387
 388 SoKG against three representative approaches in 388

Method	Qwen-2.5	GPT-4o-mini	GPT-4o	Gemini-2.5	Claude-4
Direct Extraction	66.5	68.5	78.1	84.6	86.8
GraphRAG	59.7	49.5	49.3	48.5	52.3
KGGen	56.7	44.3	66.4	62.5	69.1
SoKG (w/o 5W1H)	67.1	80.5	83.5	85.6	94.6
SoKG (Ours)	73.4	83.9	89.3	87.7	96.3

Table 1: Comparison of factual retention scores (%) on the MINE benchmark. SoKG consistently achieves the highest performance across all evaluated models. The vanilla variant (*i.e.*, SoKG w/o 5W1H) shows how the 5W1H scaffold captures procedural and causal facts to improve factual consistency even on smaller models like Qwen-2.5.

the current landscape of LLM-based KG construction. To ensure a valid comparison, we selected the comparative methods that operated in autonomous and open-domain settings without pre-defined schemas or human intervention. Figure 2 summarizes the procedures of these methods.

- **Direct Extraction:** A single-pass extraction strategy where triples are generated directly from raw text (Appendix C.2). For fair comparison, we apply the identical canonicalization procedure used in KGGen and SoKG to consolidate the extracted triples. It serves as a primary benchmark for the LLM’s implicit reasoning capability without the benefit of intermediate semantic scaffolding.
- **GraphRAG:** A prominent solution across industry and academia for global, query-focused entity indexing. We utilize Microsoft’s official implementation for hierarchical community detection and aggregation, providing a benchmark against the widely adopted text-summary-based method.
- **KGGen:** A recent state-of-the-art method focusing on entity-centric extraction and structural consolidation. It serves as a primary comparative method for evaluating factual retention and structural cohesion in open-domain.
- **SoKG (w/o 5W1H):** An ablated variant of SoKG that retains QA pairs as its intermediate representation but replaces the 5W1H-guided inquiry with generic QA. This design isolates the contribution of the 5W1H-guided scaffold to evaluate its impact.
- **SoKG (Ours):** Our proposed approach utilizing 5W1H-guided QA generation to systematically construct KGs from source documents. Unless otherwise specified, SoKG refers to this complete implementation.

Evaluation across LLMs To assess robustness across varying LLM architectures and scales, we evaluated the selected KG construction methods on five LLMs: GPT-4o, GPT-4o-mini, Gemini-2.5 (Gemini-2.5-Flash-Lite), Qwen-2.5 (Qwen2.5-7B-Instruct), and Claude-4 (Claude-4-Sonnet).

4.3 Implementation Details

For all LLMs, we set the decoding temperature to 0 to ensure reproducibility, except for GraphRAG, which follows the default stochastic configuration of its official implementation.

We adopted the canonicalization and factual retention evaluation protocol proposed by Mo et al. (2025). For the semantic clustering mentioned in Section 3.3, we partitioned entities and relations into clusters containing at most 128 elements. For the identification of potential matches, we set the candidate retrieval size to $k = 16$, which defines the number of top-ranked duplicates evaluated by the LLM. All embedding-based processes used the all-MiniLM-L6-v2 model, and factual verification was performed via an LLM-as-a-judge protocol using GPT-4o.

5 Results and Discussion

5.1 Factual Retention Performance

Table 1 summarizes the factual retention performance on the MINE benchmark. Across all compared methods and evaluated LLMs, SoKG consistently achieves the highest scores, peaking at 96.3% with Claude-4.

The comparison between Direct Extraction and SoKG (w/o 5W1H) illustrates the benefit of introducing QA as an intermediate representation. Even without 5W1H guidance, SoKG outperforms Direct Extraction on all LLM models. This advantage stems from decomposing documents into discrete, self-contained QA pairs prior to triples extraction.

Method	Qwen-2.5			GPT-4o-mini			GPT-4o			Gemini-2.5			Claude-4		
	N	E	Deg	N	E	Deg	N	E	Deg	N	E	Deg	N	E	Deg
Direct Extraction	21.7	17.3	1.60	33.5	28.1	1.69	33.9	27.4	1.62	58.4	64.1	2.20	46.4	40.8	1.77
GraphRAG	19.8	19.0	2.00	11.2	10.2	1.84	11.3	9.70	1.75	15.4	17.7	2.35	14.6	16.2	2.20
KGGen	28.1	22.1	1.56	19.3	16.7	1.75	33.2	28.9	1.74	38.1	43.2	2.23	57.2	58.9	2.07
SoKG (w/o 5W1H)	28.0	25.4	1.81	49.2	50.5	2.06	51.9	49.1	1.89	58.0	67.8	2.34	84.2	94.5	2.25
SoKG (Ours)	34.9	34.1	1.96	57.9	62.2	2.16	62.3	60.5	1.95	65.7	80.8	2.47	104.2	128.4	2.48

Table 2: Topological characteristics averaged over the 100 articles in the MINE benchmark. N, E, and Deg denote the mean count of Nodes, Edges, and Average Degree per graph, respectively. SoKG consistently expands the knowledge scale while maintaining high connectivity density across all backbones.

Method	Qwen-2.5		GPT-4o-mini		GPT-4o		Gemini-2.5		Claude-4	
	NFI	#Tri	NFI	#Tri	NFI	#Tri	NFI	#Tri	NFI	#Tri
Direct Extraction	0.162	1,955	0.145	3,100	0.172	2,941	0.038	7,315	0.127	4,417
GraphRAG	0.084	1,981	0.038	1,076	0.083	1,009	0.036	1,848	0.067	1,590
KGGen	0.187	2,375	0.091	1,942	0.112	3,089	0.030	5,301	0.052	6,391
SoKG (w/o 5W1H)	0.106	2,871	0.059	5,646	0.092	5,345	0.034	7,875	0.056	10,511
SoKG (Ours)	0.078	3,958	0.047	7,069	0.086	6,627	0.023	9,612	0.039	14,849

Table 3: Comparison of graph fragmentation averaged over the 100 articles (NFI; lower is better) and total extracted information volume summed over the 100 articles (#Tri). The results demonstrate that SoKG effectively resolves the trade-off between knowledge coverage and structural connectivity, maintaining high graph cohesion even as the volume of extracted facts increases.

Notably, both GraphRAG and KGGen underperform Direct Extraction in terms of factual retention. GraphRAG prioritizes hierarchical community structures and query-focused summarization over comprehensive fact preservation, resulting in lower coverage of atomic facts. KGGen’s entity-first bottleneck similarly leads to fact omission when initial entity identification fails, showing inconsistent performance across models.

In contrast, SoKG with 5W1H guidance further enhances performance by systematically surfacing procedural and causal dimensions. This interrogative framework ensures that latent dependencies are explicitly captured, maintaining high factual consistency regardless of the underlying LLM model’s inherent reasoning capacity.

To further validate our triple extraction strategy, we conducted additional experiments in Appendix A. By isolating the impact of the QA scaffold from the extraction strategy, these studies reveal that entity-first approaches persist as a performance bottleneck even when applied to QA-preprocessed inputs.

5.2 Graph Scale and Connectivity

The superior factual retention shown in Table 1 raises a critical question: is SoKG simply extracting more triples, or is it building a fundamentally

better graph? To address this, we examine graph scale and connectivity in Table 2.

SoKG significantly expands graph scale while maintaining or improving connectivity density across all evaluated LLMs. In contrast, Direct Extraction produces smaller graphs with lower connectivity, while GraphRAG generates compact structures that sacrifice comprehensive fact coverage for hierarchical organization. The comparison between SoKG (w/o 5W1H) and SoKG reveals that 5W1H guidance substantially increases the number of extracted entities and relations while enhancing connectivity density. This indicates that 5W1H systematically surfaces additional facts without fragmenting the graph structure.

Importantly, SoKG achieves higher connectivity than both Direct Extraction and KGGen despite using the same canonicalization procedure. This confirms that the structural advantage originates from the QA-mediated intermediate representation, enabling relevant evidence to co-locate within 2-hop neighborhoods and directly supporting the high fact recoverability in Table 1.

5.3 Factual Volume and Structural Cohesion

To further examine the relationship between knowledge coverage and graph fragmentation, we analyze triple counts and the NFI in Table 3. SoKG

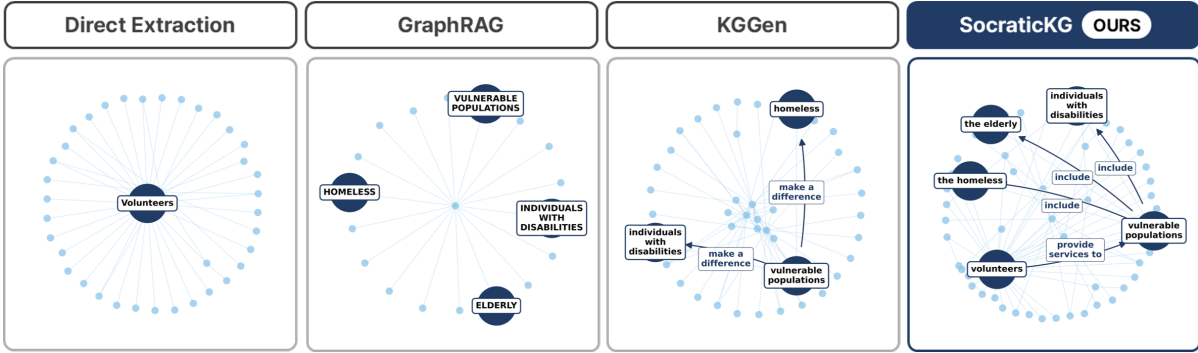


Figure 3: Comparison of extracted graphs for the example sentence: “Volunteers provide essential services and support to vulnerable populations, such as the homeless, the elderly, and individuals with disabilities.” The nested relational path implied by this text ($Volunteers \rightarrow Vulnerable\ Populations \rightarrow \{Homeless, Elderly, Individuals\ with\ disabilities\}$) is emphasized to assess relational completeness. Specifically, nodes corresponding to this path are enlarged for clear visibility, connected by thick dark blue arrows to indicate the sequence of triples, while the remaining background graph elements are displayed in light blue.

substantially expands knowledge volume while simultaneously reducing fragmentation across all evaluated LLMs. While alternative methods either limit fact extraction (GraphRAG) or exhibit higher fragmentation (KGGen and Direct Extraction), SoKG extracts substantially more triples while maintaining lower NFI values.

The comparison between SoKG (w/o 5W1H) and SoKG (Ours) further illustrates the effectiveness of 5W1H guidance. Adding 5W1H consistently increases triple extraction volume while reducing or maintaining similar fragmentation levels. This pattern indicates that 5W1H not only surfaces additional facts but also enhances their integration into the graph structure.

5.4 Qualitative Analysis

The cases in Figures 2 and 3 provide concrete examples of how SoKG’s interrogative process resolves the structural deficiencies and information loss observed in alternative methods.

Figure 2 illustrates SoKG’s capacity to preserve logical coherence in complex participle phrases, such as *facilitating cross-pollination*. While other approaches produce fragmented or oversimplified triples, SoKG articulates mediating concepts to ensure a cohesive causal chain: *bees* \rightarrow *cross-pollination* \rightarrow *genetic diversity*.

Similarly, Figure 3 demonstrates how SoKG resolves relational fragmentation in nested entity structures. Alternative methods often omit key nodes or overlook entities like *the elderly*, whereas SoKG fully reconstructs the relational tree by identifying all key entities and linking them via precise

predicates such as *provide services to* and *include*.

These examples demonstrate that QA-mediated semantic scaffolding, guided by 5W1H inquiry, systematically addresses both causal reconstruction and relational completeness, enabling more structured knowledge extraction.

6 Conclusion

We present SoKG, LLM-based KG construction method that uses QA pairs as a structured intermediate representation for document-level semantic expansion prior to triple extraction. By employing 5W1H-guided QA generation, SoKG resolves referential ambiguities and surfaces implicit relational dependencies, ensuring that subsequent structural mapping is grounded in explicit, contextualized entities rather than underspecified inferences.

Evaluation on the MINE benchmark demonstrates that SoKG achieves superior factual coverage while simultaneously improving structural cohesion across diverse LLMs. This performance stems from the QA-mediated scaffold, which systematically externalizes latent causal and relational dependencies that enhance graph connectivity even as the volume of extracted facts increases.

Our findings indicate that explicit semantic organization through QA generation is not merely an auxiliary preprocessing step but a critical component for maintaining graph fidelity in LLM-based construction. By addressing the inherent trade-off between factual coverage and structural connectivity, SoKG provides a more reliable foundation for document-grounded knowledge representation and structured reasoning.

584 Limitations

585 While SoKG enables the construction of dense
586 knowledge graphs, the QA-mediated pipeline nat-
587 urally involves higher token consumption and la-
588 tency than direct triple extraction. We prioritize
589 factual density over cost optimization, though effi-
590 ciency remains a target for future refinement. Fur-
591 thermore, as the graph quality depends on the rea-
592 soning depth of the underlying LLM, performance
593 may vary in domains requiring highly specialized
594 interrogative logic.

595 Regarding graph representation, our current use
596 of binary triples may simplify multidimensional
597 qualifiers (e.g., temporal or spatial data) that could
598 be more compactly encoded via n-ary relations. Fi-
599 nally, our evaluation focuses on factual recoverabil-
600 ity through the MINE benchmark. While this aligns
601 with our objective of preserving document-level
602 semantics, other dimensions—such as schema-
603 alignment and relation-type fidelity—are left as
604 promising avenues for the community to explore
605 as KG evaluation standards evolve.

606 Ethical Considerations

607 This study utilizes the publicly available MINE
608 benchmark and LLMs. We acknowledge that the
609 benchmark and underlying LLMs may possess in-
610 herent biases, which could be reflected in the con-
611 structed graphs. Additionally, automated extraction
612 carries a risk of hallucinating facts not present in
613 source documents. We recommend human verifica-
614 tion and validation for applications in sensitive or
615 high-stakes domains.

616 References

617 Susan A Ambrose, Michael W Bridges, Michele DiPi-
618 etro, Marsha C Lovett, and Marie K Norman. 2010.
619 *How learning works: Seven research-based princi-
620 ples for smart teaching*. John Wiley & Sons.

621 Kartikeya Aneja, Manasvi Srivastava, Subhayan Das,
622 and Nagender Aneja. 2025. Interpretable question
623 answering with knowledge graphs. *arXiv preprint*
624 *arXiv:2510.19181*.

625 Sydney Anuyah, Mehedi Mahmud Kaushik, Sri Rama
626 Krishna Reddy Dwarampudi, Rakesh Shiradkar, Ar-
627 jan Durresti, and Sunandan Chakraborty. 2025. Au-
628 tomated knowledge graph construction using large
629 language models and sentence complexity modelling.
630 In *Proceedings of the 2025 Conference on Empiri-
631 cal Methods in Natural Language Processing*, page
632 15526–15550. Association for Computational Lin-
633 guistics.

Albert-László Barabási. 2013. Network science. *634*
Philosophical Transactions of the Royal Society A:
635 Mathematical, Physical and Engineering Sciences,
636 371(1987):20120375. 637

Zhen Bi, Jing Chen, Yinuo Jiang, Feiyu Xiong, Wei Guo,
638 Huajun Chen, and Ningyu Zhang. 2024. Codekgc:
639 Code language model for generative knowledge
640 graph construction. *ACM Transactions on Asian*
641 and Low-Resource Language Information Process-
642 ing, 23(3):1–16. 643

Stephen P Borgatti. 2003. *The key player problem*. na. 644

Pere-Lluís Hugué Cabot and Roberto Navigli. 2021. 645
Rebel: Relation extraction by end-to-end language
646 generation. In *Findings of the association for compu-*
647 tational linguistics: emnlp 2021, pages 2370–2381. 648

Matthias Cetto, Christina Niklaus, André Freitas,
649 and Siegfried Handschuh. 2018. Graphene: a
650 context-preserving open information extraction sys-
651 tem. *arXiv preprint arXiv:1808.09463*. 652

Micheline TH Chi, Miriam Bassok, Matthew W Lewis,
653 Peter Reimann, and Robert Glaser. 1989. Self-
654 explanations: How students study and use examples
655 in learning to solve problems. *Cognitive science*,
656 13(2):145–182. 657

William W Cohen, Wenhui Chen, Michiel De Jong,
658 Nitish Gupta, Alessandro Presta, Pat Verga, and
659 John Wieting. 2023. Qa is the new kr: Question-
660 answer pairs as knowledge bases. In *Proceedings of*
661 the AAAI Conference on Artificial Intelligence, vol-
662 ume 37, pages 15385–15392. 663

Xinya Du and Claire Cardie. 2020. Event extraction by
664 answering (almost) natural questions. In *Proceedings*
665 of the 2020 Conference on Empirical Methods in
666 Natural Language Processing (EMNLP), pages 671–
667 683. 668

Darren Edge, Ha Trinh, Newman Cheng, Joshua
669 Bradley, Alex Chao, Apurva Mody, Steven Truitt,
670 Dasha Metropolitanaky, Robert Osazuwa Ness, and
671 Jonathan Larson. 2024. From local to global: A
672 graph rag approach to query-focused summarization.
673 *arXiv preprint arXiv:2404.16130*. 674

Oren Etzioni, Michele Banko, Stephen Soderland, and
675 Daniel S Weld. 2008. Open information extrac-
676 tion from the web. *Communications of the ACM*,
677 51(12):68–74. 678

Anthony Fader, Stephen Soderland, and Oren Etzioni.
679 2011. Identifying relations for open information ex-
680 traction. In *Proceedings of the 2011 conference on*
681 empirical methods in natural language processing,
682 pages 1535–1545. 683

Nicholas FitzGerald, Julian Michael, Luheng He, and
684 Luke Zettlemoyer. 2018. Large-scale qa-srl parsing.
685 *arXiv preprint arXiv:1805.05377*. 686

687	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. <i>arXiv preprint arXiv:2312.10997</i> , 2(1).	741
688		742
689		743
690		
691		
692	Arthur C Graesser and Natalie K Person. 1994. Question asking during tutoring. <i>American educational research journal</i> , 31(1):104–137.	744
693		745
694		746
695	Ryan Henry and Jiaqi Gong. 2025. Clare: Context-aware, interactive knowledge graph construction from transcripts. <i>Information</i> , 16(10):866.	747
696		
697		
698	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 others. 2025. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. <i>ACM Transactions on Information Systems</i> , 43(2):1–55.	748
699		749
700		750
701		
702		
703		
704		
705	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM computing surveys</i> , 55(12):1–38.	751
706		752
707		753
708		754
709		
710	Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension. <i>arXiv preprint arXiv:1706.04115</i> .	755
711		756
712		757
713		758
714	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in neural information processing systems</i> , 33:9459–9474.	759
715		760
716		761
717		762
718		763
719		
720		
721	Xiaoya Li, Fan Yin, Zijun Sun, Xiayu Li, Arianna Yuan, Duo Chai, Mingxin Zhou, and Jiwei Li. 2019. Entity-relation extraction as multi-turn question answering. <i>arXiv preprint arXiv:1905.05529</i> .	764
722		765
723		766
724		
725	Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. 2014. The stanford corenlp natural language processing toolkit. In <i>Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations</i> , pages 55–60.	767
726		768
727		769
728		770
729		771
730		
731	Dipak Meher, Carlotta Domeniconi, and Guadalupe Correa-Cabrera. 2025. Link-kg: Llm-driven coreference-resolved knowledge graphs for human smuggling networks. <i>arXiv preprint arXiv:2510.26486</i> .	772
732		773
733		774
734		775
735		776
736	Belinda Mo, Kysen Yu, Joshua Kazdan, Joan Cabezas, Proud Mpala, Lisa Yu, Chris Cundy, Charilaos Kanatsoulis, and Sanmi Koyejo. 2025. Kggen: Extracting knowledge graphs from plain text with language models. <i>arXiv preprint arXiv:2502.09956</i> .	777
737		778
738		779
739		780
740		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794

795 Bowen Zhang and Harold Soh. 2024. Extract, define,
796 canonicalize: An LLM-based framework for knowl-
797 edge graph construction. In *Proceedings of the 2024*
798 *Conference on Empirical Methods in Natural Lan-*
799 *guage Processing*, pages 9820–9836. Association for
800 Computational Linguistics.

801 Zexuan Zhong and Danqi Chen. 2021. A frustrat-
802 ingly easy approach for entity and relation extrac-
803 tion. In *Proceedings of the 2021 conference of the*
804 *North American chapter of the association for com-*
805 *putational linguistics: human language technologies*,
806 pages 50–61.

807 Yuqi Zhu, Xiaohan Wang, Jing Chen, Shuofei Qiao,
808 Yixin Ou, Yunzhi Yao, Shumin Deng, Huajun Chen,
809 and Ningyu Zhang. 2024. Llms for knowledge graph
810 construction and reasoning: Recent capabilities and
811 future opportunities. *World Wide Web*, 27(5):58.

A Ablation Studies

Prompt Archetype	Qwen-2.5		GPT-4o-mini		GPT-4o		Gemini-2.5		Claude-4	
	w/o 5W1H	Full	w/o 5W1H	Full	w/o 5W1H	Full	w/o 5W1H	Full	w/o 5W1H	Full
Role-Oriented (RO)	67.1	73.4	80.5	83.9	83.5	89.3	85.6	87.7	94.6	96.3
Procedural-Step (PS)	60.7	64.5	79.6	83.0	82.1	87.1	86.7	87.3	93.5	95.8
Instructional-Direct (ID)	65.3	68.6	77.9	83.5	79.7	81.2	83.9	88.7	91.3	96.3
Average	64.4	68.8	79.3	83.5	81.8	85.9	85.4	87.9	93.1	96.1

Table 4: Effectiveness of 5W1H-guided expansion across prompt archetypes. Scores represent factual retention (%) on the MINE benchmark. The results demonstrate that the 5W1H framework is a robust cognitive guide independent of stylistic framing.

Method	Qwen-2.5	GPT-4o-mini	GPT-4o	Gemini-2.5	Claude-4
KGGen	56.7	44.3	66.4	62.5	69.1
SoKG-EF	72.1	79.7	76.9	87.0	92.5
SoKG (Ours)	73.4	83.9	89.3	87.7	96.3

Table 5: Ablation study on input representation and triple extraction strategy. While SoKG-EF demonstrates the foundational impact of the QA scaffold, SoKG achieves peak performance by removing the entity-first bottleneck to maximize factual retention across all LLMs.

A.1 Robustness of Prompt Designs

To evaluate the contribution of the 5W1H framework, Table 4 compares the full 5W1H-integrated pipeline (Full) against the version omitting 5W1H-guided expansion (w/o 5W1H) across three distinct prompt archetypes:

- **Role-Oriented (RO):** Assigns a specific persona (e.g., Knowledge Archivist) and uses 5W1H as analytical lenses to guide deep exploration. This prompt design was adopted as the primary setting for our main experiments.
- **Procedural-Step (PS):** Defines a systematic workflow (Read → Segment → Generate) to ensure atomic factual extraction.
- **Instructional-Direct (ID):** Employs standard task-based instructions without complex role-play or multi-step procedures.

As shown in Table 4, the 5W1H framework provides a universal performance lift across all LLMs regardless of the underlying prompt structure. While the RO archetype generally yields the highest retention, peaking at 96.3% with Claude-4, even the more concise PS and ID templates show significant improvements once the interrogative scaffold is present. These results indicate that the 5W1H constraint functions as a fundamental cognitive guide that systematically surfaces procedural and causal dimensions.

A.2 The Entity-First Constraint

We evaluate the respective impacts of the QA scaffold and extraction strategy by comparing three configurations: KGGen, SoKG-EF (Entity-First), and SoKG. SoKG-EF incorporates both the extraction and consolidation logic of KGGen, applying this entity-centric pipeline to our QA-mediated scaffold. This setup allows us to evaluate the benefit of the scaffold independently while preserving the underlying entity-first logic.

As shown in Table 5, the superior performance of SoKG-EF over KGGen confirms that a QA scaffold effectively mitigates the complexity of raw text. However, SoKG’s even greater success reveals that rigid entity-first filtering acts as a restrictive bottleneck, limiting the model’s ability to capture full relational depth even when supported by a comprehensive QA scaffold.

Although intermediate stages introduce overhead, this rich interrogative structure justifies the investment by providing a dense semantic foundation for superior factual recoverability. Unlike isolated entity extraction which often incurs information loss through restrictive filtering, the QA-mediated scaffold preserves a richer semantic context. These results demonstrate that for structural refinement, a QA-driven approach offers a more systematic and inclusive foundation for knowledge construction than traditional entity-centric methods.

B QA Generation Prompt Details

869

B.1 Role-Oriented (RO), w/ 5W1H

870

```
## ROLE
You are a Comprehensive Knowledge Archivist who converts the [Full Document] into detailed,
document-grounded QA pairs.

## OBJECTIVE
Extract as many meaningful Question-Answer pairs as possible from the document.
Use the 5W1H perspectives (Who, What, When, Where, Why, How) as analytical lenses to help you
identify and expand potential questions, but do NOT restrict yourself to producing only
5W1H-type questions.
Your goal is to maximize informational coverage, capturing every explicit fact, relation, event,
definition, rationale, and process described in the document.

## INPUT
Full Document: "{document_text}"

## CONSTRAINTS
1. Context-Independent
  - Each QA must be self-contained and understandable without referencing the original text.
  - Replace pronouns with explicit entities.

2. No Hallucination
  - Use only facts explicitly stated in the document.

3. Expansion-Oriented Thinking
  - For each sentence or factual unit, consider the 5W1H perspectives as prompts to explore:
    - WHO is involved?
    - WHAT happened or is described?
    - WHEN did it occur?
    - WHERE did it occur?
    - WHY did it occur?
    - HOW was it carried out?
  - These perspectives are guides to inspire multiple possible QA pairs, even if they are
    implicit or only partially expressed.

4. Coverage
  - Extract all possible QA pairs that can be reasonably derived from the document.

## OUTPUT FORMAT
Return a JSON list of QA objects:

[
  {"question": "...", "answer": "..."},
  ...
]
```

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917 **B.2 Role-Oriented (RO), w/o 5W1H**

```
918 ## ROLE
919 You are a Comprehensive Knowledge Archivist who converts the [Full Document] into precise and
920 meaningful QA pairs.
921
922 ## OBJECTIVE
923 Extract as many high-quality Question-Answer pairs as needed to fully represent the document's
924 explicit information.
925 Use the following analytical perspectives as guides to discover potential questions, but do NOT
926 restrict your output to only these categories:
927
928 1. Entities & Definitions - Identify and clarify key terms, objects, roles, or concepts.
929 2. Properties & Characteristics - Extract notable features, attributes, components, or
930 qualities.
931 3. Events & Stated Facts - Capture actions, processes, or explicit factual statements.
932 4. Relationships & Dependencies - Identify connections, comparisons, or dependencies between
933 entities or ideas.
934
935 These perspectives are guides for expanding coverage, not mandatory categories.
936
937 ## INPUT
938 Full Document: "{document_text}"
939
940 ## CONSTRAINTS
941 1. Context-Independent
942 - Each QA must be self-contained and understandable without referencing the original text.
943 - Replace pronouns with explicit entities when needed.
944
945 2. No Hallucination
946 - Use only facts explicitly stated in the document.
947
948 3. Coverage without Inflation
949 - Extract all meaningful QA pairs that can be reasonably derived from the document.
950
951 ## OUTPUT FORMAT
952 Return a JSON list:
953
954 [
955   {"question": "...", "answer": "..."},
956   ...
957 ]
```

B.3 Procedural-Step (PS), w/ 5W1H

```
## ROLE
You are a **Document-Grounded QA Extractor**.

## OBJECTIVE
Convert the full document into high-coverage, explicit-fact QA pairs.

## PROCEDURE
1. Read the document end-to-end.
2. Segment into atomic factual units.
3. For each unit:
  - Generate QAs that capture all explicit information it contains.
  - When forming questions, view the unit through the 5W1H angles (Who, What, When, Where, Why, How) so that different aspects of the same fact can be covered.
4. Merge duplicates and keep the most precise wording.

## INPUT
Full Document: "{document_text}"

## CONSTRAINTS
- Context-Independent QAs only.
- No Hallucination.
- Prefer concise but complete answers.

## OUTPUT FORMAT
Return a JSON list:
[
  {"question": "...", "answer": "..."},
  ...
]
```

960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990

B.4 Procedural-Step (PS), w/o 5W1H

```
## ROLE
You are a **Document-Grounded QA Extractor**.

## OBJECTIVE
Convert the full document into high-coverage, explicit-fact QA pairs.

## PROCEDURE
1. Read the document end-to-end.
2. Segment into atomic factual units.
3. For each unit, generate QAs that capture all explicit information it contains.
4. Merge duplicates and keep the most precise wording.

## INPUT
Full Document: "{document_text}"

## CONSTRAINTS
- Context-Independent QAs only.
- No Hallucination.
- Prefer concise but complete answers.

## OUTPUT FORMAT
Return a JSON list:
[
  {"question": "...", "answer": "..."},
  ...
]
```

992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1020

1021 **B.5 Instructional-Direct (ID), w/ 5W1H**

1022 Read the following document and generate question-answer pairs based on its content.
1023 Generate as many high-quality questions as needed to cover the information explicitly stated in
1024 the document.
1025 For the same piece of information, consider the 5W1H dimensions (Who, What, When, Where, Why, How)
1026 and generate separate questions whenever different aspects are supported by the text.
1027 Do not stop at a single question if multiple 5W1H aspects can be identified.
1028 If different parts of the document support different questions, include all of them.
1029 Each question should be answerable using information explicitly stated in the document and written
1030 in a clear and self-contained manner.
1031
1032
1033 Input Document:
1034 "{document_text}"
1035
1036 Output Format:
1037 Return a JSON list of objects in the following form:
1038 [
1039 {"question": "...", "answer": "..."},
1040 ...
1041]

1043 **B.6 Instructional-Direct (ID), w/o 5W1H**

1044 Read the following document and generate question-answer pairs based on its content.
1045 Generate as many high-quality questions as needed to cover the information explicitly stated in
1046 the document.
1047 If different parts of the document support different questions, include all of them.
1048 Each question should be answerable using information explicitly stated in the document and written
1049 in a clear and self-contained manner.
1050
1051
1052 Input Document:
1053 "{document_text}"
1054
1055 Output Format:
1056 Return a JSON list of objects in the following form:
1057 [
1058 {"question": "...", "answer": "..."},
1059 ...
1060]

C Triple Extraction Prompt Details

1062

C.1 Triple Extraction from QA Pairs

1063

```
## ROLE
You are a Semantic Knowledge Graph Builder.
Extract every structured triples (entity1, relation, entity2) from the Q&A pair, following the
rules below.

## GOAL
From the question-answer pair, extract only useful, knowledge-ready triples that can serve as
entries in a semantic knowledge graph.

## RULES
Extract clean (subject, relation, object) triples following the rules:

1. Split every stated or clearly implied fact into minimal triples; integrate question and answer
context when needed.

2. Entities (entity1, entity2) must be short, concrete noun phrases.
- No pronouns (this, that, it, its, these, those, etc.).
- Entities must not be unresolved or reference-based pronouns (\eg those, they, someone,
anyone, whoever); if such a pronoun appears, rewrite it into a specific, explicit noun phrase
or skip the triple.
- No clauses or relative clauses (no "who/that/which/what/as it ..." inside an entity).
- No long gerund or sentence-like phrases. If a phrase contains a verb or clause marker,
rewrite it into a concise noun concept or skip the triple.

3. Relations must be short, canonical verbs or verb phrases.
- Express a single semantic link between the two entities (\eg causes, leads to, supports,
believes, opposes).
- Must be a compact predicate, not a sentence fragment.
- No pronouns or clause markers inside the relation (no "its", "that", "as it", "what", etc.).
- If the source uses an idiomatic or long expression, rewrite it into a simple canonical
relation without pronouns or embedded clauses, or skip the triple.

4. Include a fact if it can be clearly rewritten into a concise, explicit triple that fits the
rules above; otherwise skip it.

5. Output only concise, interpretable, knowledge-ready triples.

## INPUT
Q: {question}
A: {answer}

## OUTPUT FORMAT (JSON List)
- Return a list of JSON objects.
- Return [] if no valid triples exist.

[
  {"entity1": "Specific_Noun", "relation": "precise_verb_phrase", "entity2": "Specific_Noun"}
]
```

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

C.2 Triple Extraction from Raw Text (Direct Extraction)

```
1114
1115 ## ROLE
1116 You are a Semantic Knowledge Graph Builder.
1117
1118 Extract every structured triples (entity1, relation, entity2) from the text, following the rules
1119 below.
1120
1121 ## GOAL
1122 From the given text, extract only useful, knowledge-ready triples that can serve as entries in a
1123 semantic knowledge graph.
1124
1125 ## RULES
1126 Extract clean (subject, relation, object) triples following the rules:
1127
1128 1. Split every stated or clearly implied fact into minimal triples.
1129
1130 2. Entities (entity1, entity2) must be short, concrete noun phrases.
1131 - No pronouns (this, that, it, its, these, those, etc.).
1132 - Entities must not be unresolved or reference-based pronouns (\eg those, they, someone,
1133 anyone, whoever); if such a pronoun appears, rewrite it into a specific, explicit noun phrase
1134 or skip the triple.
1135 - No clauses or relative clauses (no "who/that/which/what/as it ..." inside an entity).
1136 - No long gerund or sentence-like phrases. If a phrase contains a verb or clause marker,
1137 rewrite it into a concise noun concept or skip the triple.
1138
1139 3. Relations must be short, canonical verbs or verb phrases.
1140 - Express a single semantic link between the two entities (\eg causes, leads to, supports,
1141 believes, opposes).
1142 - Must be a compact predicate, not a sentence fragment.
1143 - No pronouns or clause markers inside the relation (no "its", "that", "as it", "what", etc.).
1144 - If the source uses an idiomatic or long expression, rewrite it into a simple canonical
1145 relation without pronouns or embedded clauses, or skip the triple.
1146
1147 4. Include a fact if it can be clearly rewritten into a concise, explicit triple that fits the
1148 rules above; otherwise skip it.
1149
1150 5. Output only concise, interpretable, knowledge-ready triples.
1151
1152 ## INPUT
1153 Text: {document_text}
1154
1155 ## OUTPUT FORMAT (JSON List)
1156 - Return a list of JSON objects.
1157 - Return [] if no valid triples exist.
1158
1159 [
1160   [{"entity1": "Specific_Noun", "relation": "precise_verb_phrase", "entity2": "Specific_Noun"}]
1161 ]
```